

# *X/Open Preliminary Specification*

## **DCE 1.1: Distributed File Service Specification**

*X/Open Company Ltd.*



© September 1996, X/Open Company Limited

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior permission of the copyright owners.

X/Open Preliminary Specification  
DCE 1.1: Distributed File Service Specification  
X/Open Document Number: P409

Published by X/Open Company Ltd., U.K.

Any comments relating to the material contained in this document may be submitted to:

The Open Group  
Apex Plaza  
Forbury Road  
Reading  
Berkshire, RG1 1AX  
United Kingdom

or by Electronic Mail to:

XoSpecs@tog.org

# Contents

<b>Part</b>	<b>1</b>	<b>Distributed File System Introduction .....</b>	<b>1</b>
<b>Chapter</b>	<b>1</b>	<b>Introduction to the DFS Specification .....</b>	<b>3</b>
	1.1	Portability.....	3
	1.1.1	Overview.....	3
	1.1.2	IDLs, Data Types, Constants and Flags .....	3
	1.2	Document Organization .....	4
	1.2.1	RPC Interfaces for <b>DCE DFS</b> .....	4
	1.2.2	Access Control Lists ( <b>ACLs</b> ) in DFS.....	4
	1.2.3	The DCE DFS VFS+ Interface Specification.....	5
	1.3	Terminology .....	6
	1.4	Conformance Requirements .....	9
<b>Part</b>	<b>2</b>	<b>DFS Interface Definition Language and Data.....</b>	<b>11</b>
<b>Chapter</b>	<b>2</b>	<b>File Exporter IDL Declarations .....</b>	<b>13</b>
	2.1	The File Exporter Interface .....	13
	2.2	Mask Values for afsStoreStatus Structure .....	14
	2.3	AFS_GetTime Sync Constants.....	14
	2.4	Flag Parameters for ACL Type Parameter .....	14
	2.5	Definitions for Flag Used in Getting a Token.....	14
	2.6	Definitions for Flag in Token Recovery (TSR).....	15
	2.7	Definitions for Flag in Fileset Operations .....	15
	2.8	Definitions for Flags in AFS_SetParams() .....	16
	2.9	Definitions for Flags in AFS_SetContext() .....	16
	2.10	Definitions for Client-only Attribute Flags .....	16
	2.11	Data Types .....	16
	2.11.1	Define Generic Network Address Information .....	16
	2.11.2	Define afsVolSync Structure for Tracking Replicas.....	16
	2.11.3	Define the afsFetchStatus Structure .....	17
	2.11.4	Define the afsStoreStatus Structure.....	17
	2.11.5	Structure for Physical File System Type.....	18
	2.11.6	Structure for Statistics .....	18
	2.12	Various Bulk Data Type Definitions for RPC Operations.....	20
	2.12.1	The afsFidExp Structure .....	20
	2.12.2	The afsBulkFEX Structure.....	20
	2.12.3	The afsACL Structure.....	20
	2.12.4	The afsQuota Structure .....	20
	2.12.5	The afsBulkVVs Structure.....	20
	2.12.6	The afsBulkVolIDs Structure.....	21
	2.12.7	The afsBulkStats Structure .....	21

2.13	Definitions of File Server Exported Operations.....	22
2.13.1	AFS_SetContext.....	22
2.13.2	AFS_LookupRoot.....	23
2.13.3	AFS_FetchData.....	24
2.13.4	AFS_FetchACL.....	25
2.13.5	AFS_FetchStatus.....	26
2.13.6	AFS_StoreData.....	27
2.13.7	AFS_StoreACL.....	28
2.13.8	AFS_StoreStatus.....	29
2.13.9	AFS_RemoveFile.....	30
2.13.10	AFS_CreateFile.....	31
2.13.11	AFS_Rename.....	32
2.13.12	AFS_Symlink.....	33
2.13.13	AFS_HardLink.....	34
2.13.14	AFS_MakeDir.....	35
2.13.15	AFS_RemoveDir.....	36
2.13.16	AFS_Readdir.....	37
2.13.17	AFS_Lookup.....	38
2.13.18	AFS_GetToken.....	39
2.13.19	AFS_ReleaseTokens.....	40
2.13.20	AFS_GetTime.....	41
2.13.21	AFS_MakeMountPoint.....	42
2.13.22	AFS_GetStatistics.....	43
2.13.23	AFS_BulkFetchVV.....	44
2.13.24	AFS_BulkKeepAlive.....	45
2.13.25	AFS_ProcessQuota.....	46
2.13.26	AFS_GetServerInterfaces.....	47
2.13.27	AFS_SetParams.....	48
2.14	The File Exporter Interface End.....	49
<b>Chapter 3</b>	<b>Cache Manager IDL Declarations.....</b>	<b>51</b>
3.1	The Token Manager Interface.....	51
3.2	Definitions for Flags Field.....	52
3.3	Definitions of Token Manager Exported Operations.....	52
3.3.1	TKN_Probe.....	52
3.3.2	TKN_InitTokenState.....	52
3.3.3	TKN_TokenRevoke.....	53
3.3.4	TKN_GetCellName.....	53
3.3.5	TKN_GetLock.....	53
3.3.6	TKN_GetCE.....	54
3.3.7	TKN_GetServerInterfaces.....	54
3.3.8	TKN_SetParams.....	54
3.3.9	TKN_AsyncGrant.....	55
3.4	The Token Manager Interface End.....	55

<b>Chapter</b>	<b>4</b>	<b>Common IDL Data.....</b>	<b>57</b>
	4.1	Interface Common Data.....	57
	4.2	General AFS Constants .....	57
	4.3	Constants for Cell and Hosts .....	57
	4.4	AFS Object Types Used by AFS_Mount .....	58
	4.5	Quota Types for Quota Setting Commands .....	58
	4.6	Quota Opcodes for Quota Setting Commands.....	58
	4.7	Physical File System Types (for afsFStype).....	58
	4.8	Volume Types Used for afsVolumeType.....	58
	4.9	Values for the afsRevokeDesc Flags.....	59
	4.10	Values Used in afsRevokeDesc's outFlags .....	59
	4.11	Data Types .....	59
	4.11.1	General Definitions for AFS Data Structures .....	59
	4.11.2	General Tagged-name for a Pathname Component.....	59
	4.11.3	Define the afsTaggedName Structure .....	59
	4.11.4	Define the afsTaggedPath Structure .....	60
	4.11.5	Define the afsNetAddr Structure .....	60
	4.11.6	Define the afsTimeval Structure.....	60
	4.11.7	Define the afsHyper Structure.....	60
	4.11.8	Define the afsFid Structure.....	60
	4.11.9	Define the afsFidTaggedName Structure .....	60
	4.11.10	Define the afsToken Structure.....	61
	4.11.11	Define the afsRecordLock Structure.....	61
	4.11.12	Define the afsRevokeDesc Structure .....	61
	4.11.13	Define the afsReturnDesc Structure .....	62
	4.11.14	Define the afsConnParams Structure .....	62
	4.11.15	Define the afsDBLockDesc Structure .....	62
	4.11.16	Define the afsDBCACHEEntry Structure .....	62
	4.11.17	Define the afsDBLock Structure .....	63
	4.12	Various Bulk Typedefs From Primitive Structures.....	63
	4.12.1	Define the afsRevokes Structure .....	63
	4.12.2	Define the afsReturns Structure .....	63
	4.12.3	Define the afsFids Structure.....	63
	4.12.4	Define the afsTokens Structure.....	63
	4.12.5	Define the afsStrings Structure.....	64
	4.13	Data Types for DFS RPC Versioning Scheme .....	64
	4.13.1	Constants for RPC Versioning Scheme.....	64
	4.13.2	dfs_interfaceDescription Structure.....	64
	4.13.3	dfs_interfaceList Structure .....	64
	4.14	Interface Common Data End .....	65

<b>Part</b>	<b>3</b>	<b>DFS Versioning Scheme .....</b>	<b>67</b>
<b>Chapter</b>	<b>5</b>	<b>DFS RPC Versioning Scheme .....</b>	<b>69</b>
	5.1	Data Structures .....	69
	5.2	DFS Versioning API .....	70
	5.2.1	Register an Interface with Versioning Mechanism .....	71
	5.2.2	Printing a Returned List of Interfaces .....	71
	5.2.3	Function Called at Startup to Identify the Server Interface .....	71
	5.2.4	Function to Compare Interfaces.....	72
	5.3	Example IDL File .....	73
	5.3.1	The UPDATE Server Interface .....	73
	5.3.2	Constants for the UPDATE Interface .....	73
	5.3.3	Define the updateFileStatS Structure .....	74
	5.3.4	UPDATE_GetServerInterfaces.....	74
	5.3.5	UPDATE_FetchInfo.....	74
	5.3.6	UPDATE_FetchFile.....	75
	5.3.7	UPDATE_FetchObjectInfo .....	75
	5.3.8	The UPDATE Server Interface End.....	75
	5.4	Example Client Application.....	76
	5.4.1	Typical Client Headers.....	76
	5.4.2	Headers for Serviceability .....	76
	5.4.3	Constants for the Client Application .....	76
	5.4.4	RPC prewrap and postwrap functions .....	77
	5.4.5	Typical Serviceability initialization .....	77
	5.4.6	ANSI C Declaration .....	78
	5.4.7	Client's import list.....	78
	5.4.8	Client Globals.....	78
	5.4.9	Using dfs_selectInterface to Select One.....	78
	5.4.10	Example RPC Wrappers .....	80
	5.5	Example Server Application .....	83
	5.5.1	Typical Server Headers .....	83
	5.5.2	Headers for Serviceability .....	83
	5.5.3	Constants for the Server Application.....	84
	5.5.4	Define IS_COMM_ERR Function.....	84
	5.5.5	Typical Serviceability Initialization.....	84
	5.5.6	Skeleton Server Main.....	84
	5.6	Example Manager Application.....	87
	5.6.1	Typical MANAGER Headers.....	87
	5.6.2	Headers for Serviceability .....	87
	5.6.3	Typical Initialization .....	88
	5.6.4	ANSI C Declaration .....	88
	5.6.5	Example Mutex Lock.....	88
	5.6.6	Example Mutex Unlock.....	88
	5.6.7	Example Error Exit.....	88
	5.6.8	UPDATE_GetServerInterfaces.....	89
	5.6.9	UPDATE_FetchFile.....	90
	5.6.10	UPDATE_FetchInfo.....	91
	5.6.11	UPDATE_FetchObjectInfo .....	93

5.6.12	UPDATE_v4_0_manager_epv .....	94
5.7	Extending the DFS Interface .....	95
5.7.1	Case 1 — An Owner Extends a DFS Interface.....	95
5.7.2	Case 2 — A Non-owner Extends a DFS Interface .....	96
<b>Part</b>	<b>4 DFS Client/Server and Token Manager Interfaces....</b>	<b>97</b>
<b>Chapter</b>	<b>6 DCE DFS File Exporter Interface .....</b>	<b>99</b>
	<i>AFS_SetContext()</i> .....	103
	<i>AFS_LookupRoot()</i> .....	104
	<i>AFS_FetchData()</i> .....	105
	<i>AFS_FetchACL()</i> .....	107
	<i>AFS_FetchStatus()</i> .....	108
	<i>AFS_StoreData()</i> .....	109
	<i>AFS_StoreACL()</i> .....	111
	<i>AFS_StoreStatus()</i> .....	113
	<i>AFS_RemoveFile()</i> .....	114
	<i>AFS_Lookup()</i> .....	116
	<i>AFS_CreateFile()</i> .....	118
	<i>AFS_Rename()</i> .....	120
	<i>AFS_Symlink()</i> .....	122
	<i>AFS_MakeMountPoint()</i> .....	124
	<i>AFS_HardLink()</i> .....	126
	<i>AFS_MakeDir()</i> .....	128
	<i>AFS_RemoveDir()</i> .....	130
	<i>AFS_Readdir()</i> .....	132
	<i>AFS_GetToken()</i> .....	134
	<i>AFS_GetStatistics()</i> .....	136
	<i>AFS_ReleaseTokens()</i> .....	137
	<i>AFS_GetTime()</i> .....	138
	<i>AFS_BulkFetchVV()</i> .....	139
	<i>AFS_ProcessQuota()</i> .....	140
	<i>AFS_BulkKeepAlive()</i> .....	141
	<i>AFS_SetParams()</i> .....	142
	<i>AFS_GetServerInterfaces()</i> .....	143
<b>Chapter</b>	<b>7 Cache Manager Service Interface .....</b>	<b>145</b>
	<i>TKN_InitTokenState()</i> .....	146
	<i>TKN_Probe()</i> .....	147
	<i>TKN_TokenRevoke()</i> .....	148
	<i>TKN_GetCellName()</i> .....	149
	<i>TKN_SetParams()</i> .....	150
	<i>TKN_GetServerInterfaces()</i> .....	151
	<i>TKN_AsyncGrant()</i> .....	152

<b>Part</b>	<b>5</b>	<b>Access Control Lists (ACLs) in DFS .....</b>	<b>153</b>
<b>Chapter</b>	<b>8</b>	<b>Access Control List Overview .....</b>	<b>155</b>
	8.1	ACL Types .....	155
	8.2	ACL Entry Types .....	156
	8.2.1	Rules Governing Entries for Filesystem ACLs.....	156
	8.2.2	Optional ACL Entry Mask Types for Authenticated Users.....	157
	8.2.3	ACL Entry Types for Unauthenticated Users.....	157
	8.2.4	Simple and Complex Entry Types .....	158
	8.3	Delegates and DCE LFS Objects .....	158
	8.4	ACL Entry Types for Delegation.....	159
	8.4.1	ACL Access Types for Delegation .....	159
	8.5	Access Types.....	159
	8.6	Interaction of Filesystem ACLs with UNIX Permission Bits .....	160
	8.7	Access Check Algorithm.....	161
	8.8	Access Check Algorithm for Delegation .....	162
	8.8.1	Delegation and non-DCE LFS Objects .....	163
<b>Chapter</b>	<b>9</b>	<b>ACL Storage Format.....</b>	<b>165</b>
	9.1	Principal ID Formats.....	165
	9.2	Principal Identifier Format .....	165
	9.3	Foreign Cell Principal Identifier Format.....	166
	9.4	Access Type Format .....	166
	9.4.1	Access Type Definitions.....	166
	9.5	ACL Entry Type Format.....	167
	9.6	Simple ACL Entry Format .....	168
	9.7	Complex ACL Entry Format .....	168
	9.7.1	Extended Complex Entry Format .....	169
	9.8	ACL Structure .....	169
	9.9	The Structures for Reading Lists of sec_id_t .....	170
	9.9.1	Define the epi_sec_id Structure.....	170
	9.9.2	Define the epi_sec_id_foreign Structure .....	170
<b>Chapter</b>	<b>10</b>	<b>ACL Interface Functions.....</b>	<b>171</b>
		<i>dacl_CheckAccessId()</i> .....	172
		<i>dacl_DetermineAccessAllowed()</i> .....	174
		<i>dacl_CheckAccessPac()</i> .....	176
		<i>dacl_epi_CheckAccessPac()</i> .....	177
		<i>dacl_CheckAccessAllowedPac()</i> .....	178
		<i>dacl_epi_CheckAccessAllowedPac()</i> .....	179
		<i>dacl_PacFromUcred()</i> .....	180
		<i>dacl_FlattenAclWithModeBits()</i> .....	181
		<i>dacl_FlattenAcl()</i> .....	182
		<i>dacl_epi_FlattenAcl()</i> .....	183
		<i>dacl_ParseAclDiskOption()</i> .....	184
		<i>dacl_ParseSyscallAcl()</i> .....	185
		<i>dacl_ParseAcl()</i> .....	186
		<i>dacl_ExtractPermBits()</i> .....	187



<i>dacl_ChmodAcl()</i> .....	188
<i>dacl_FreeAclEntries()</i> .....	189
<i>dacl_PrintAclEntry()</i> .....	190
<i>dacl_PrintAcl()</i> .....	191
<i>dacl_WriteToDisk()</i> .....	192
<i>dacl_CreateAclOnDisk()</i> .....	193
<i>dacl_ReadFromDisk()</i> .....	194
<i>dacl_AddEntryToAcl()</i> .....	195
<i>dacl_ModifyAclEntry()</i> .....	196
<i>dacl_DeleteAclEntry()</i> .....	197
<i>dacl_DeleteAllEntries()</i> .....	198
<i>sec_acl_FlattenAcl()</i> .....	199
<i>sec_acl_ParseAcl()</i> .....	200
<i>dacl_InitAclEntryFromStrings()</i> .....	201
<i>dacl_NameAndTypeStringsFromEntry()</i> .....	202
<i>dacl_EntryType_ToString()</i> .....	203
<i>dacl_EntryType_FromString()</i> .....	204
<i>dacl_Permsset_ToString()</i> .....	205
<i>dacl_Permsset_FromString()</i> .....	206
<i>dacl_ValidateBuffer()</i> .....	207
<i>dacl_InitEpiAcl()</i> .....	208
<i>dacl_AreObjectEntriesRequired()</i> .....	209
<i>dacl_AreObjectUuidsRequiredOnAccessCheck()</i> .....	210
<i>dacl_ArePermBitsRequiredOnAccessCheck()</i> .....	211
<i>dacl_AclMgrName()</i> .....	212
<i>Epi_PrinId_ToUuid()</i> .....	213
<i>Epi_PrinId_FromUuid()</i> .....	214
<i>Epi_PrinId_Cmp()</i> .....	215
<i>hton_epi_uuid()</i> .....	216
<i>ntoh_epi_uuid()</i> .....	217
<i>hton_epi_principal_id()</i> .....	218
<i>ntoh_epi_principal_id()</i> .....	219

<b>Appendix A</b>	<b>Mapping DFS ACLs to UNIX Mode Bits .....</b>	<b>221</b>
A.1	Relationship Between ACLs and UNIX Mode Bits .....	221
A.1.1	ACLs and UNIX Mode Bits .....	221
A.1.2	Changing other Mode Bits.....	222
A.1.3	Changing group Mode Bits on a File with ACL with No mask_obj.....	222
A.1.4	Changing the user_obj Entry .....	222
A.1.5	Changing the group_obj Entry .....	223
A.1.6	Changing the group_obj Entry in an ACL with a mask_obj Entry..	223
A.1.7	The Results of chmod of the group Bits .....	224
A.1.8	The Results of acl_edit of the mask_obj Bits .....	224
A.1.9	Changing the group_obj Entry for an ACL with a mask_obj .....	225
A.2	File Creation Mask and ACLs.....	226
A.2.1	File Creation with no Initial Object ACL.....	226
A.2.2	File Creation with an Initial Object ACL.....	226
A.2.3	File ACL Creation with an initial object ACL .....	227

<b>Appendix B</b>	<b>Access Control List Package Error List.....</b>	<b>229</b>
B.1	Access Control List Return Values .....	229
B.2	Filesystem Access Control List Return Values.....	230
<b>Part 6</b>	<b>The DCE DFS VFS+ Interface Specification .....</b>	<b>231</b>
<b>Chapter 11</b>	<b>DCE DFS VFS+ Interface Introduction.....</b>	<b>233</b>
11.1	Definition of Terms .....	233
11.2	VFS+ Interface Goals and Constraints.....	234
11.3	Overview of Interfaces to the LFS .....	235
11.3.1	Locking.....	236
11.3.2	Credentials.....	237
11.3.3	Operations on Filesets and Aggregates.....	237
11.4	Organization of the VFS+ Switch.....	238
11.4.1	The O-ops.....	238
11.4.2	The N-ops.....	238
11.4.3	The X-ops .....	238
11.5	Basic Operation.....	239
11.6	The Fileset Registry.....	239
11.7	The Fileset/Vnode Interfaces.....	239
11.8	References from the LFS back into DFS.....	240
11.8.1	LFS Introducing Itself to DFS.....	240
11.8.2	Operating on DFS Lock Structures .....	240
11.8.3	Allocating and Freeing Memory .....	241
11.8.4	Fileset Creation Assistance.....	241
11.8.5	Obtaining a Fileset Structure .....	242
11.8.6	Releasing a Fileset .....	242
11.8.7	Conversion of Operations Vectors.....	242
11.8.8	Getting the Local Cell ID .....	243
11.8.9	Getting the Administrative Group ID.....	243
11.8.10	Obtaining the Identity of the Principal.....	243
<b>Appendix C</b>	<b>Components of a Typical VFS+ Package.....</b>	<b>245</b>
C.1	The VFS Vector.....	245
C.2	Naming Conventions .....	245
C.3	Glue Functions .....	246
C.3.1	xglue_rename() Function .....	246
C.3.2	xglue_lookup() Function.....	247
C.4	Extended Functions .....	248
C.4.1	xufs_lookup() Function .....	248
C.5	Extended Vnode Attributes.....	249
<b>Chapter 12</b>	<b>The DCE DFS ACL Model for an LFS.....</b>	<b>251</b>
12.1	Overview.....	251
12.2	Definition of Terms .....	251
12.3	Primitive Data Types .....	252
12.4	Local Realm .....	252
12.5	DFS Administrator.....	252

12.6	DFS vs non-DFS requests.....	252
12.7	ACL Contents.....	252
12.8	Define the External ACL Representation .....	254
12.8.1	Formats of ACL Types.....	255
12.9	Relationship — ACLs and UNIX Protections .....	256
12.9.1	Episode Visibility above VFS+ .....	257
12.10	ACL Creation From Mode Bits Algorithm.....	257
12.11	Initial ACL and File Creation Algorithm .....	258
12.12	User and Group Identities .....	259
12.13	Algorithm for Obtaining a Principal Identity.....	259
12.14	Algorithm for Generating PAC from Ucred Structure .....	261
12.15	perm_control Access Right .....	261
12.16	Access Rights Algorithm .....	262
12.16.1	No Object ACL Exists Rights Algorithm.....	263
12.16.2	Object ACL Exists Rights Algorithm .....	263
12.17	mask_obj ACL Entry Algorithm Impact .....	265
12.18	Miscellaneous Topics and Suggestions.....	266
<b>Chapter 13</b>	<b>VFS+ Data Types.....</b>	<b>267</b>
13.1	Primitive Data Types .....	267
13.2	Aggregates and Aggregate Registry Data Types.....	269
13.2.1	Aggregate Static Status .....	269
13.2.2	Aggregate Dynamic Status.....	269
13.2.3	Aggregate Status .....	270
13.2.4	Valid Aggregate Types .....	270
13.2.5	Adding of New Filesystem Types.....	270
13.2.6	Aggregate Structure.....	270
13.2.7	Aggregate Field Definitions as a_* Items .....	271
13.2.8	Aggregate States .....	271
13.2.9	Fileset Creation Flags .....	271
13.2.10	Aggregate Attach Flags.....	271
13.2.11	Aggregate Sync Flags .....	271
13.2.12	Aggregate Operations Vector .....	272
13.2.13	Exported Aggregate Registry Items .....	272
13.2.14	Vnode Ops Classification.....	273
13.2.15	Concurrency for vol_stat_st.....	273
13.2.16	Aggregate Table Entry Format .....	273
13.2.17	UFS Entry Extra Data.....	274
13.2.18	Values for Filesystem Type.....	274
13.2.19	Adding New Filesystem Types .....	274
13.3	Fileset Data Types .....	274
13.3.1	Define the Fileset Function Array.....	274
13.3.2	Volume Operations Definitions.....	276
13.3.3	The Fileset NextHole Structure .....	277
13.3.4	Define the vol_stat_st Structure .....	277
13.3.5	Define the vol_stat_dy Structure .....	279
13.3.6	Define the vol_status Structure .....	282
13.3.7	Define the vol_statusDesc Structure .....	282

13.3.8	Define Transient Error .....	282
13.3.9	Volume States (for vol_stat_st) .....	282
13.3.10	On-disk States .....	283
13.3.11	Kernel Maintained State Bits.....	284
13.3.12	Useful Macros in Volume Pointers .....	284
13.3.13	Mask Bits for VOL_SETSTATUS MASK Argument.....	284
13.3.14	Define the Fileset (volume) Structure .....	285
13.3.15	Define Volume Fields as v_*.....	286
13.3.16	Define the vol_Dirent Structure .....	286
13.3.17	Define the Fileset Error Codes .....	286
13.3.18	Define the Fileset Open Operation Types .....	286
13.3.19	Define the Fileset Operations VOL_SYS_XXX .....	287
13.3.20	Sync Type Values for vol_sync().....	288
13.3.21	Define the Fileset Handle Structure .....	288
13.3.22	Define the Root Anode Index .....	289
13.3.23	Define the Maximum Quota Size.....	289
13.4	Extended Credential Data Types .....	289
13.4.1	Define the xcred_PListEntry_t Structure.....	289
13.4.2	Define the xcred_t Structure .....	290
13.4.3	Reserved Attributes .....	290
13.5	The VFS+ Switch .....	290
13.5.1	DCE LFS VNOPS Vector Organization.....	290
13.5.2	LFS Generating Extended Vnode Operations .....	291
13.5.3	Converting Old Vnode Operations.....	292
13.6	Vnode Data Types .....	293
13.6.1	Preliminary Items .....	293
13.6.2	Converted Vnode Indication.....	293
13.6.3	Define the xvfs_attr Structure.....	294
13.6.4	Define the Vnode Operations Function Array .....	296
13.6.5	Define the Enhanced Vnode Operations .....	298
13.6.6	Define VOPX_XXX for Extended Vnodeops.....	298
13.6.7	Define VOPX_UPDATE Flags.....	299
13.6.8	Define the Enhanced Operations Vector .....	299
13.6.9	Define the VFS Operations .....	300
13.6.10	Define the Vnode Operation Classifications.....	300
13.6.11	Directory Entry Formats .....	301
<b>Chapter 14</b>	<b>Aggregate Operations Interface.....</b>	<b>303</b>
14.1	Initialization .....	303
14.1.1	Identifying a New LFS to DCE DFS .....	303
14.1.2	Registering Aggregate Operations .....	304
14.2	Exporting volumes to DFS .....	304
14.3	Aggregate Mounts.....	305
14.4	Aggregate Array Functions .....	305
	ag_hold() .....	306
	ag_rele() .....	307
	ag_lock() .....	308
	ag_unlock() .....	309

	<i>ag_stat()</i> .....	310
	<i>ag_volCreate()</i> .....	311
	<i>ag_volInfo()</i> .....	312
	<i>ag_detach()</i> .....	313
	<i>ag_attach()</i> .....	314
	<i>ag_sync()</i> .....	315
<b>Appendix D</b>	<b>Information Pertinent to <i>ag_volCreate()</i></b> .....	<b>317</b>
D.1	Static Status .....	317
D.2	Dynamic Status .....	317
D.3	Other Items .....	318
D.4	Attaching the New Fileset .....	318
<b>Appendix E</b>	<b>Information Pertinent to <i>ag_volInfo()</i></b> .....	<b>321</b>
E.1	Volume Structure Fields .....	321
E.2	Static Status .....	321
<b>Appendix F</b>	<b>Information Pertinent to <i>ag_[de,at]tach()</i></b> .....	<b>323</b>
F.1	Making an Aggregate Available .....	323
F.2	Private Storage .....	323
F.3	Aggregate Structure .....	323
<b>Chapter 15</b>	<b>Fileset (Volume) Operations Interface</b> .....	<b>325</b>
15.1	Overview .....	325
15.1.1	Classes of Fileset Operations .....	327
15.1.2	Fileset Registry Functions .....	327
15.2	Fileset Types Overview .....	327
15.3	Fileset Clone Algorithms .....	328
15.3.1	Requirements on Cloning .....	329
15.3.2	Uses for Clones .....	330
15.3.3	Some Fileset and Clone Requirements .....	331
15.4	Fileset Indices .....	331
15.5	LFS Modification of Fileset Status .....	332
15.6	Zero Link Count Files .....	332
15.7	Fileset Quotas .....	333
15.8	Anode Generation Numbers .....	334
15.9	File Identifiers (afsFids) .....	334
15.10	Looping Operation Considerations .....	335
15.11	Vnode to LFS Association .....	336
15.11.1	Determining whether a Fileset is Local .....	336
15.11.2	Handling Dis-sociations and Re-associations .....	337
15.11.3	Complications in Dis-sociations and Re-associations .....	338
15.11.4	Fileset Moves .....	338
15.12	Private LFS Fileset Data .....	339
15.13	Fileset Operations .....	339
15.14	Fileset Array Functions .....	340
	<i>vol_hold()</i> .....	341
	<i>vol_rele()</i> .....	342

	<i>vol_lock()</i> .....	343
	<i>vol_unlock()</i> .....	344
	<i>vol_open()</i> .....	345
	<i>vol_seek()</i> .....	347
	<i>vol_tell()</i> .....	348
	<i>vol_scan()</i> .....	349
	<i>vol_close()</i> .....	350
	<i>vol_destroy()</i> .....	351
	<i>vol_deplete()</i> .....	352
	<i>vol_attach()</i> .....	353
	<i>vol_detach()</i> .....	354
	<i>vol_getstatus()</i> .....	355
	<i>vol_setstatus()</i> .....	356
	<i>vol_create()</i> .....	357
	<i>vol_read()</i> .....	358
	<i>vol_write()</i> .....	359
	<i>vol_truncate()</i> .....	360
	<i>vol_delete()</i> .....	361
	<i>vol_getattr()</i> .....	362
	<i>vol_setattr()</i> .....	363
	<i>vol_getacl()</i> .....	365
	<i>vol_setacl()</i> .....	367
	<i>vol_clone()</i> .....	369
	<i>vol_reclone()</i> .....	371
	<i>vol_unclone()</i> .....	373
	<i>vol_vget()</i> .....	375
	<i>vol_root()</i> .....	376
	<i>vol_isroot()</i> .....	377
	<i>vol_getvv()</i> .....	378
	<i>vol_setdystat()</i> .....	379
	<i>vol_freedystat()</i> .....	380
	<i>vol_setnewvid()</i> .....	381
	<i>vol_copyacl()</i> .....	382
	<i>vol_concurr()</i> .....	384
	<i>vol_swapids()</i> .....	385
	<i>vol_sync()</i> .....	386
	<i>vol_pushstatus()</i> .....	387
	<i>vol_readdir()</i> .....	388
	<i>vol_appenddir()</i> .....	390
	<i>vol_bulksetstatus()</i> .....	392
	<i>vol_getzlc()</i> .....	394
	<i>vol_getnextholes()</i> .....	395
15.15	Fileset Registry Array Functions .....	396
	<i>volreg_Enter()</i> .....	397
	<i>volreg_Delete()</i> .....	398
	<i>volreg_Lookup()</i> .....	399

<b>Appendix G</b>	<b>Exporting the Filesets in an Aggregate</b> .....	<b>401</b>
G.1	Fileset Export Steps.....	401
<b>Appendix H</b>	<b>Filled Values for vol_dirent Fields</b> .....	<b>403</b>
H.1	Returned Values.....	403
<b>Appendix I</b>	<b>Values for vol_open</b> .....	<b>405</b>
I.1	The Type Argument.....	405
I.1.1	Concurrency.....	406
I.1.2	Handling Inconsistent State.....	406
I.2	The Fileset Handle .....	407
I.2.1	Preparing for Fileset Operations .....	407
<b>Appendix J</b>	<b>Status Returned for vol_getstatus</b> .....	<b>409</b>
J.1	Fileset Status Set .....	409
<b>Appendix K</b>	<b>Status Set for vol_setstatus</b> .....	<b>411</b>
K.1	Status Set .....	411
<b>Appendix L</b>	<b>Processing for vol_create</b> .....	<b>413</b>
L.1	Processing Accomplished.....	413
<b>Chapter 16</b>	<b>VFS (Vnode) Interface and Operations</b> .....	<b>415</b>
16.1	Overview.....	415
16.1.1	Enhanced Vnode Operations Vector .....	415
16.1.2	Converted Vnodes.....	418
16.1.3	Enhanced Vfs Vector.....	419
16.2	Administrative Rights .....	420
16.3	Copy-on-Write Impacts .....	421
16.4	Swap Files .....	421
16.5	Synchronization Between Vnode and Fileset Operations.....	422
16.5.1	Directory Offsets.....	423
16.6	Vfs Operations .....	423
16.6.1	<i>vfs_getvolume()</i> .....	424
16.7	Base Vnode Interface.....	424
	<i>vn_open()</i> .....	427
	<i>vn_close()</i> .....	428
	<i>vn_rdwr()</i> .....	429
	<i>vn_ioctl()</i> .....	430
	<i>vn_select()</i> .....	431
	<i>vn_getattr()</i> .....	432
	<i>vn_setattr()</i> .....	434
	<i>vn_access()</i> .....	436
	<i>vn_lookup()</i> .....	437
	<i>vn_create()</i> .....	438
	<i>vn_remove()</i> .....	440
	<i>vn_link()</i> .....	441
	<i>vn_rename()</i> .....	442

	<i>vn_mkdir()</i> .....	444
	<i>vn_rmdir()</i> .....	446
	<i>vn_readdir()</i> .....	447
	<i>vn_symlink()</i> .....	450
	<i>vn_readlink()</i> .....	451
	<i>vn_fsync()</i> .....	452
	<i>vn_inactive()</i> .....	453
	<i>vn_bmap()</i> .....	454
	<i>vn_strategy()</i> .....	455
	<i>vn_bread()</i> .....	456
	<i>vn_brelse()</i> .....	457
	<i>vn_lockctl()</i> .....	458
	<i>vn_fid()</i> .....	459
16.8	Extended Vnode Interface .....	460
16.8.1	Generalized Differences.....	460
	<i>vn_getvolume()</i> .....	461
	<i>vn_afsid()</i> .....	462
	<i>vn_getacl()</i> .....	463
	<i>vn_setacl()</i> .....	465
<b>Chapter 17</b>	<b>DCE DFS VFS+ Extended Credential Package.....</b>	<b>469</b>
17.1	The xcred Package.....	469
17.2	Package Overview.....	469
17.2.1	Interface Overview .....	469
17.3	xcred Functions.....	470
	<i>xcred_Init()</i> .....	471
	<i>xcred_Create()</i> .....	472
	<i>xcred_Hold()</i> .....	473
	<i>xcred_Delete()</i> .....	474
	<i>xcred_Release()</i> .....	475
	<i>xcred_AssociateCreds()</i> .....	476
	<i>xcred_UCredToXCred()</i> .....	477
	<i>xcred_FindByPag()</i> .....	478
	<i>xcred_PutProp()</i> .....	479
	<i>xcred_GetProp()</i> .....	480
	<i>xcred_GetUFlags()</i> .....	481
	<i>xcred_SetUFlags()</i> .....	482
	<i>xcred_EnumerateProp()</i> .....	483
	<i>xcred_DeleteEntry()</i> .....	484
	<b>Index.....</b>	<b>485</b>
 <b>List of Examples</b>		
5-1	Interface Description for an Interface.....	69
5-2	List of Interface Descriptions for an Interface .....	70
5-3	Usage of IMPORT definition with _TAKES notation.....	70
5-4	Construction of the <i>getenv</i> utility.....	70



A-1	ACLs and UNIX Mode Bits.....	221
A-2	Changing <i>other</i> Mode Bits .....	222
A-3	Changing <i>group</i> Mode Bits on a File with <b>ACL</b> and No <i>mask_obj</i> .....	222
A-4	Changing the <i>user_obj</i> Entry .....	222
A-5	Changing the <i>group_obj</i> Entry .....	223
A-6	Changing the <i>group_obj</i> Entry in an ACL with a <i>mask_obj</i> Entry .....	223
A-7	The Results of <b>chmod</b> of the <i>group</i> Bits .....	224
A-8	The Results of <b>acl_edit</b> of the <i>mask_obj</i> Bits .....	224
A-9	Changing the <i>group_obj</i> Entry for an ACL with a <i>mask_obj</i> .....	225
A-10	File Creation with no <b>IO ACL</b> .....	226
A-11	File Creation with an <b>IO ACL</b> .....	226
A-12	File <b>ACL</b> Creation with an <b>initial object ACL</b> .....	227
C-1	xglue_rename() Function.....	246
C-2	xglue_lookup() Function .....	247
C-3	xufs_lookup() Function.....	248
13-1	Initializing the Extended VNOPS Vector.....	290
13-2	Generating New Style Vnode Operations From Original .....	291
13-3	GetNewVnodeOpsFromOld Routine .....	292
16-1	Protocol Exporter Access to a Fully Functional LFS.....	417
16-2	Protocol Exporter Access to a non-LFS Filesystem .....	417
16-3	Local OS to a Fully Functional LFS .....	417
16-4	Local OS Access to a Non-LFS Filesystem Such as UFS.....	418

**List of Tables**

12-1	<b>ACL Entry Types</b> .....	254
------	------------------------------	-----



# Preface

## **The Open Group**

In February 1996, X/Open and the Open Software Foundation (OSF) joined forces to become The Open Group, which represents one of the leading authorities in open systems. It is supported by most of the world's largest information systems suppliers, user organisations and software companies. By combining their complementary strengths — X/Open in providing specifications and its trade mark branding scheme, and OSF in facilitating collaboration among customers and system/software vendors toward the development of new open systems technologies — The Open Group is well positioned to assist vendors and users to develop and implement products which support adoption and proliferation of open systems.

Established in 1984, X/Open is an organisation dedicated to the identification, agreement and widescale adoption of Information Technology standards which provide compatibility (portability and interoperability) between software products, and so help users realise the business benefits of open information systems — lower costs, increased choice and greater flexibility. It achieves this by combining existing and emerging standards into an evolving set of integrated, high-value and usable open system specifications, which form a Common Applications Environment (CAE). It licenses a trade mark — called the X/Open Brand — for use on products which vendors have demonstrated are conformant to this CAE. The X/Open brand is recognised by users worldwide as the guarantee of compliance to open systems standards. X/Open is also responsible for the management of the UNIX trade mark on behalf of the industry.

Founded in 1988, the Open Software Foundation (OSF) delivers technology innovations in all areas of open systems, including interoperability, scalability, portability and usability. OSF is a worldwide coalition of vendors and customers in industry, government and academia, who work together to provide advanced open-systems technology solutions for use in a distributed computing environment. It runs programmes in collaborative research and development, to deliver vendor-neutral technology (software source code and supporting documentation) in key IT areas. These include OSF/1, Distributed Computing Environment (DCE), OSF/Motif and Common Desktop Environment (CDE).

## **X/Open Technical Publications**

X/Open publishes a wide range of technical literature, the main part of which is focused on specification development, but which also includes Guides, Snapshots, Technical Studies, Branding/Testing documents, industry surveys and business titles.

There are two types of X/Open specification:

- *CAE Specifications*

CAE (Common Applications Environment) specifications are the stable specifications that form the basis for X/Open-branded products. These specifications are intended to be used widely within the industry for product development and procurement purposes.

Anyone developing products that implement an X/Open CAE specification can enjoy the benefits of a single, widely supported industry standard. In addition, they can demonstrate product compliance through the X/Open brand trade mark. CAE specifications are published as soon as they are developed, so enabling vendors to proceed with development of conformant products without delay.

- *Preliminary Specifications*

These specifications usually address an emerging area of technology and consequently are not yet supported by multiple sources of stable conformant implementations. They are published for the purpose of validation through implementation of products. A Preliminary specification is not a draft specification; rather, it is as stable as can be achieved, through applying X/Open's rigorous development and review procedures.

Preliminary specifications are analogous to the *trial-use* standards issued by formal standards organisations, and developers are encouraged to develop products on the basis of them. However, experience through implementation work may result in significant (possibly upwardly incompatible) changes before its progression to becoming a CAE specification. While the intent is to progress Preliminary specifications to corresponding CAE specifications, the ability to do so depends on consensus among X/Open members.

In addition, X/Open publishes:

- *Guides*

These provide information that X/Open believes is useful in the evaluation, procurement, development or management of open systems, particularly those that are X/Open-compliant. X/Open Guides are advisory, not normative, and should not be referenced for purposes of specifying or claiming X/Open conformance.

- *Technical Studies*

X/Open Technical Studies present results of analyses performed on subjects of interest in areas relevant to X/Open's Technical Programme. They are intended to communicate the findings to the outside world so as to stimulate discussion and activity in other bodies and the industry in general.

- *Snapshots*

These provide a mechanism for X/Open to disseminate information on its current direction and thinking, in advance of possible development of a Specification, Guide or Technical Study. The intention is to stimulate industry debate and prototyping, and solicit feedback. A Snapshot represents the interim results of a technical activity. Although at the time of its publication, there may be an intention to progress the activity towards development of a Specification, Guide or Technical Study, the ability to do so depends on consensus among X/Open members.

### **Versions and Issues of Specifications**

As with all *live* documents, CAE Specifications require revision to align with new developments and associated international standards. To distinguish between revised specifications which are fully backwards compatible and those which are not:

- A new *Version* indicates there is no change to the definitive information contained in the previous publication of that title, but additions/extensions are included. As such, it *replaces* the previous publication.
- A new *Issue* indicates there is substantive change to the definitive information contained in the previous publication of that title, and there may also be additions/extensions. As such, both previous and new documents are maintained as current publications.

### Corrigenda

Readers should note that Corrigenda may apply to any publication. Corrigenda information is published on the World-Wide Web at <http://www.xopen.org> under Sales and Ordering.

### Ordering Information

Full catalogue and ordering information on all Open Group publications is available on the World-Wide Web at <http://www.xopen.org> under Sales and Ordering.

### This Document

This document is a Preliminary Specification (see above). It specifies Distributed File System (DFS) services, interface, protocols, encoding rules and the Interface Definition Language (IDL).

The purpose of this document is to provide a portability guide for DFS application programs and a conformance specification for DFS implementations, particularly those with extensions to the layer known as the Virtual File System (VFS). These extensions will be known as the Virtual File System additions (VFS+).

### Structure

This document is organised into six parts.

Part 1, *Distributed File System Introduction*, describes this volume in detail, covering the DFS Interface Definition Language and data, DFS Versioning Scheme, DFS Client/Server and Token Manager Interfaces, the Access Control List Model used, the operations on aggregates (collections of filesets), the operations on volumes (filesets), and the vnode (VFS) interface and its operations. It contains material relevant to both application programmers and implementors.

Part 2, *DFS Interface Definition Language and Data*, describes the File Exporter IDL Declarations, Cache Manager IDL Declarations, and Common IDL Data used by DFS.

Part 3, *DFS Versioning Scheme*, describes the DFS RPC Versioning Scheme that permits providers to supply enhancements to the RPC functions of DFS. It includes data structures; the DFS versioning API; a set of examples including IDL files, Client, Server, and Manager applications; and a set of rules to follow for extending the DFS interface.

Part 4, *DFS Client/Server and Token Manager Interfaces*, provides a set of manpages for both the DCE DFS File Exporter and the Token Manager.

Part 5, *Access Control Lists (ACLs) in DFS*, provides an overview of Access Control Lists in DFS, their formats and the interface provided to manipulate them.

Part 6, *DCE DFS VFS+ Interface Specification*, specifies a portable Application Programmer's Interface (API). It contains material on the DFS aggregate operations, volume operations, and vnode (VFS) operations that are relevant to both application programmers and implementors.

### Intended Audience

This document is written for DFS application programmers and developers of DFS implementations.

## Typographical Conventions

The following typographical conventions are used throughout this document:

- **Bold** font is used for text for options to commands, filenames, keywords, type names, data structures and their members.
- *Italic* strings are used for emphasis or to identify the first instance of a word requiring definition. Italics in text also denote:
  - variable names; for example, substitutable argument prototypes
  - environment variables, which are also shown in capitals
  - utility names
  - external variables, such as *errno*
  - functions; these are shown as follows: *name()*.
- Normal font is used for the names of constants and literals.
- The notation **<file.h>** indicates a header file.
- The notation [EABCD] is used to identify an error value EABCD.
- Syntax, code examples and user input in interactive examples are shown in `fixed width font`.
- Variables within syntax statements are shown in *italic fixed width font*.

## *Trade Marks*

OSF™ is a trade mark of The Open Software Foundation, Inc.

Transarc®, Encina® and AFS® are registered trade marks of Transarc Corporation.

UNIX® is a registered trade mark in the United States and other countries, licensed exclusively through X/Open Company Limited.

X/Open® is a registered trade mark, and the “X” device is a trade mark, of X/Open Company Limited.

This list represents, as far as possible, those products that are trade marked. X/Open acknowledges that there may be other products that might be covered by trade mark protection and advises the reader to verify them independently.

# *Acknowledgements*

*Part 6, The DCE DFS VFS+ Interface Specification of the DCE1.1:DistributedFileServiceSpecification* is based upon the specification of the extensions to the Virtual File System (VFS) from Hewlett-Packard known as the VFS+ Specification, which was originally written by Daryl Kinney and given to the Open Software Foundation in July, 1993, with subsequent updates by Hewlett-Packard.



## *Referenced Documents*

The following documents are referenced in this specification:

DCE RPC

X/Open CAE Specification, August 1994, X/Open DCE: Remote Procedure Call (ISBN: 1-85912-041-5, C309).

OSF DCE DFS Administration Guide and Reference  
for OSF DCE 1.1.

OSF DCE Administration Guide — Core Components  
for OSF DCE 1.1.



# ***X/Open Preliminary Specification***

## **Part 1:**

### **Distributed File System Introduction**



# Introduction to the DFS Specification

This document specifies both portability and interoperability for the Distributed File System (DFS). The specification contains material directed at two audiences:

- It provides a portability guide for certain application programmers.
- It provides both portability and interoperability specifications for those who are implementing their own Virtual File System Extended (VFS+) or porting the DFS VFS+.

This document may be thought of as an implementation specification, covering both portability and interoperability, that contains within it a portability guide for those that are porting the DFS VFS+ LFS, as well as an implementation specification for those who choose not to port the DFS VFS+ portion.

The implementation specification for the DFS VFS+ LFS is contained in a manner that is intended to be self sufficient. It is mostly self describing, although it does reference other parts of the portability specification; for instance, the information on ACLs and data types.

**Note:** As a result of the intent to be self sufficient, there is some repetition of information, particularly ACL fields, and access algorithms. Thus, chapter 12 describes ACL information that was presented in chapters 8 and 9 from the point of view of an LFS (Local File System).

## 1.1 Portability

### 1.1.1 Overview

This part of the specification consists of chapters 2 through 7. These chapters describe all RPC interfaces used within the kernel resident components of the DCE DFS distributed file system, namely, the file exporter, the cache manager, and the rules for adding new or modified RPC functions to the file exporter. Chapter 6 on page 99 describes the RPC interface to the file exporter, which is used by cache managers to access and manipulate files. Chapter 7 on page 145 describes the RPC interface the cache manager exports in order to handle token management and other administrative inquiries from the file exporter. Chapter 5 on page 69 describes the mechanism for adding new or modified RPC functions to the file exporter.

### 1.1.2 IDLs, Data Types, Constants and Flags

In addition, a set of chapters provides related information such as data types, definitions of constants and flags. Chapter 2 on page 13 describes the IDL definitions and related constants and flags associated with the *DCE DFS File Exporter Interface*. In addition, Chapter 3 on page 51 describes the same kinds of information associated with the *Cache Manager Service Interface*. There is additional data type, flag, and constant information relating to both interfaces in Chapter 4 on page 57, titled *Common IDL Data*.

## 1.2 Document Organization

The **DCE 1.1: Distributed File Service** specification is organised as described in the following. This organisation follows the pattern that the first part is at the lowest level (RPC) and proceeds to higher levels, such as Access Control Lists, and finally, the extensions to the VFS known as the VFS+ Interface.

### 1.2.1 RPC Interfaces for DCE DFS

Parts 2 through 4 of the **DCE 1.1: Distributed File Service** specification describe the exported interface used within the kernel resident portion of the **DCE DFS** distributed file system. This interface description includes all RPC interfaces used within the kernel and the rules for adding new or modified RPC functions to the file exporter (sometimes called the protocol exporter) and consists of:

- Part 2, *DFS Interface Definition Language and Data*, consisting of Chapter 2 on page 13, the *File Exporter IDL Declarations* for both clients and servers, Chapter 3 on page 51, *Cache Manager IDL Declarations*, and Chapter 4 on page 57, *Common IDL Data* shared by both the File Exporter and the Cache Manager
- Part 3, *DFS Versioning Scheme*, consisting of Chapter 5 on page 69, *DFS RPC Versioning Scheme*, which provides a mechanism and rules for extending the DFS interface at the RPC level
- Part 4, *DFS Client/Server and Token Manager Interfaces*, consisting of Chapter 6 on page 99, *DCE DFS File Exporter Interface* and Chapter 7 on page 145, *Cache Manager Service Interface* describing the manpages for these interfaces.

### 1.2.2 Access Control Lists (ACLs) in DFS

Part 5 of the **DCE 1.1: Distributed File Service** specification describes the Access Control Lists used in **DCE DFS**. It consists of:

- Chapter 8 on page 155, *Access Control List Overview*, describing the *ACL* types, and rules governing entries in them for filesystem *ACLs*, including delegation, as well as access checking algorithms for principals and delegation
- Chapter 9 on page 165, *ACL Storage Format*, describing the structured formats used for *ACLs* in memory (which differs from the unstructured byte-stream format used on disk)
- Chapter 10 on page 171, *Access Control List Interface Functions*, describing the functions available in the *ACL* interface by which the *ACLs* can be manipulated. The functions are used for both in-kernel filesystem *ACLs* and Administrative Lists *ACLs*, which are used in user space.

In addition, there are two appendixes in Part 5:

- Appendix A on page 221, *Mapping DFS ACLs to UNIX mode bits*, which discusses the relationship between *ACLs* and UNIX mode bits
- Appendix B on page 229, *Access Control List Package Error List*, which lists the return values that can occur in this interface.

### 1.2.3 The DCE DFS VFS+ Interface Specification

Part 6 of the **DCE 1.1: Distributed File Service** specification describes the **DCE DFS VFS+** interface specification for a **DCE LFS**. It consists of chapters 11 through 17, and 10 appendixes, organized as follows:

- Chapter 11 on page 233, *DCE DFS VFS+ Interface Introduction*, consists of an overview of the interfaces for an LFS, the organization of the VFS+ switch for extending the VFS interface, the fileset registry, fileset and vnode interfaces, and LFS references back into DFS.

It has an appendix, Appendix C, *Components of a Typical VFS+ Package*, which provides information on items such as naming conventions in this interface extension, example wrapper routines called "glue" functions, that can be used to transparently extend the VFS vector so that the extensions can be invoked, and extensions to the **vnode** attributes.

- Chapter 12 on page 251, *The DCE DFS ACL Model for an LFS*, discusses **ACLs**, security and protection checking as they relate to a DFS physical file system (known as a local file system, or LFS).
- Chapter 13 on page 267, *VFS+ Data Types*, describes the data types used in a **DCE LFS**, including primitive data types, as well as data types for aggregates, filesets and vnodes.
- Chapter 14 on page 303, *Aggregate Operations Interface*, describes the portion of the VFS+ interface which allows operations on aggregates, the **DCE DFS** data containers which house filesets, as introduced in Section 11.3 on page 235. The aggregate array functions found in the **struct aggrops** structure, are also described in this chapter.

Additional information for the aggregate interface is found in Appendix D through Appendix F.

- Chapter 15 on page 325, *Fileset (Volume) Operations Interface*, describes the facilities provided for the VFS+ Fileset Operations, including fileset types, clones, fileset and clone requirements, fileset indices, LFS modification of fileset status, zero link count files, quotas, anode generation numbers, file identifiers, vnode to LFS association, and lastly, fileset and *fileset registry* operations.

Appendix G through Appendix L provide additional information about filesets, such as values for opening a fileset, getting status, and creating a fileset.

- Chapter 16 on page 415, *VFS (Vnode) Interface and Operations*, describes the enhanced **vnode** operations vector, converted **vnodes**, synchronization between **vnode** and **fileset** operations, as well as the functions provided for operating on **vnodes**. It also includes the interface description for the extended **vnode** interface.
- Chapter 17 on page 469, *DCE DFS VFS+ Extended Credential Package*, describes the interface used for authentication within the VFS+ layer. It presents an extensible and upwards-compatible interpretation of the VFS UNIX **credential** structure.

It describes what is essentially a subroutine library allowing VFS functions to interpret the VFS credentials structure in a more general manner.

## 1.3 Terminology

This is a list of commonly used terms in DFS.

### ACL

Short for *access control list*; a description of which users (*principals*) and groups (groups of principals) can use a file or directory and what types of uses are permitted (*rights*). Every file can have an **ACL**. Directories have three **ACL** types; one governs operations on the directory itself, a second provides an initial **ACD** for files created in the directory, and a third provides the **ACL** for sub-directories created within the directory. These **ACL** types are known as the *initial object*, *object*, and *initial container ACLs*, respectively.

### Aggregate

An expanse of physical disk which is managed by a DFS LFS. It is similar in context to a UNIX partition and is identified by both a device number and a device file. In many cases, the storage provided to an aggregate will be by a Logical Volume Manager (LVM).

### Anode

A 252-byte area of disk which is contiguous. It is the primitive data structure in the **DCE LFS** file system. All objects in LFS reside in a *container*; a logical variable-sized object, sometimes called a "bucket", that can hold data. An *anode* stores status and disk addressing information -- called meta-data as well as information describing user files.

### Anode Index

An index which can be used to refer to the *anode*. This index is used in many situations by an LFS to name and locate *anodes*. A *fileset* is implemented as an array of *anodes*.

### Backing Anode

This is an *anode* whose storage is shared with other *COW anodes* and whose contents are fixed. Backing anodes are **read-only**. Multiple *COW anodes* can refer to the same backing anodes.

### Block

The larger unit in a system, for allocating disk space to *containers* and doing I/O and buffering. Block size must be a power of two in **DCE DFS**.

### Cell

A collection of DCE machines administered as a single entity. In this context, the key point is that a cell is serviced by a single registry with a single, consistent set of user and group identifiers.

### Clone

A loosely-used term to describe a *fileset* which is related to another *fileset*. Sometimes it refers to a read-write (R/W) fileset and sometimes to a *replica*. The term clone is also applied to the operation on individual anodes during the fileset clone process. A R/W fileset is comprised of *COW anodes* each of which was "cloned" from the *backing anode* in the corresponding *anode index* in a read-only snapshot of a fileset. A snapshot is a consistent, read-only version of a fileset.

### COW Anode

Also known as a Copy-on-Write anode. This is an *anode* whose storage is shared with a *backing anode* and whose contents can be changed. Modifications to a COW anode are made by allocating new storage, initializing the storage from the backing anode, then making the desired modifications. Both containers (anodes) must reside in the same *aggregate*.

### Container

An open-ended, ordered collection of disk space; the basic storage unit for data on a disk in a **DCE LFS**. It is variable-sized and is represented by an anode. A container can contain



other containers. For instance, an *aggregate* is a container that holds other containers, such as filesets and also representations of free space, or a log. A *fileset*, in turn, is a container holding files.

**DFS**

The DCE Distributed File System.

**EFS**

The optional portion of DFS which deals with extended fileset operations (cloning, backup, replicas, and so on).

**Xcred**

Also known as an *extended credential*. A credential structure, like that used in the UNIX kernel, but augmented in various ways - including a *realm* identifier; it distinguishes a *local group* from a *foreign group*. The xcred credential identifies a *principal*.

**File**

A *container* used to implement UNIX files, directories and symbolic links. It has a status area in which the standard information returned by **stat(2)** is stored. The status area also contains pointers to information for access control.

**Fileset**

A set of related files, connected via a sub-tree with a single root point, that is administered as an entity by DFS. Although filesets cannot span aggregates, a single aggregate can hold many filesets. Filesets which have an associated maximum size or quota, grow independently of each other and compete for space within their aggregate. Individual filesets can be mounted locally, backed up and restored, cloned or transparently moved to another aggregate.

The terms fileset and volume are used interchangeably, although the term volume is becoming archaic.

**Fragment**

The smaller unit in a system, for allocating disk space to *containers*. The number of fragments in a *block* must be a power of two.

**Group**

A membership list, identified by a **uid\_t**, which contains a list of principals (by **uid\_t**) which are members of that group. Again, in a non-DCE environment, they are identified by a 32-bit (at best) **gid**.

**LFS**

Local File System: a physical file system, which in the context of this document, provides DFS semantics.

**Episode** is the name of the LFS that Transarc Corporation provides as part of the DFS EFS package.

**Magic Cookie**

A piece of data, passed back and forth between two functions which is opaque to the higher-level function. It is often used as a token of progress in an iteration.

**PAC**

Privilege Attribute Certificate. A data structure, generated by the DCE Privilege Server, that contains the identity of a DCE authenticated principal. To a first approximation, it contains the principal's **uid\_t** along with the **uid\_t** of each group of which it is a member.

**Principal**

An entity that is interacting with DFS. Although it usually refers to a person, it could as easily correspond to a DCE server or a computer within a DCE cell. In the DCE environment, principals are identified by a **uuid\_t**. In a non-DCE, standard UNIX, environment, they are identified by a 32-bit (at best) **uid**.

**Quota**

Maximum space allowed to be allocated to a *fileset*, independent of how much space remains in the *aggregate*.

**Realm**

Equivalent to a Cell. In security DCE security discussions, the term realm is typically used instead of cell.

**Replica**

A read-only *fileset* whose contents match that of a consistent read-only version of a fileset on a different *aggregate*. A replica differs from a snapshot in that the snapshot is really linked to the primary fileset by a copy on write relationship and shares storage with it, while a replica is an independent copy of a snapshot in another aggregate.

**Rights**

The classes of operations whose access can be separately controlled. For traditional UNIX file systems, these are designated **rwX**. A **DCE LFS** has more rights for directories which allow separate control of *insert* (**i**) and *delete* (**d**) operations. In addition, the operations traditionally reserved to a file's owner are under the control of right (**c**).

**Token**

A data structure passed between a DFS file server and client (user of the file server) that guarantees to the client that it has permission to perform the operation described in the token. The DFS file servers issue and revoke tokens granted to clients, and clients follow this convention. Tokens are used to provide a coherent view of data across all DFS clients as if they were in a single system.

**UFS**

The **UNIX File System**. The name of the file system originally provided by vendors of **UNIX** systems. In this file system the *inode* is a container for user file data and directory data. It does not contain metadata such as found in an *anode*.

**Uniquifier**

Synonym used interchangeably with *unique*. It distinguishes between different uses (over time) of an entity. It is a number incremented each time an entity is created, and can be thought of as a generation number. Its primary use is in the construction of file identifiers. In a *fileset*, it is a number incremented each time a file is created, which is used to distinguish between instantiations of a file using the same *anode index*. In a **vnnode**, for each new use, a new *uniquifier* is used in the **afsFid**, the DFS file identifier.

**VFS+ Interface**

A VFS interface, including **vnnode** operations and **VFS** operations, extended by adding **vnnode** operations for portability, for **Episode**, and for other **DCE LFS** file systems. An extension to this interface usually includes **fileset** operations and **aggregate** operations.

In addition, for portability, this interface requires portability "glue" or wrapper functions to adapt the *vnnode* operations to a particular kernel. This interface also requires token "glue" in order to enable a file system to be exported.

## **1.4 Conformance Requirements**

To conform to this document, implementations must meet the following requirements:

- Implementations of file servers must support the interfaces defined in Chapter 6 on page 99.
- Implementations of file servers must support the AFS4Int interface defined in Chapter 2 on page 13.
- Implementations of file servers must support the common data structures defined in Chapter 4 on page 57.
- Implementations of cache managers must support the interfaces defined in Chapter 7 on page 145.
- Implementations of cache managers must support the TKN4Int interface defined in Chapter 3 on page 51.
- Implementations of cache managers must support the common data structures defined in Chapter 4 on page 57.
- Implementations of Local File Systems must support the interfaces and data structures defined in Part 6 on page 237.
- Implementations that extend either the file exporter interface defined in Chapter 2 on page 13 or the token manager interface defined in Chapter 3 on page 51 (or both) must conform to the mechanism defined in Chapter 5 on page 69.



# ***X/Open Preliminary Specification***

## **Part 2:**

### **DFS Interface Definition Language and Data**



## *File Exporter IDL Declarations*

This chapter specifies the RPC interfaces used by the AFS (4.0) Client and Cache Manager to the AFS (4.0) Server and File Exporter Interface.

The File Exporter interface consists of a set of structures, mask values and flags contained in an IDL file with a provider version number of 1 to indicate that it is an original interface. In addition, each RPC defined in this structure also has a provider version number of 1, indicating that it also is an original RPC. See Chapter 5 on page 69 for information on versioning which will explain in more detail about provider version numbers.

### 2.1 The File Exporter Interface

The File Exporter Interface is contained within the interface, **AFS4Int**, which starts out as shown below and contains all the items shown in this chapter exclusive of the definitions of the parameters in the RPC calls which are included from Chapter 6 on page 99 and the common data definitions in the configuration file **common\_data.idl** which are contained in Chapter 4 on page 57 for clarity. There are twenty seven remote procedure calls defined that the client and cache manager can use in this interface which is defined as provider version 1. Chapter 5 on page 69 provides information for adding new interfaces, or new remote procedure calls to this interface.

```
[
uuid(4d37f2dd-ed93-0000-02c0-37cf1e000000),
version(4.0)
/* provider_version(1) */
]

interface AFS4Int
{

import "dcedfs/common_data.idl";

.
.
.
/* } at end of this chapter concludes this interface. */
```

## 2.2 Mask Values for afsStoreStatus Structure

```

const unsigned32 AFS_SETMODTIME = 1;
const unsigned32 AFS_SETOWNER = 2;
const unsigned32 AFS_SETGROUP = 4;
const unsigned32 AFS_SETMODE = 8;
const unsigned32 AFS_SETACCESSTIME = 0x10;
const unsigned32 AFS_SETCHANGETIME = 0x20;
                                /* don't allow this for now */
const unsigned32 AFS_SETLENGTH = 0x40;
const unsigned32 AFS_SETTYPEUUID = 0x80;
const unsigned32 AFS_SETDEVNUM = 0x100;
const unsigned32 AFS_SETMODEXACT = 0x200;
                                /* allow setting back */
const unsigned32 AFS_SETTRUNCLength = 0x400;
                                /* truncate to this size first */
const unsigned32 AFS_SETCLIENTSPARE = 0x800;
                                /* set clientSpare1 */

```

## 2.3 AFS\_GetTime Sync Constants

These are distance and dispersion constants.

```

const unsigned32 AFS_SYNCUNSYNC = 0x7fffffff0;
const unsigned32 AFS_SYNCINITIAL = 0x7fffffff1;

```

## 2.4 Flag Parameters for ACL Type Parameter

These are flag parameters for the high order 16 bits of the getacl and setacl ACL type parameter.

```

const unsigned32 AFS_ACLFLAG_COPY = 1;
    /* copy the ACL from the specified fid */

```

## 2.5 Definitions for Flag Used in Getting a Token

These are the definitions for the flags used in getting a token for each afs4int call.

```

const unsigned32 AFS_FLAG_RETURNTOKEN = 1;
    /* Return tokens to caller not to the token manager */

const unsigned32 AFS_FLAG_TOKENJUMPQUEUE = 2;
    /* in AFS_GetToken, call tkset_AddTokenSet with
       TKSET_ATS_WANTJUMPQUEUE flag */

const unsigned32 AFS_FLAG_SKIPTOKEN = 4;
    /* don't obtain tokens, client already has them */

const unsigned32 AFS_FLAG_NOOPTIMISM = 0x8;
    /* return only the exact asking token, no "optimism" tokens */

const unsigned32 AFS_FLAG_TOKENID = 0x10;

```



```

    /* to say get-the-token-with-this-ID */

const unsigned32 AFS_FLAG_RETURNBLOCKER = 0x20;
    /* To return blocker's info for Sys V record lock */

const unsigned32 AFS_FLAG_ASYNCGRANT = 0x40;
    /* permit async grant in AFS_GetToken if can't get now */

const unsigned32 AFS_FLAG_NOREVOKE = 0x80;
    /* Ask server not to revoke other's token, if there is a
       conflict. Just return a failure */

```

## 2.6 Definitions for Flag in Token Recovery (TSR)

These are the definitions for the flag used in performing token recovery (TSR) operations.

```

const unsigned32 AFS_FLAG_MOVE_REESTABLISH = 0x100;
    /* to say you're reestablishing tokens after a move */

const unsigned32 AFS_FLAG_SERVER_REESTABLISH = 0x200;
    /* to say you're reestablishing tokens after a server crash */

const unsigned32 AFS_FLAG_NO_NEW_EPOCH = 0x400;
    /* to prevent advancing the indicated token into a new epoch */

const unsigned32 AFS_FLAG_MOVE_SOURCE_OK = 0x800;

```

## 2.7 Definitions for Flag in Fileset Operations

These are the definitions for the flag used in performing fileset operations.

```

const unsigned32 AFS_FLAG_SYNC = 0x1000;
    /* fsync all modified files after operation */

const unsigned32 AFS_FLAG_ZERO = 0x2000;
    /* Data is all zero instead of in pipe */

const unsigned32 AFS_FLAG_SKIPSTATUS = 0x4000;
    /* Dont bother filling out file status */

const unsigned32 AFS_FLAG_FORCEREVOCATIONS = 0x8000;
    /* In AFS_GetToken, insist on getting this token, and make conflicting
       tokens be revoked with the AFS_REVOKE_FORCED_REVOCATION flag. */

const unsigned32 AFS_FLAG_FORCEVOLQUIESCE = 0x10000;
    /* Ensure that this whole-fileset token blocks all per-file tokens. */

```

## 2.8 Definitions for Flags in AFS\_SetParams()

These are the definitions for the flags used only in calling *AFS\_SetParams()*.

```
const unsigned32 AFS_PARAM_RESET_CONN    = 0x1;
const unsigned32 AFS_PARAM_TSR_COMPLETE = 0x2;
```

## 2.9 Definitions for Flags in AFS\_SetContext()

```
const unsigned32 AFS_FLAG_SEC_SERVICE    = 0x1;
    /* client needs a secondary service */

const unsigned32 AFS_FLAG_CONTEXT_NEW_IF = 0x2;
    /* this is a new-interface call, with different DOWN semantics */

const unsigned32 AFS_FLAG_CONTEXT_DO_RESET = 0x4;
    /* client will reset all token state itself, to a
       new-interface server */
const unsigned32 AFS_FLAG_CONTEXT_NEW_ACL_IF = 0x8;
    /* client understands new delegation acl entry types */
```

## 2.10 Definitions for Client-only Attribute Flags

```
const unsigned32 AFS_CLIENTATTR_HIDDEN = 0x1;
const unsigned32 AFS_CLIENTATTR_SYSTEM = 0x2;
const unsigned32 AFS_CLIENTATTR_ARCHIVE = 0x4;
const unsigned32 AFS_CLIENTATTR_READONLY = 0x8;
```

## 2.11 Data Types

The following data types are used to define structures used by the remote procedure calls in provider version 1 of the File Exporter Interface.

### 2.11.1 Define Generic Network Address Information

```
typedef struct afsNetData {
    afsNetAddr sockAddr;
    NameString_t principalName;
} afsNetData;
```

### 2.11.2 Define afsVolSync Structure for Tracking Replicas

This is the structure by which replicas are tracked.

```
typedef struct afsVolSync {
    afsHyper VolID;
    afsHyper VV;          /* volume's version */
    unsigned32 VVAge;     /* age, in seconds, of the knowledge that the
                           given VolVers is current */
    unsigned32 VVPingAge; /* age, in seconds, of the last probe from
                           the callee (the secondary) to the primary */
    unsigned32 vv_spare1;
```

```

    unsigned32 vv_spare2;
} afsVolSync;

```

### 2.11.3 Define the afsFetchStatus Structure

```

typedef struct afsFetchStatus {
    unsigned32    interfaceVersion;
    unsigned32    fileType;
    unsigned32    linkCount;
    afsHyper      length;
    afsHyper      dataVersion;
    unsigned32    author;
    unsigned32    owner;
    unsigned32    group;
    unsigned32    callerAccess;
    unsigned32    anonymousAccess;
    unsigned32    aclExpirationTime;
    unsigned32    mode;
    unsigned32    parentVnode;
    unsigned32    parentUnique;
    afsTimeval    modTime;
    afsTimeval    changeTime;
    afsTimeval    accessTime;
    afsTimeval    serverModTime;
    afsUUID       typeUUID;
    afsUUID       objectUUID;
    unsigned32    deviceNumber;
    unsigned32    blocksUsed;
    unsigned32    clientSpare1;    /* client-only attrs */
    unsigned32    deviceNumberHighBits;
    unsigned32    spare0;
    unsigned32    spare1;
    unsigned32    spare2;
    unsigned32    spare3;
    unsigned32    spare4;
    unsigned32    spare5;
    unsigned32    spare6;
} afsFetchStatus;

```

### 2.11.4 Define the afsStoreStatus Structure

```

typedef struct afsStoreStatus {
    unsigned32    mask;
    afsTimeval    modTime;
    afsTimeval    accessTime;
    afsTimeval    changeTime;
    unsigned32    owner;
    unsigned32    group;
    unsigned32    mode;
    afsHyper      truncLength;    /* applied first */
    afsHyper      length;
    afsUUID       typeUUID;
    unsigned32    deviceType;    /* character or block */
}

```

```

    unsigned32    deviceNumber;
    unsigned32    cmask;
    unsigned32    clientSpare1;    /* client-only attrs */
    unsigned32    deviceNumberHighBits;
    unsigned32    spare1;
    unsigned32    spare2;
    unsigned32    spare3;
    unsigned32    spare4;
    unsigned32    spare5;
    unsigned32    spare6;
} afsStoreStatus;

```

### 2.11.5 Structure for Physical File System Type

```

typedef struct afsDisk { /* used by Statistics */
    unsigned32    BlocksAvailable;
    unsigned32    TotalBlocks;
    unsigned32    spare1;
    unsigned32    spare2;
    unsigned32    spare3;
    afsDiskName    Name;
} afsDisk;

```

### 2.11.6 Structure for Statistics

```

typedef struct afsStatistics {
    unsigned32    CurrentMsgNumber;
    unsigned32    OldestMsgNumber;
    unsigned32    CurrentTime;
    unsigned32    BootTime;
    unsigned32    StartTime;
    unsigned32    CurrentConnections;
    unsigned32    TotalAfsCalls;
    unsigned32    TotalFetchs;
    unsigned32    FetchDatAs;
    unsigned32    FetchedBytes;
    unsigned32    HighFetchedBytes;
    unsigned32    FetchDataRate;
    unsigned32    TotalStores;
    unsigned32    StoreDatAs;
    unsigned32    StoredBytes;
    unsigned32    HighStoredBytes;
    unsigned32    StoreDataRate;
    unsigned32    TotalRPCBytesSent;
    unsigned32    HighTotalRPCBytesSent;
    unsigned32    TotalRPCBytesReceived;
    unsigned32    HighTotalRPCBytesReceived;
    unsigned32    TotalRPCPacketsSent;
    unsigned32    TotalRPCPacketsReceived;
    unsigned32    TotalRPCPacketsLost;
    unsigned32    TotalRPCBogusPackets;
    unsigned32    SystemCPU;
    unsigned32    UserCPU;
}

```

```
    unsigned32    NiceCPU;
    unsigned32    IdleCPU;
    unsigned32    TotalIO;
    unsigned32    ActiveVM;
    unsigned32    TotalVM;
    unsigned32    EtherNetTotalErrors;
    unsigned32    EtherNetTotalWrites;
    unsigned32    EtherNetTotalInterupts;
    unsigned32    EtherNetGoodReads;
    unsigned32    EtherNetTotalBytesWritten;
    unsigned32    EtherNetTotalBytesRead;
    unsigned32    ProcessSize;
    unsigned32    WorkStations;
    unsigned32    ActiveWorkStations;
    unsigned32    Spare1;
    unsigned32    Spare2;
    unsigned32    Spare3;
    unsigned32    Spare4;
    unsigned32    Spare5;
    unsigned32    Spare6;
    unsigned32    Spare7;
    unsigned32    Spare8;
    afsDisk       Disk1;
    afsDisk       Disk2;
    afsDisk       Disk3;
    afsDisk       Disk4;
    afsDisk       Disk5;
    afsDisk       Disk6;
    afsDisk       Disk7;
    afsDisk       Disk8;
    afsDisk       Disk9;
    afsDisk       Disk10;
    afsDisk       Disk11;
    afsDisk       Disk12;
    afsDisk       Disk13;
    afsDisk       Disk14;
    afsDisk       Disk15;
    afsDisk       Disk16;
} afsStatistics;
```

## 2.12 Various Bulk Data Type Definitions for RPC Operations

### 2.12.1 The afsFidExp Structure

This structure is a parameter for BulkKeepAlive, so that the PX doesn't have to look up the lifetimes.

```
typedef struct afsFidExp {
    afsFid fid;
    unsigned32 keepAliveTime;
} afsFidExp;
```

A special value for keepAliveTime to force a re-check of the reclaimDally value.

```
const unsigned32 AFS_KA_TIME_RECHECK = 1;
```

### 2.12.2 The afsBulkFEX Structure

```
typedef struct afsBulkFEX {
    unsigned32 afsBulkFEX_len;
    [length_is(afsBulkFEX_len)] afsFidExp afsBulkFEX_val[AFS_BULKMAX];
} afsBulkFEX;
```

### 2.12.3 The afsACL Structure

```
typedef struct afsACL {
    unsigned32 afsACL_len;
    [length_is(afsACL_len)] byte afsACL_val[AFS_ACLMAX];
} afsACL;
```

### 2.12.4 The afsQuota Structure

```
typedef struct afsQuota {
    unsigned32 afsQuotaType;    /* BSD, Episode, etc */
    unsigned32 afsQuotaOp;     /* set or get */
    unsigned32 afsQuota_len;
    [length_is(afsQuota_len)] unsigned32 afsQuota_val[AFS_BULKMAX];
} afsQuota;
```

### 2.12.5 The afsBulkVVs Structure

```
typedef struct afsBulkVVs {
    unsigned32 afsBulkVVs_len;
    [length_is(afsBulkVVs_len)] afsVolSync afsBulkVVs_val[AFS_BULKMAX];
} afsBulkVVs;
```

**2.12.6 The afsBulkVolIDs Structure**

```
typedef struct afsBulkVolIDs {  
    unsigned32 afsBulkVolIDs_len;  
    [length_is(afsBulkVolIDs_len)] afsHyper afsBulkVolIDs_val[AFS_BULKMAX];  
} afsBulkVolIDs;
```

**2.12.7 The afsBulkStats Structure**

```
typedef struct afsBulkStats {  
    unsigned32 afsBulkStats_len;  
    [length_is(afsBulkStats_len)] afsFetchStatus afsBulkStats_val[AFS_BULKMAX];  
} afsBulkStats;
```

## 2.13 Definitions of File Server Exported Operations

**Note:** The Opcode of each operation is implicitly assigned based on the order of where the operation is placed. The recommendation is to always place the new operation at the end of this list.

### 2.13.1 AFS\_SetContext

```
error_status_t AFS_SetContext
(/* provider_version(1) */
 [in] handle_t           h,
 [in] unsigned32        epochTime,
 [in] afsNetData        *callbackAddr,
 [in] unsigned32        Flags,
 [in] afsUUID           *secObjectID,
 [in] unsigned32        clientSizeAttrs,
 [in] unsigned32        parm7
);
```

**AFS\_SetContext** parameters are:

<i>h</i>	The RPC binding handle.
<i>epochTime</i>	The restart time of the DFS client.
<i>callbackAddr</i>	The RPC endpoint of the client, for token revocation purposes.
<i>Flags</i>	If 0 then the call only defines a primary interface UUID. If the flag, AFS_FLAG_SEC_SERVICE, is set then a secondary interface is defined and is stored in the parameter, <i>secObjectID</i> .
<i>secObjectID</i>	When the AFS_FLAG_SEC_SERVICE flag is set, this parameter holds the secondary interface UUID.
<i>clientSizeAttrs</i>	For 64/32-bit compatibility. Through DCE1.1, this is a spare parameter, <i>parm6</i> .
<i>parm7</i>	A spare parameter.



## 2.13.2 AFS\_LookupRoot

```

error_status_t AFS_LookupRoot
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *InFidp,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsFid          *OutFidp,
 [out] afsFetchStatus  *OutFidStatusp,
 [out] afsToken        *OutTokenp,
 [out] afsVolSync      *Syncp
);

```

AFS\_LookupRoot parameters are:

<i>h</i>	The RPC binding handle.
<i>InFidp</i>	The file identifier specifying the fileset whose root directory will be retrieved from the file server.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier describing the root directory of the selected fileset.
<i>OutFidStatusp</i>	Returns the status of that directory after the current operation.
<i>OutTokenp</i>	The promise the file server returns to the cache manager about the provided data; this is only returned if the file resides in a read/write fileset.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.3 AFS\_FetchData

```

error_status_t AFS_FetchData
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *Fidp,
 [in] afsHyper         *minVVp,
 [in] afsHyper         *Position,
 [in] signed32         Length,
 [in] unsigned32       Flags,
 [out] afsFetchStatus  *OutStatusp,
 [out] afsToken        *OutTokenp,
 [out] afsVolSync      *Syncp,
 [out] pipe_t          *fetchStream
);

```

**AFS\_FetchData** parameters are:

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file whose contents will be retrieved from the file server.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Position</i>	Specifies the first byte to be fetched by this call with zero (0) being the first byte in the file.
<i>Length</i>	Specifies the number of bytes desired with the value 0xFFFFFFFF indicating the entire file contents.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutStatusp</i>	Returns the status of the file after the current operation.
<i>OutTokenp</i>	The promise the file server returns to the cache manager about the provided data; this is only returned if the file resides in a read/write fileset.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.
<i>fetchStream</i>	The character pipe parameter returning the data from the file. (See <b>X/Open DCE: Remote Procedure Call</b> specification for implementation details.)

## 2.13.4 AFS\_FetchACL

```

error_status_t AFS_FetchACL
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *Fidp,
 [in] unsigned32       aclType,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsACL          *AccessListp,
 [out] afsFetchStatus  *OutStatusp,
 [out] afsVolSync      *Syncp
);

```

AFS\_FetchACL parameters are:

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file or directory whose access control list (ACL) will be retrieved from the file server.
<i>aclType</i>	The type of the access list being modified. One of VNX_ACL_REGULAR_ACL, VNX_ACL_DEFAULT_ACL or VNX_ACL_INITIAL_ACL.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this ACL.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>AccessListp</i>	The access control list returned by the file server for the specified file or directory.
<i>OutStatusp</i>	Returns the current status of the file or directory.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.5 AFS\_FetchStatus

```

error_status_t AFS_FetchStatus
(/* provider_version(1) */
 [in] handle_t           h,
 [in] afsFid            *Fidp,
 [in] afsHyper          *minVVp,
 [in] unsigned32        Flags,
 [out] afsFetchStatus   *OutStatusp,
 [out] afsToken         OutTokenp,
 [out] afsVolSync       *Syncp
);

```

AFS\_FetchStatus parameters are:

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file or directory whose status information will be retrieved from the file server.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutStatusp</i>	Returns the current status of the file or directory.
<i>OutTokenp</i>	Returns a token from the file server reflecting guarantees granted by the file server.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.6 AFS\_StoreData

```

error_status_t AFS_StoreData
(* provider_version(1) *)
[in] handle_t          h,
[in] afsFid           *Fidp,
[in] afsStoreStatus   *InStatusp,
[in] afsHyper         *Position,
[in] signed32         Length,
[in] afsHyper         *minVVp,
[in] unsigned32       Flags,
[in] pipe_t           *storeStream,
[out] afsFetchStatus  *OutStatusp,
[out] afsVolSync      *Syncp
);

```

AFS\_StoreData parameters are:

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file whose status information will be updated from the file server.
<i>InStatusp</i>	The new status information that should be recorded for this file.
<i>Position</i>	Represents the position of the first byte of the data block.
<i>Length</i>	Represents the total length of the transferred data block.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>storeStream</i>	The data stream containing the file updates. (See the <b>X/Open DCE: Remote Procedure Call</b> specification for implementation details.)
<i>OutStatusp</i>	Returns the status of the file.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.7 AFS\_StoreACL

```

error_status_t AFS_StoreACL
(* provider_version(1) */
[in] handle_t          h,
[in] afsFid           *Fidp,
[in] afsACL           *AccessListp,
[in] unsigned32       aclType,
[in] afsFid           *aclFidp,
[in] afsHyper         *minVVp,
[in] unsigned32       Flags,
[out] afsFetchStatus *OutStatusp,
[out] afsVolSync      *Syncp
);

```

AFS\_StoreACL parameters are:

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file or directory whose access control information will be updated from the file server.
<i>AccessListp</i>	The access control list sent to the file server for the specified file or directory.
<i>aclType</i>	The type of access control list being modified (VNX_ACL_REGULAR_ACL, VNX_ACL_DEFAULT_ACL or VNX_ACL_INITIAL_ACL, in the low-order 8 bits. In the next-higher-order 8 bits, the type of access control list being copied from the file described by <i>aclFidp</i> .
<i>aclFidp</i>	The file identifier specifying the file or directory whose access control information will be copied.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutStatusp</i>	Returns the updated file or directory status.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.8 AFS\_StoreStatus

```

error_status_t AFS_StoreStatus
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *Fidp,
 [in] afsStoreStatus   *InStatusp,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsFetchStatus  *OutStatusp,
 [out] afsVolSync      *Syncp
);

```

AFS\_StoreStatus parameters are:

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file or directory whose access control information will be updated from the file server.
<i>InStatusp</i>	Contains the new status information for the specified file or directory.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutStatusp</i>	Contains the updated status information for the file or directory.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.9 AFS\_RemoveFile

```

error_status_t AFS_RemoveFile
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *DirFidp,
 [in] afsFidTaggedName *Namep,
 [in] afsHyper         *returnTokenIDp,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsFetchStatus  *OutDirStatusp,
 [out] afsFetchStatus  *OutFileStatusp,
 [out] afsFid          *OutFileFidp,
 [out] afsVolSync      *Syncp
);

```

AFS\_RemoveFile parameters are:

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory from which to remove the file.
<i>Namep</i>	The complex name of the file to delete.
<i>returnTokenIDp</i>	A token ID being returned, if any.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutDirStatusp</i>	Contains the updated directory status information.
<i>OutFileStatusp</i>	Contains the updated file status information.
<i>OutFileFidp</i>	The file id of the file which was actually deleted.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.



## 2.13.10 AFS\_CreateFile

```

error_status_t AFS_CreateFile
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *DirFidp,
 [in] afsTaggedName    *Namep,
 [in] afsStoreStatus   *InStatusp,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsFid          *OutFidp,
 [out] afsFetchStatus  *OutFidStatusp,
 [out] afsFetchStatus  *OutDirStatusp,
 [out] afsToken        *OutTokenp,
 [out] afsVolSync      *Syncp
);

```

AFS\_CreateFile parameters are:

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which to create the requested file.
<i>Namep</i>	The character string name of the file to create.
<i>InStatusp</i>	Specifies the initial status fields for the new file.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier of the newly created file.
<i>OutFidStatusp</i>	The status fields for the newly created file.
<i>OutDirStatusp</i>	The status information of the specified directory.
<i>OutTokenp</i>	A new token granted against the new file.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.11 AFS\_Rename

```

error_status_t AFS_Rename
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid            *OldDirFidp,
 [in] afsFidTaggedName *OldNamep,
 [in] afsFid            *NewDirFidp,
 [in] afsFidTaggedName *NewNamep,
 [in] afsHyper          *returnTokenIDp,
 [in] afsHyper          *minVVp,
 [in] unsigned32        Flags,
 [out] afsFetchStatus   *OutOldDirStatusp,
 [out] afsFetchStatus   *OutNewDirStatusp,
 [out] afsFid            *OutOldFileFidp,
 [out] afsFetchStatus   *OutOldFileStatusp,
 [out] afsFid            *OutNewFileFidp,
 [out] afsFetchStatus   *OutNewFileStatusp,
 [out] afsVolSync       *Syncp
);

```

AFS\_Rename parameters are:

<i>h</i>	The RPC binding handle.
<i>OldDirFidp</i>	The file identifier specifying the directory in which the file is currently located.
<i>OldNamep</i>	The complex name of the file or directory to rename.
<i>NewDirFidp</i>	The file identifier specifying the directory into which the file is to be moved.
<i>NewNamep</i>	The complex name of the file or directory after it is moved.
<i>returnTokenIDP</i>	A token ID being returned, if any.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutOldDirStatusp</i>	The status information of the old directory (the one from which the file is being moved) upon termination of the call.
<i>OutNewDirStatusp</i>	The status information of the new directory (the one to which the file is being moved) upon termination of the call.
<i>OutOldFileFidp</i>	The file identifier for the file which was moved.
<i>OutOldFileStatusp</i>	The status information of the file which was moved.
<i>OutNewFileFidp</i>	The file identifier for the file to which the file identified in <i>OutOldFileFidp</i> was, in fact, moved.
<i>OutNewFileStatusp</i>	The status information of the file identified by <i>OutNewFileFidp</i> .
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.12 AFS\_Symlink

```

error_status_t AFS_Symlink
(* provider_version(1) */
[in] handle_t          h,
[in] afsFid           *DirFidp,
[in] afsTaggedName   *Namep,
[in] afsTaggedPath   *LinkContentsp,
[in] afsStoreStatus  *InStatusp,
[in] afsHyper        *minVVp,
[in] unsigned32      Flags,
[out] afsFid         *OutFidp,
[out] afsFetchStatus *OutFidStatusp,
[out] afsFetchStatus *OutDirStatusp,
[out] afsToken       *OutTokenp,
[out] afsVolSync     *Syncp
);

```

AFS\_Symlink parameters are:

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the symbolic link is to be created.
<i>Namep</i>	The name of the link to create.
<i>LinkContentsp</i>	The target of the new symbolic link.
<i>InStatusp</i>	This specifies the clientModTime field and unixModeBits of the new link.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier for the newly created symbolic link.
<i>OutFidStatusp</i>	The status information of the newly created symbolic link upon termination of the call.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the symbolic link was created) upon termination of the call.
<i>OutTokenp</i>	The token returned against the specified directory.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.13 AFS\_HardLink

```

error_status_t AFS_HardLink
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *DirFidp,
 [in] afsTaggedName    *Namep,
 [in] afsFid           *ExistingFidp,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsFetchStatus  *OutFidStatusp,
 [out] afsFetchStatus  *OutDirStatusp,
 [out] afsVolSync      *Syncp
);

```

AFS\_HardLink parameters are:

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the file is currently located.
<i>Namep</i>	The name to use for the new hard link.
<i>ExistingFidp</i>	The file identifier specifying the file identifier of the file to which the hard link should be made.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidStatusp</i>	The status information of the newly created hard link upon termination of the call.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the hard link was created) upon termination of the call.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.14 AFS\_MakeDir

```

error_status_t AFS_MakeDir
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *DirFidp,
 [in] afsTaggedName    *Namep,
 [in] afsStoreStatus   *InStatusp,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsFid          *OutFidp,
 [out] afsFetchStatus  *OutFidStatusp,
 [out] afsFetchStatus  *OutDirStatusp,
 [out] afsToken        *OutTokenp,
 [out] afsVolSync      *Syncp
);

```

AFS\_MakeDir parameters are:

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the new directory is to be created.
<i>Namep</i>	The name of the directory to create.
<i>InStatusp</i>	This specifies the new status information, including clientModTime field and unixModeBits, of the new directory point.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier for the newly created directory.
<i>OutFidStatusp</i>	The status information of the newly created directory upon termination of the call.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the directory was created) upon termination of the call.
<i>OutTokenp</i>	A new token granted against the newly created directory.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.15 AFS\_RemoveDir

```

error_status_t AFS_RemoveDir
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *DirFidp,
 [in] afsFidTaggedName *Namep,
 [in] afsHyper         *returnTokenIDp,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsFetchStatus  *OutDirStatusp,
 [out] afsFid           *OutFidp,
 [out] afsFetchStatus  *OutDelStatusp,
 [out] afsVolSync      *Syncp
);

```

AFS\_RemoveDir parameters are:

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the directory to be deleted is located.
<i>Namep</i>	The complex name of the directory to delete.
<i>returnTokenIDp</i>	A token ID being returned, if any.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the directory was deleted) upon termination of the call.
<i>OutFidp</i>	
<i>OutDelStatusp</i>	
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## 2.13.16 AFS\_Readdir

```

error_status_t AFS_Readdir
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *DirFidp,
 [in] afsHyper         *Offsetp,
 [in] unsigned32       Size,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsHyper        *NextOffsetp,
 [out] afsFetchStatus  *OutDirStatusp,
 [out] afsToken        *OutTokenp,
 [out] afsVolSync      *Syncp,
 [out] pipe_t          *dirStream
);

```

AFS\_Readdir parameters are:

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file descriptor for the specified directory.
<i>Offsetp</i>	The offset into the directory for this entry.
<i>Size</i>	Number of bytes to read.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>NextOffsetp</i>	The offset into the directory for the following entry.
<i>OutDirStatusp</i>	Status information for the directory pointed to by <i>DirFidp</i> .
<i>OutTokenp</i>	Token granted against <i>DirFidp</i> .
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.
<i>dirStream</i>	The array of bytes making up the external representation of this part of the directory.

## 2.13.17 AFS\_Lookup

```

error_status_t AFS_Lookup
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *DirFidp,
 [in] afsTaggedName    *Namep,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsFid          *OutFidp,
 [out] afsFetchStatus  *OutFidStatusp,
 [out] afsFetchStatus  *OutDirStatusp,
 [out] afsToken        *OutTokenp,
 [out] afsVolSync      *Syncp
);

```

AFS\_Lookup parameters are:

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory from which to obtain the directory information.
<i>Namep</i>	The character string name of the file for which information is being requested.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier of the requested file.
<i>OutFidStatusp</i>	Status information for the specified file.
<i>OutDirStatusp</i>	Status information of the specified directory.
<i>OutTokenp</i>	The token against the directory (allowing the directory entry to be cached).
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.



## 2.13.18 AFS\_GetToken

```

error_status_t AFS_GetToken
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *Fidp,
 [in] afsToken         *MinTokenp,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsToken        *OutTokenp,
 [out] afsRecordLock   *OutBlockerp,
 [out] afsFetchStatus  *OutStatusp,
 [out] afsVolSync      *Syncp
);

```

**AFS\_GetToken** parameters are:

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	File identifier of the file to obtain a token against.
<i>MinTokenp</i>	Specification of the minimum requested token.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutTokenp</i>	Actual token granted.
<i>OutBlockerp</i>	Information about the possessor of the token that prevents the granting of the requested token (valid only for lock-family tokens).
<i>OutStatusp</i>	Status information on the file specified by <i>Fidp</i> .
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**2.13.19 AFS\_ReleaseTokens**

```
error_status_t AFS_ReleaseTokens
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsReturns       *Tokens_Arrayp,
 [in] unsigned32       Flags
);
```

**AFS\_ReleaseTokens** parameters are:

*h*                    The RPC binding handle.

*Tokens\_Arrayp*    Tokens granted to this file server which are to be relinquished.

*Flags*              Any of the AFS\_FLAG\_\* flags defined in the <afs4int.h> file.

**2.13.20 AFS\_GetTime**

```
error_status_t AFS_GetTime
(/* provider_version(1) */
 [in] handle_t          h,
 [out] unsigned32      *Secondsp,
 [out] unsigned32      *USecondsp,
 [out] unsigned32      *SyncDistance,
 [out] unsigned32      *SyncDispersion
);
```

**AFS\_GetTime** parameters are:

<i>h</i>	The RPC binding handle.
<i>Secondsp</i>	Number of seconds since January 1, 1970 UTC.
<i>USecondsp</i>	Number of microseconds into the current segment.
<i>SyncDistance</i>	Estimated path length to source of reliable time.
<i>SyncDispersion</i>	Measure of SyncDistance's variance.

## 2.13.21 AFS\_MakeMountPoint

```

error_status_t AFS_MakeMountPoint
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *DirFidp,
 [in] afsTaggedName    *Namep,
 [in] afsTaggedName    *cellNamep,
 [in] afsFStype        Type,
 [in] afsTaggedName    *volumeNamep,
 [in] afsStoreStatus   *InStatusp,
 [in] afsHyper         *minVVp,
 [in] unsigned32       Flags,
 [out] afsFid          *OutFidp,
 [out] afsFetchStatus  *OutFidStatusp,
 [out] afsFetchStatus  *OutDirStatusp,
 [out] afsVolSync      *Syncp
);

```

AFS\_MakeMountPoint parameters are:

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the new mount point is to be created.
<i>Namep</i>	The name of the mount point.
<i>cellNamep</i>	The string name of the cell in which the mount point is being created.
<i>Type</i>	The type of the file system mount point.
<i>volumeNamep</i>	The name of the fileset to be mounted on the newly created mount point.
<i>InStatusp</i>	This specifies the clientModTime field and unixModeBits of the new mount point.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier for the newly created mount point.
<i>OutFidStatusp</i>	The status information of the newly created mount point upon termination of the call.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the mount point was created) upon termination of the call.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**2.13.22 AFS\_GetStatistics**

```
error_status_t AFS_GetStatistics
(/* provider_version(1) */
 [in] handle_t          h,
 [out] afsStatistics    *Statisticsp
);
```

**AFS\_GetStatistics** parameters are:

*h*                    The RPC binding handle.  
*Statisticsp*        File server statistics.

## 2.13.23 AFS\_BulkFetchVV

```

error_status_t AFS_BulkFetchVV
(* provider_version(1) */
[in] handle_t          h,
[in] afsHyper         *cellIdp,
[in] afsBulkVolIDs    *VolIDsp,
[in] unsigned32       NumVols,
[in] unsigned32       Flags,
[in] unsigned32       spare1,
[in] unsigned32       spare2,
[out] afsBulkVVs      *VolVVsp,
[out] unsigned32      *spare4
);

```

AFS\_BulkFetchVV parameters are:

<i>h</i>	The RPC binding handle.
<i>cellIdp</i>	The cell identifier for the filesets whose Volume Version numbers are desired.
<i>VolIDsp</i>	The fileset identifiers (within <i>cellIdp</i> ) whose Volume Version numbers are desired.
<i>NumVols</i>	(Redundant) number of identifiers in <i>VolIDsp</i> .
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>spare1, spare2</i>	Spare IN parameters.
<i>VolVVsp</i>	The Volume Version numbers for the filesets whose identifiers were passed in.
<i>spare4</i>	Spare OUT parameter.

## 2.13.24 AFS\_BulkKeepAlive

```

error_status_t AFS_BulkKeepAlive
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsBulkFEX       *KAFEXp,
 [in] unsigned32       numExecFids,
 [in] unsigned32       Flags,
 [in] unsigned32       spare1,
 [in] unsigned32       spare2,
 [out] unsigned32      *spare4
);

```

AFS\_BulkKeepAlive parameters are:

<i>h</i>	The RPC binding handle.
<i>KAFEXp</i>	Collection of afsFid structures.
<i>numExecFids</i>	(Redundant) the number of afsFid structures in <i>KAFEXp</i> .
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>spare1, spare2</i>	Spare IN parameters.
<i>spare4</i>	Spare OUT parameter.

## 2.13.25 AFS\_ProcessQuota

```

error_status_t AFS_ProcessQuota
(/* provider_version(1) */
 [in] handle_t           h,
 [in] afsFid            *Fidp,
 [in] afsHyper          *minVVp,
 [in] unsigned32        Flags,
 [in,out]afsQuota       *quotaListp,
 [out] afsFetchStatus    *OutStatusp,
 [out] afsVolSync       *Syncp
);

```

AFS\_ProcessQuota parameters are:

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	File ID of file or directory whose quota is being changed.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>quotaLisp</i>	The quota information being provided or obtained.
<i>OutStatusp</i>	The updated file status after the operation.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.



**2.13.26 AFS\_GetServerInterfaces**

```
error_status_t AFS_GetServerInterfaces
(/* provider_version(1) */
 [in] handle_t          h,
 [in,out]dfs_interfaceList  *serverInterfacesP
);
```

**AFS\_GetServerInterfaces** parameters are:

*h*                    The RPC binding handle.

*serverInterfacesP* Will contain the interface information when the call returns. Currently there is only one interface defined.

**2.13.27 AFS\_SetParams**

```
error_status_t AFS_SetParams
(/* provider_version(1) */
 [in] handle_t          h,
 [in] unsigned32       Flags,
 [in,out]afsConnParams *paramsP
);
```

**AFS\_SetParams** parameters are:

<i>h</i>	The RPC binding handle.
<i>Flags</i>	Either AFS_PARAM_TSR_COMPLETE or AFS_PARAM_RESET_CONN.
<i>paramsP</i>	A block of connection-oriented parameters.

**2.14 The File Exporter Interface End**  
}



## Cache Manager IDL Declarations

This chapter specifies the RPC interfaces used by the AFS (4.0) Token Manager Interface. It contains the declaration of all structures and remote procedure calls (rpc) calls required for the AFS (4.0) Token Manager Interface.

The Token Manager Interface consists of a set of structures, mask values and flags contained in an IDL file with a provider version number of 1 to indicate that it is an original interface. In addition, each RPC defined in this structure also has a provider version number of 1, indicating that it also is an original RPC. See the topic, *DFS RPC Versioning Scheme*, Chapter 5 on page 69, for information on versioning which will explain in more detail about provider version numbers.

### 3.1 The Token Manager Interface

The Token Manager Interface is contained within the interface, **TKN4Int**, which starts out as shown below and contains all the items shown in this chapter exclusive of the definitions of the parameters in the RPC calls which are included from Chapter 7 on page 145 and the common data definitions in the configuration file **common\_data.idl** which are contained in Chapter 4 on page 57 for clarity. There are nine remote procedure calls defined that the client and cache manager can use in this interface which is defined as provider version 1. Chapter 5 on page 69 provides information for adding new interfaces, or new remote procedure calls to this interface.

```
[
uuid(4d37f2dd-ed96-0000-02c0-37cf1e000000),
version(4.0)
/* provider_version(1) */
]

interface TKN4Int
{

import "dcedfs/common_data.idl";

.
.
.
/* } at end of this chapter concludes this interface. */
```

## 3.2 Definitions for Flags Field

These are the definitions for the flags field used in this interface.

```
const unsigned32 TKN_FLAG_BACK_UP = 0x001;
    /* this is an unannounced TKN_InitTokenState call from a
       recovering server */

const unsigned32 TKN_FLAG_CRASHED = 0x002;
    /* the file exporter is currently in a post-crash
       token-recovery phase */

const unsigned32 TKN_FLAG_DISALLOW_SAME = 0x004;
    /* this recovery cannot accept AFS_FLAG_TOKENID AFS_GetToken calls */

const unsigned32 TKN_FLAG_NEW_IF = 0x008;
    /* this call is made to set the timer values; the tokens are OK */
```

## 3.3 Definitions of Token Manager Exported Operations

**Note:** The Opcode of each operation is implicitly assigned based on the order of where the operation is placed. The recommendation is to always place the new operation at the end of this list.

### 3.3.1 TKN\_Probe

```
error_status_t TKN_Probe /* May be defined as "idempotent" */
    /* provider_version(1) */
    [in] handle_t h
);
```

TKN\_Probe parameters are:

*h*                      The RPC binding handle.

### 3.3.2 TKN\_InitTokenState

```
error_status_t TKN_InitTokenState
    /* provider_version(1) */
    [in] handle_t                      h,
    [in] unsigned32                    Flags,
    [in] unsigned32                    hostLifeGuarantee,
    [in] unsigned32                    hostRPCGuarantee,
    [in] unsigned32                    deadServerTimeout,
    [in] unsigned32                    serverRestartEpoch,
    [in] unsigned32                    spare1,
    [in] unsigned32                    spare2,
    [in] unsigned32                    spare3,
    [out] unsigned32                   *spare4,
    [out] unsigned32                   *spare5,
    [out] unsigned32                   *spare6
);
```

TKN\_InitTokenState parameters are:

*h*                      The RPC binding handle.

<i>Flags</i>	Flags (TKN_FLAG*) indicating the server's likely token state.
<i>hostLifeGuarantee</i>	Duration of the host lifetime on the server.
<i>hostRPCGuarantee</i>	Duration of the promise of RPC attempts from the server.
<i>deadServerTimeout</i>	Suggested maximum interval between probes from Cache Managers to unresponsive servers.
<i>serverRestartEpoch</i>	Numeric tag for this run of the File Server, to allow Cache Managers to distinguish network partitions from server crashes.
<i>spare1 .. spare6</i>	Spares.

### 3.3.3 TKN\_TokenRevoke

```
error_status_t TKN_TokenRevoke
(/* provider_version(1) */
 [in] handle_t          h,
 [in,out]afsRevokes    *revokeDescp
);
```

TKN\_TokenRevoke parameters are:

<i>h</i>	The RPC binding handle.
<i>revokeDescp</i>	The tokens which are being revoked.

### 3.3.4 TKN\_GetCellName

```
error_status_t TKN_GetCellName
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsHyper          *Cellp,
 [out] NameString_t     CellNamep
);
```

TKN\_GetCellName parameters are:

<i>h</i>	The RPC binding handle.
<i>Cellp</i>	The ID for a cell.
<i>CellNamep</i>	The string name corresponding to <b>cellp</b> .

### 3.3.5 TKN\_GetLock

```
error_status_t TKN_GetLock
(/* provider_version(1) */
 [in] handle_t          h,
 [in] unsigned32       index,
 [out] afsDBLock       *lock
);
```

TKN\_GetLock parameters are:

<i>h</i>	The RPC binding handle.
----------	-------------------------

*index*                   Entry to return.

*lock*                    Lock information returned, including name, wait states, exclusive locks, readers reading and number waiting on it.

### 3.3.6 TKN\_GetCE

```
error_status_t TKN_GetCE
(* provider_version(1) */
[in] handle_t                    h,
[in] unsigned32                 index,
[out] afsDBCACHEEntry         *ce
);
```

TKN\_GetCE parameters are:

*h*                        The RPC binding handle.

*index*                   Entry to return.

                          [out] afsDBCACHEEntry         \*ce

*ce*                      The cache entry with token information, number of opens, readers, writers, shared, exclusives; locks, , expiration, mask, reference count, network part of the fid, and so on.

### 3.3.7 TKN\_GetServerInterfaces

```
error_status_t TKN_GetServerInterfaces
(* provider_version(1) */
[in] handle_t                    h,
[in,out]dfs_interfaceList       *serverInterfacesP
);
```

TKN\_GetServerInterfaces parameters are:

*h*                        The RPC binding handle.

*serverInterfacesP*   List of supported interfaces.

### 3.3.8 TKN\_SetParams

```
error_status_t TKN_SetParams
(* provider_version(1) */
[in] handle_t                    h,
[in] unsigned32                 Flags,
[in,out]afsConnParams         *paramsP
);
```

TKN\_SetParams parameters are:

*h*                        The RPC binding handle.

*Flags*                   For future expansion.

*paramsP*                A block of connection-oriented parameters.



### 3.3.9 TKN\_AsyncGrant

```
error_status_t TKN_AsyncGrant
(/* provider_version(1) */
 [in] handle_t          h,
 [in] afsFid           *grantedFileIDp,
 [in] afsToken         *grantedTokenP,
 [in] long              grantedRequestID
);
```

**TKN\_AsyncGrant** parameters are:

*h* The RPC binding handle.

*grantedFileIDp* File identifier of the file or directory to which the granted token applies.

*grantedTokenP* The token that was granted.

*grantedRequestID* Identifier of the original async grant request.

## 3.4 The Token Manager Interface End

```
}
```



## Common IDL Data

This chapter specifies the common data structures required for the AFS (4.0) Client and Server. The definitions here are necessary for both the File Exporter Interface and the Token Manager Interface and are imported into those interfaces.

### 4.1 Interface Common Data

```
interface common_data

{
. /* All the information in this chapter is part of the
. * common_data interface, from here to the ending brace
. * following the heading, "Interface Common Data end" */
/* } at the end of this chapter concludes the interface common data. */
```

### 4.2 General AFS Constants

```
const unsigned32 AFS_UNDEF_AFSID = -1;
/* Generic undefined AFS id */
const unsigned32 AFS_NAMEMAX = 256;
const unsigned32 AFS_PATHMAX = 1024;
const unsigned32 AFS_MAXHOSTS = 8;
const unsigned32 AFS_TOKENDEAD = 1235;
const unsigned32 AFS_ACLMAX = 8188;
/* 8k - sizeof(long), so afsACL in 8k; was 1024 */
const unsigned32 AFS_BULKMAX = 32;
const unsigned32 AFS_DISKNAME_SIZE = 32;
const unsigned32 AFS_NAMEMAXLEN = 112;
/* Max len for fileset names */
```

### 4.3 Constants for Cell and Hosts

```
const unsigned32 MAXCELLCHARS = 128;
const unsigned32 MAXHOSTCHARS = 128;
const unsigned32 MAXHOSTSPERCELL = 64;
```

#### 4.4 AFS Object Types Used by AFS\_Mount

```
const unsigned32 Invalid = 0;
const unsigned32 File = 1;
const unsigned32 Directory = 2;
const unsigned32 SymbolicLink = 3;
const unsigned32 MountPoint = 4;
const unsigned32 FIFO = 5;
const unsigned32 Socket = 6;
const unsigned32 BlockDev = 7;
const unsigned32 CharDev = 8;
const unsigned32 FETCHSTATUS_VERSION = 2;
```

#### 4.5 Quota Types for Quota Setting Commands

```
const unsigned32 AFS_FILESYS_BSD = 1;
const unsigned32 AFS_FILESYS_EPISODE = 2;
```

#### 4.6 Quota Opcodes for Quota Setting Commands

```
const unsigned32 AFS_QUOTA_GET = 1;
const unsigned32 AFS_QUOTA_SET = 2;
```

#### 4.7 Physical File System Types (for afsFStype)

```
const unsigned32 AFS_FS = 0;
const unsigned32 DEC_FS = 1;
const unsigned32 UFS_FS = 2;
const unsigned32 AIX_FS = 4;
```

#### 4.8 Volume Types Used for afsVolumeType

```
const unsigned32 ReadOnly = 0;
const unsigned32 ReadWrite = 1;
```

## 4.9 Values for the *afsRevokeDesc* Flags

```
const unsigned32 AFS_REVOKE_COL_A_VALID = 0x1;
const unsigned32 AFS_REVOKE_COL_B_VALID = 0x2;
const unsigned32 AFS_REVOKE_FORCED_REVOCATION = 0x4;
const unsigned32 AFS_REVOKE_DUE_TO_GC = 0x8;
```

## 4.10 Values Used in *afsRevokeDesc*'s *outFlags*

```
const unsigned32 AFS_REVOKE_LOCKDATA_VALID = 0x1;
```

## 4.11 Data Types

```
typedef uuid_t    afsUUID;
typedef long      afsFStype;
typedef long      afsVolumeType;

typedef pipe      byte    pipe_t;

typedef unsigned32 codesetTag;
const unsigned32 AFS_TAG_ORIGASCII = 0;
```

### 4.11.1 General Definitions for AFS Data Structures

```
typedef [string] byte    NameString_t[AFS_NAMEMAX];
typedef [string] byte    PathString_t[AFS_PATHMAX];
typedef byte    afsDiskName[AFS_DISKNAME_SIZE];
```

### 4.11.2 General Tagged-name for a Pathname Component

Use of the *tn\_tag* value *AFS\_TAG\_ORIGASCII* defined above in Section 4.11 means that the *tn\_length* field doesn't need to be filled in.

**Note:** *tn\_chars* is assumed to be null-terminated.

### 4.11.3 Define the *afsTaggedName* Structure

```
typedef struct afsTaggedName {
    codesetTag    tn_tag;
    unsigned16    tn_length;
    byte          tn_chars[AFS_NAMEMAX+1];
} afsTaggedName;
```

**4.11.4 Define the afsTaggedPath Structure**

```
typedef struct afsTaggedPath {
    codesetTag    tp_tag;
    unsigned16    tp_length;
    byte          tp_chars[AFS_PATHMAX+1];
} afsTaggedPath;
```

**4.11.5 Define the afsNetAddr Structure**

```
typedef struct afsNetAddr {
    unsigned16    type;
    unsigned8     data[14];
} afsNetAddr;
```

**4.11.6 Define the afsTimeval Structure**

```
typedef struct afsTimeval {
    unsigned32    sec;
    unsigned32    usec;
} afsTimeval;
```

**4.11.7 Define the afsHyper Structure**

**Note:** This is a 64-bit structure used by objects such as cells, volumes, and so on.

```
typedef struct afsHyper {
    unsigned32    high;
    unsigned32    low;
} afsHyper;

/*
 * Generic operations on afsHypers are defined in dcedfs/stds.h.
 */
```

**4.11.8 Define the afsFid Structure**

```
typedef struct afsFid {
    afsHyper    Cell;
    afsHyper    Volume;
    unsigned32  Vnode;
    unsigned32  Unique;
} afsFid;
```

**4.11.9 Define the afsFidTaggedName Structure**

```
typedef struct afsFidTaggedName {
    afsFid    fid;
    afsTaggedName    name;
} afsFidTaggedName;
```

**4.11.10 Define the afsToken Structure**

```
typedef struct afsToken {
    afsHyper tokenID;
    unsigned32 expirationTime;
    afsHyper type;
    unsigned32 beginRange;
    unsigned32 endRange;
    unsigned32 beginRangeExt;
    unsigned32 endRangeExt;
} afsToken;
```

**4.11.11 Define the afsRecordLock Structure**

**Note:** This is a Sys V Record Lock to return blocker's info.

```
typedef struct afsRecordLock {
    signed16    l_type;
    signed16    l_whence;
    unsigned32  l_start_pos;
    unsigned32  l_end_pos;
    unsigned32  l_pid;
    unsigned32  l_sysid;
    unsigned32  l_fstype;
    unsigned32  l_start_pos_ext;
    unsigned32  l_end_pos_ext;
} afsRecordLock;
```

**4.11.12 Define the afsRevokeDesc Structure**

**Note:** This structure contains information about the token the Cache Manager (cm) already holds.

```
typedef struct afsRevokeDesc {

    afsFid fid;           /* useful hint */
    afsHyper tokenID;
    afsHyper type;       /* mask */
    unsigned32 flags;    /* input flags to RPC*/
    unsigned32 outFlags; /* output flags from RPC */
    afsHyper errorIDs;   /* (mask[dude] == 1) <==> dude not revoked */
    /*
     * Info about the tokens to be offered for swapping
     * The first one
     */
    afsToken columnA;    /* the offer */
    afsHyper colAChoice; /* the accepted token types */
    /*
     * the other token to be offered
     */
    afsToken columnB;    /* the offer */
    afsHyper colBChoice; /* the accepted token types */
    afsRecordLock recordLock; /* the blocker's info */
} afsRevokeDesc;
```

**4.11.13 Define the afsReturnDesc Structure**

**Note:** This structure contains information for the return of tokens.

```
typedef struct afsReturnDesc {
    afsFid fid;           /* useful hint */
    afsHyper tokenID;
    afsHyper type;       /* mask */
    unsigned32 flags;    /* just in case */
} afsReturnDesc;
```

**4.11.14 Define the afsConnParams Structure**

**Note:** This is the method of adjusting connection parameters:

1. Low-order bit of Mask (0x1) declares that Values[0] contains information.
2. Next bit in Mask (0x2) declares that Values[1] contains information.

```
typedef struct afsConnParams {
    unsigned32 Mask;
    unsigned32 Values[20];
} afsConnParams;
const unsigned32 AFS_CONN_PARAM_HOSTLIFE = 0;
const unsigned32 AFS_CONN_PARAM_HOSTRPC = 1;
const unsigned32 AFS_CONN_PARAM_DEADSERVER = 2;
const unsigned32 AFS_CONN_PARAM_EPOCH = 3;
```

**4.11.15 Define the afsDBLockDesc Structure**

**Note:** This is a temporary structure for debugging purposes.

```
typedef struct afsDBLockDesc {
    unsigned8 waitStates;
    unsigned8 exclLocked;
    unsigned8 readersReading;
    unsigned8 numWaiting;
} afsDBLockDesc;
```

**4.11.16 Define the afsDBCACHEEntry Structure**

**Note:** This is a temporary structure for debugging purposes.

```
typedef struct afsDBCACHEEntry {
    unsigned32 addr;
    afsFid fid;           /* network part of the fid */
    afsHyper length;
    afsHyper dataVersion;
    afsDBLockDesc lock;
    unsigned32 tokenBaseID;
    unsigned32 tokenMask;
    unsigned32 tokenExpiration;
    signed16 refCount;
    signed16 opens;
    signed16 writers;
    signed16 readers;
```



```

    signed16 shareds;
    signed16 exclusives;
    unsigned8 mvstat;
    unsigned8 states;
} afsDBCacheEntry;

```

#### 4.11.17 Define the afsDBLock Structure

**Note:** This is a temporary structure for debugging purposes.

```

typedef struct afsDBLock {
    byte name[16];
    afsDBLockDesc lock;
} afsDBLock;

```

## 4.12 Various Bulk Typedefs From Primitive Structures

These are from the primitive structures above.

### 4.12.1 Define the afsRevokes Structure

```

typedef struct afsRevokes { /* needed also by tkn4int.idl */
    long afsRevokes_len;
    [length_is(afsRevokes_len)] afsRevokeDesc afsRevokes_val[AFS_BULKMAX];
} afsRevokes;

```

### 4.12.2 Define the afsReturns Structure

```

typedef struct afsReturns { /* needed also by tkn4int.idl */
    long afsReturns_len;
    [length_is(afsReturns_len)] afsReturnDesc afsReturns_val[AFS_BULKMAX];
} afsReturns;

```

### 4.12.3 Define the afsFids Structure

```

typedef struct afsFids {
    unsigned32 afsFids_len;
    [length_is(afsFids_len)] afsFid afsFids_val[AFS_BULKMAX];
} afsFids;

```

### 4.12.4 Define the afsTokens Structure

```

typedef struct afsTokens {
    unsigned32 afsTokens_len;
    [length_is(afsTokens_len)] afsToken afsTokens[AFS_BULKMAX];
} afsTokens;

```

### 4.12.5 Define the afsStrings Structure

**Note:** afsStrings, used in AFS\_BulkLookup, used to be identical to NameString\_t itself.

```
typedef struct afsStrings {
    unsigned32 afsStrings_len;
    [length_is(afsStrings_len)] afsTaggedName afsStrings_val[AFS_BULKMAX];
} afsStrings;
```

## 4.13 Data Types for DFS RPC Versioning Scheme

### 4.13.1 Constants for RPC Versioning Scheme

```
const unsigned32 MAXSPARETEXT = 50;
const unsigned32 MAXINTERFACESPERVERSION = 10;
```

### 4.13.2 dfs\_interfaceDescription Structure

```
typedef struct dfs_interfaceDescription {

    uuid_t      interface_uuid; /* i/f UUID of the supported i/f */
    unsigned16  vers_major;     /* i/f major version number */
    unsigned16  vers_minor;     /* i/f minor version number */
    unsigned32  vers_provider;  /* provider version number */

    unsigned32  spare0;         /* some long spares */
    unsigned32  spare1;
    unsigned32  spare2;
    unsigned32  spare3;
    unsigned32  spare4;
    unsigned32  spare5;
    unsigned32  spare6;
    unsigned32  spare7;
    unsigned32  spare8;
    unsigned32  spare9;

    byte spareText[MAXSPARETEXT]; /* spare text */

} dfs_interfaceDescription;
```

### 4.13.3 dfs\_interfaceList Structure

```
typedef struct dfs_interfaceList {

    unsigned32  dfs_interfaceList_len;
    [length_is(dfs_interfaceList_len)] dfs_interfaceDescription
        dfs_interfaceList_val[MAXINTERFACESPERVERSION];

} dfs_interfaceList;
```

#### **4.14 Interface Common Data End**

}



# ***X/Open Preliminary Specification***

## **Part 3:**

### **DFS Versioning Scheme**



## DFS RPC Versioning Scheme

This is the formal specification document for the DFS client/server versioning scheme. This mechanism should be supported by all DFS clients and servers. This specification is organized into six sections.

Data Structures	Defined in Chapter 4 on page 57 with specific definitions for versioning found in Section 4.13 on page 64, and demonstrated in Section 5.1.
Versioning API	Defined in Section 5.2 on page 70.
Example IDL file	Defined in Section 5.3 on page 73.
Example Client Application	Defined in Section 5.4 on page 76.
Example Server Application	Defined in Section 5.5 on page 83.
Example Manager Application	Defined in Section 5.6 on page 87.
Extending DFS interfaces	Defined in Section 5.7 on page 95.

### 5.1 Data Structures

This example structure defines the interface Description for each interface that is supported by the server. The server maintains an array of such structures for each interface that it exports. Also, upon startup the server fills these arrays with the interfaces that it registers with the runtime.

The actual structure is defined in Section 4.13.2 on page 64 and is repeated here so that this chapter is self-explanatory.

#### Example 5-1 Interface Description for an Interface

```
const long MAXSPARETEXT = 50;
const long MAXINTERFACESPERVERSION = 10;

typedef struct dfs_interfaceDescription {

    uuid_t      interface_uuid; /* i/f UUID of the supported i/f */
    unsigned16  vers_major;     /* i/f major version number */
    unsigned16  vers_minor;     /* i/f minor version number */
    unsigned32  vers_provider;  /* provider version number */

    unsigned32  spare0;         /* some long spares */
    unsigned32  spare1;
    unsigned32  spare2;
    unsigned32  spare3;
    unsigned32  spare4;
    unsigned32  spare5;
    unsigned32  spare6;
    unsigned32  spare7;
    unsigned32  spare8;
```

```

unsigned32 spare9;

byte spareText[MAXSPARETEXT];/* spare text */
} dfs_interfaceDescription;

```

The following example structure shows a list of interface descriptions for a particular interface. A Client requests a list of different interface versions for a particular RPC interface using the `RPC_GetServerInterfaces()` remote procedure call.

The actual structure is defined in Section 4.13.3 on page 64 and is repeated here so that this chapter is self-explanatory.

**Example 5-2** List of Interface Descriptions for an Interface

```

typedef struct dfs_interfaceList {

    unsigned32 dfs_interfaceList_len;
    [length_is(dfs_interfaceList_len)] dfs_interfaceDescription \
        dfs_interfaceList_val[MAXINTERFACESPERVERSION];
} dfs_interfaceList;

```

This example list (defined in Section 4.13.3 on page 64) shows the structure that servers fill with the interface(s) they support. They return a set of such (interface description) structures (defined in Section 4.13.2 on page 64) back to the client which is requesting it. This list is used by the versioning code described in Section 5.2 to find the best match between the facilities supported by a server and the facilities desired by a client.

## 5.2 DFS Versioning API

The versioning API for DFS makes use of some definitions whose syntax is unique to DFS. They are *IMPORT* and *\_TAKES*. The syntax for *IMPORT* is:

```
#define IMPORT extern
```

The syntax for *\_TAKES* is shown in its usage via the following example:

**Example 5-3** Usage of *IMPORT* definition with *\_TAKES* notation

```

IMPORT char * getenv _TAKES((
    char * varName
));

```

As can be seen, *IMPORT* is defined as *extern*. The *getenv* function is a utility function defined in the header `<stdlib.h>`, that is used to declare functions for number conversion, storage allocation and similar items. *\_TAKES* is used simply as notation showing what the shown *IMPORT* definition *takes* as a construct.

**Example 5-4** Construction of the *getenv* utility

*getenv* is an implementation-dependent function. It has the following format:

```
char * getenv(const char *name)
```



*getenv* returns the environment string associated with name, or NULL if no string exists.

### 5.2.1 Register an Interface with Versioning Mechanism

This function is used to register an interface with the versioning mechanism for DFS. It is called by server applications only.

```
IMPORT void dfs_installInterfaceDescription _TAKES((
    rpc_if_handle_t if_spec,
    rpc_if_handle_t orig_spec,
    unsigned32 vers_provider,
    unsigned_char_t *text,
    error_status_t *code
));
```

**dfs\_installInterfaceDescription** parameters are:

<i>if_spec</i>	Interface specification that is to be registered.
<i>orig_spec</i>	Interface specification that is extended (that is, Transarc's interface).
<i>vers_provider</i>	Provider version number of this interface.
<i>text</i>	Text describing the new interface (maximum of 50 characters).
<i>code</i>	Status that is returned from this call.

### 5.2.2 Printing a Returned List of Interfaces

This function is useful only for debugging purposes. It prints in human-readable form, the list of interfaces returned by the server. It is called by client/server applications.

```
IMPORT void dfs_printInterfaceDescription _TAKES ((
    dfs_interfaceList *interfaces,
    error_status_t *code
));
```

**dfs\_printInterfaceDescription** parameters are:

<i>interfaces</i>	A list of interfaces returned by the server. See Section 4.13.2 on page 64 and Section 4.13.3 on page 64 for more details.
<i>code</i>	Status that is returned from this call.

### 5.2.3 Function Called at Startup to Identify the Server Interface

This function is called by the client at startup to identify the server interface that the client is to use for RPCs. It is called by Client applications only.

```
IMPORT void dfs_GetServerInterfaces _TAKES((
    rpc_if_handle_t if_spec,
    dfs_interfaceList *interfaces,
    error_status_t *code
));
```

**dfs\_GetServerInterfaces** parameters are:

<i>if_spec</i>	Interface specification that is to be examined.
<i>interfaces</i>	A list of server supported interfaces returned by the server. See Section 4.13.2 on page 64 and Section 4.13.3 on page 64 for more details.

*code*                    Status that is returned from this call.

#### 5.2.4 Function to Compare Interfaces

This function is called by the client to compare the interface spec that it desires to get access to and the interfaces in the list of interface descriptions returned by the `RPC_GetServerInterfaces` call. It returns:

**MATCH\_GOOD**    UUID and major version numbers of the client and server interfaces are identical. (Refer to Example 5-1 on page 69 for more information on version numbers, minor and major, and UUIDs, if desired.)

This means that either:

1. Client and server interfaces are identical - the minor version numbers and the provider version numbers are equal.
2. Server minor version is greater than client's.
3. Server provider version is greater than client's.

**MATCH\_MEDIUM**

UUID and major version numbers of client and server are identical, and the minor version of client is less than or equal to server's. In addition the provider version number of client is greater than server's.

**MATCH\_BAD**    The UUID or the major version number of the client is different from the server's. Client-Server binding in this case is impossible and hence this case is least desirable.

```
IMPORT long dfs_sameInterface _TAKES ((
    rpc_if_handle_t if1,
    unsigned int if1_provider,
    dfs_interfaceList *serverInterfaces
));
```

**dfs\_sameInterface** parameters are:

*if1*                    Interface specification that is to be used for a comparison.

*if1\_provider*        Provider version number.

*serverInterfaces*    A list of server supported interfaces returned by the server. See Section 4.13.2 on page 64 and Section 4.13.3 on page 64 for more details.

### 5.3 Example IDL File

Provider numbers for the interface and for each RPC are added as a comment. The provider version number for the interface (in the IDL header) signifies a new version of the interface. The provider version number associated with each RPC signifies the provider version number of the interface at which time this RPC was added to the interface. In this example, the IDL file is the original interface and so, the provider number for the interface and for the RPC are both 1.

This example is taken from the UPDATE Server (AFS 4.0) Interface Definition which is owned by Transarc Corporation. The portion necessary in order to show how versioning works is the exported procedure for fetching the interfaces supported by the upserver, namely, `UPDATE_GetServerInterfaces()`. However, for completeness, other procedures are also shown since they are included in the example manager application in Section 5.6 on page 87.

#### 5.3.1 The UPDATE Server Interface

```
[
uuid(4d37f2dd-ed43-0000-02c0-37cf1e001000),
version(4.0)
/* provider_version(1) */
]

interface UPDATE
{
/*
 * Include the AFS basic Header files
 */

/**
 **  import "cellconfig.idl" ;
 **/
import "dcedfs/common_data.idl";
    ...
    ...

/* } at end of this section ends this interface. */
```

#### 5.3.2 Constants for the UPDATE Interface

```
const    long    TIMEOUT = 300;
const    long    MAX_PROTSEQ_LEN = 32;
const    long    BINDING_LEN = 256;
const    long    MAX_BINDING_LEN = 1024;
const    long    MAX_NAME_SIZE = 256;    /* include null byte at end */
```

### 5.3.3 Define the updateFileStatS Structure

```
typedef struct updateFileStatS
{
    afsHyper fileLength;
    unsigned32 mode;
    unsigned32 uid;
    unsigned32 gid;
    unsigned32 mtime;    /* modify time */
    unsigned32 atime;    /* access time */
    unsigned32 spare1;
    unsigned32 spare2;
    unsigned32 spare3;
    unsigned32 spare4;
} updateFileStatT;
```

### 5.3.4 UPDATE\_GetServerInterfaces

This exported procedure is called by the upclient to fetch the interfaces supported by the upserver. A bulk parameter of upto MAXINTERFACES number of interfaces is returned.

```
error_status_t UPDATE_GetServerInterfaces
(/* provider_version(1) */
 [in] handle_t          h,
 [in,out]dfs_interfaceList *serverInterfaces
);
```

UPDATE\_GetServerInterfaces parameters are:

*h*                    The RPC binding handle.

*serverInterfaces*   A list of server supported interfaces returned by the server. See Section 4.13.2 on page 64 and Section 4.13.3 on page 64 for more details.

### 5.3.5 UPDATE\_FetchInfo

This exported procedure is called by a client to fetch the information (that is, the status) of a specified directory in the remote host.

An AFS client calls this remote procedure periodically to bring object files at its site update-to-date.

```
error_status_t UPDATE_FetchInfo
(/* provider_version(1) */
 [in] handle_t          h, /* Necessary for explicit_handle */
 [in] NameString_t     DirName,
 [out] pipe_t          *Stream
);
```

UPDATE\_FetchInfo parameters are:

*h*                    The RPC binding handle.

*DirName*            The string name corresponding to the specified directory.

*Stream*             The character pipe parameter returning the data from the directory. (See **X/Open DCE: Remote Procedure Call** specification for implementation details.)

### 5.3.6 UPDATE\_FetchFile

This exported procedure is called by a client to fetch the whole data from a specified file in the remote host.

An AFS client calls this remote procedure periodically to bring object files at the client site update-to-date.

```
error_status_t UPDATE_FetchFile
(* provider_version(1) */
  [in]  handle_t          h, /* Necessary for explicit_handle */
  [in]  NameString_t     FileName,
  [out] pipe_t           *Stream
);
```

UPDATE\_FetchFile parameters are:

<i>h</i>	The RPC binding handle.
<i>FileName</i>	The string name corresponding to the specified file.
<i>Stream</i>	The character pipe parameter returning the data from the file. (See <b>X/Open DCE: Remote Procedure Call</b> specification for implementation details.)

### 5.3.7 UPDATE\_FetchObjectInfo

This exported procedure is called by a client to fetch the information (that is, the status) of a specified file in the remote host.

An AFS client calls this remote procedure periodically to bring object files at its site update-to-date.

```
error_status_t UPDATE_FetchObjectInfo
(* provider_version(1) */
  [in]  handle_t          h, /* Necessary for explicit_handle */
  [in]  NameString_t     objectName,
  [out] updateFileStatT  *fileStatP
);
```

UPDATE\_FetchObjectInfo parameters are:

<i>h</i>	The RPC binding handle.
<i>objectName</i>	The string name corresponding to the specified file object.
<i>fileStatP</i>	Returns the data from the specified object. See Section 5.3.3 on page 74 for details.

### 5.3.8 The UPDATE Server Interface End

```
}
```

## 5.4 Example Client Application

The client application, upon startup retrieves server supported interfaces and sets a global to identify the interface (stub) that the client would use to make RPCs. The `dfs_selectInterface()` function sets the global and each RPC wrapper first checks the global and calls the appropriate stub.

### 5.4.1 Typical Client Headers

```
#include <dcedfs/param.h>
#include <dcedfs/stds.h>
#include <dcedfs/osi.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <netdb.h>
#include <dirent.h>
#include <pthread.h>
#include <dce/rpc.h>
#include <dcedfs/compat.h>
#include <dcedfs/cmd.h>
#include <update.h>
#include <upcommon.h>
```

### 5.4.2 Headers for Serviceability

The following headers get included for serviceability. The last three are specific to this subcomponent. They would need to be defined for each subcomponent that has implemented serviceability, and the names would be chosen to be specific to the function. The first three are included as is regardless of what subcomponent they are in, as they are needed for general serviceability and are not specific to the subcomponent.

The explanation of serviceability is beyond the scope of this document and this example.

```
/* For serviceability */
#include <dce/dce.h>
#include <dce/dce_msg.h>
#include <dce/dcesvcmsg.h>
#include <dfsudtmac.h>
#include <dfsudtsvc.h>
#include <dfsudtmsg.h>
```

### 5.4.3 Constants for the Client Application

```
/* flag to check if UPDATE interface accessible */
unsigned int UPDATE_INTERFACE = 0;
unsigned int UPDATE_CHECK_PROVIDER_OR_RANGE = 0;

/* List of interfaces supported */
#define DFS_UPDATE 1
#define UPDATE_PROVIDER 1

/* should be NIDL_PIPE_BUFF_SIZE defined in stubbase.h */
#define PREFER_BUFFSIZE 2048
#define UPDATE_MAXBINDINGS 32
#define MAX_RPC_FAILURES 10
```



**5.4.6 ANSI C Declaration**

```
/* Forward declaration for ANSI C */
long isCompatibleObject();
```

**5.4.7 Client's import list**

```
/* filestr - client's list of items to import */

struct filestr {
    struct filestr *next;
    int type;
    char *name;
};
```

**5.4.8 Client Globals**

```
/* client globals */
upLogP logP = 0; /* upclient's handle on it's log file */

void pipeAlloc();
void ReceiveFile();
int AddToList();
int ZapList();
int IsCompatible();
int NotOnHost();
```

**5.4.9 Using dfs\_selectInterface to Select One**

This function checks to see if client has access to the desired interface. It can associate priorities with interfaces in the case of servers exporting multiple interfaces.

```
static error_status_t
dfs_selectInterface (h)
    rpc_binding_handle_t h;
{
    rpc_if_id_t if_id;
    static dfs_interfaceList serverInterfaces;
    unsigned int i, j, result;
    error_status_t code;

    serverInterfaces.dfs_interfaceList_len = 0;
    bzero((char *)serverInterfaces.dfs_interfaceList_val,
          MAXINTERFACESPERVERSION * sizeof(dfs_interfaceDescription));

    DFS_RPC_PREWRAP
    code = UPDATE_GetServerInterfaces (h, &serverInterfaces);
    DFS_RPC_POSTWRAP(code)

    if ( IS_COMM_ERR(code) )
        return code;

#ifdef DEBUG
```



```

dfs_printInterfaceDescription (&serverInterfaces, &code);
#endif

/* Check to see if you have access to desired interface. The
following are the different classes from most desired to
least.

Good:  UUID and major version numbers of the client and server
interfaces are identical.

      1. client and server interfaces are identical -- the minor
         version numbers and the provider version numbers are equal.
      2. Server minor version is greater than client's.
      3. Server provider version is greater than client's.

Medium: UUID and major version numbers of client and server are
identical, and the minor version of client is less than
or equal to server's. In addition the provider version
number of client is greater than server's.

Bad:   The UUID or the major version number of the client is
different from the server's. Client-Server binding in this
case is impossible and hence this case is least desirable. */

if (( result = dfs_sameInterface (UPDATE_v4_0_c_ifspec,
UPDATE_PROVIDER,
&serverInterfaces)) == MATCH_GOOD ||
    result == MATCH_MEDIUM) {

    /* check for new interfaces derived from the original interface
    should be in here, if multiple versions of a DFS interface is
    supported. Here we only have one version of the interface.
    */
    UPDATE_INTERFACE = DFS_UPDATE; /* original interface */

    if (result == MATCH_MEDIUM)
        UPDATE_CHECK_PROVIDER_OR_RANGE = 1;
    }
    if (!UPDATE_INTERFACE)
        return UP_BAD_INTERFACE;
    else return 0;
}

```

### 5.4.10 Example RPC Wrappers

The following are example RPC wrappers that a client can use.

#### The update\_Fetchfile Wrapper

```

error_status_t
update_FetchFile (h, filename, pipeP)
    rpc_binding_handle_t h;
    NameString_t filename;
    pipe_t *pipeP;
{
    error_status_t code;
    static unsigned long FUNCTION_PROVIDER = 1;
        /* provider version # this RPC is supported */
    static unsigned long FUNCTION_NOOP = 0;
        /* Is this function supported by server or not? */

    if (FUNCTION_NOOP) /* this RPC not supported */
        return (rpc_s_op_rng_error);

    /* Check to see if the desired interface provider is supported. If not
       we may try alternate interfaces. Also check if this RPC is supported
       by the server interface. This is done by comparing provider versions
       */

    if (UPDATE_INTERFACE == DFS_UPDATE &&
        UPDATE_PROVIDER >= FUNCTION_PROVIDER) {

#ifdef UPDATE_CANCEL_TEST
        Log(udt_s_call_upd_fetchfile);
#endif /* UPDATE_CANCEL_TEST */
        DFS_RPC_PREWRAP
        code = UPDATE_FetchFile (h, filename, pipeP);
        DFS_RPC_POSTWRAP(code)

#ifdef UPDATE_CANCEL_TEST
        Log(udt_s_cmpltd_upd_fetchfile);
#endif /* UPDATE_CANCEL_TEST */

        if ( IS_COMM_ERR(code) ) {
            if (UPDATE_CHECK_PROVIDER_OR_RANGE && code == rpc_s_op_rng_error){
                code = UP_BAD_INTERFACE;
                FUNCTION_NOOP = 1;
            }
        }
        return code;
    }
    else
        return(UP_BAD_INTERFACE);
}

```

**The update\_FetchInfo Wrapper**

```

error_status_t
update_FetchInfo (h, name, pipeP)
    rpc_binding_handle_t h;
    NameString_t name;
    pipe_t *pipeP;
{
    error_status_t code;
    static unsigned long FUNCTION_PROVIDER = 1;
        /* provider version # this RPC is supported */
    static unsigned long FUNCTION_NOOP = 0;
        /* Is this function supported by server or not? */

    if (FUNCTION_NOOP) /* this RPC not supported */
        return (rpc_s_op_rng_error);

    if (UPDATE_INTERFACE == DFS_UPDATE &&
        UPDATE_PROVIDER >= FUNCTION_PROVIDER) {

#ifdef UPDATE_CANCEL_TEST
        Log(udt_s_call_upd_fetchinfo);
#endif /* UPDATE_CANCEL_TEST */
        DFS_RPC_PREWRAP
        code = UPDATE_FetchInfo (h, name, pipeP);
        DFS_RPC_POSTWRAP(code)

#ifdef UPDATE_CANCEL_TEST
        Log(udt_s_cmpltd_upd_fetchinfo);
#endif /* UPDATE_CANCEL_TEST */

        if ( IS_COMM_ERR(code) ) {
            if (UPDATE_CHECK_PROVIDER_OR_RANGE && code == rpc_s_op_rng_error){
                code = UP_BAD_INTERFACE;
                FUNCTION_NOOP = 1;
            }
        }
        return code;
    }
    else
        return(UP_BAD_INTERFACE);
}

```

**The update\_FetchObjectInfo Wrapper**

```

error_status_t
update_FetchObjectInfo (h, name, buf)
    rpc_binding_handle_t h;
    char *name;
    updateFileStatT *buf;
{
    error_status_t code;
    static unsigned long FUNCTION_PROVIDER = 1;

```

```

        /* provider version # this RPC is supported */
static unsigned long FUNCTION_NOOP = 0;
        /* Is this function supported by server or not? */
NameString_t object;
        /* For use in UPDATE_FetchObjectInfo call */

if (FUNCTION_NOOP) /* this RPC not supported */
    return (rpc_s_op_rng_error);

if (UPDATE_INTERFACE == DFS_UPDATE &&
    UPDATE_PROVIDER >= FUNCTION_PROVIDER) {

    strcpy((char *)object, name);

#ifdef UPDATE_CANCEL_TEST
    Log(udt_s_call_upd_fobjinfo);
#endif /* UPDATE_CANCEL_TEST */

    DFS_RPC_PREWRAP
    code = UPDATE_FetchObjectInfo (h, object, buf);
    DFS_RPC_POSTWRAP(code)

#ifdef UPDATE_CANCEL_TEST
    Log(udt_s_cmpltd_upd_fobjinfo);
#endif /* UPDATE_CANCEL_TEST */

    if (IS_COMM_ERR(code)) {
        if (UPDATE_CHECK_PROVIDER_OR_RANGE && code == rpc_s_op_rng_error){
            code = UP_BAD_INTERFACE;
            FUNCTION_NOOP = 1;
        }
    }
    return code;
}
else
    return(UP_BAD_INTERFACE);
}

```

**Skeleton Client Main**

```

int main(argc,argv)
    int   argc;
    char  **argv;
{
    ...

    /* Initialize the serviceability function */

    initialize_svc();
    ...
}

```

## 5.5 Example Server Application

The server registers the interfaces with the runtime and with the versioning mechanism.

### 5.5.1 Typical Server Headers

```
#include <dcedfs/param.h>
#include <dcedfs/stds.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <sys/file.h>
#include <sys/time.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <netdb.h>
#include <stdio.h>

#include <dce/rpc.h>
#include <pthread.h>
#include <dcedfs/osi.h>
#include <dcedfs/compat.h>
#include <dcedfs/dfsauth.h>
#include <dcedfs/cmd.h>
#include <dirent.h>

#include <update.h>
#include <upcommon.h>
```

### 5.5.2 Headers for Serviceability

The following headers get included for serviceability. The last three are specific to this subcomponent. They would need to be defined for each file that has implemented serviceability, and the names would be chosen to be specific to the function. The first three are included as is regardless of what subcomponent they are in, as they are needed for general serviceability and are not specific to the subcomponent.

The explanation of serviceability is beyond the scope of this document and this example.

```
/* For serviceability */

#include <dce/dce.h>
#include <dce/dce_msg.h>
#include <dce/dcesvcmsg.h>
#include <dfsudtmac.h>
#include <dfsudtsvc.h>
#include <dfsudtmsg.h>
```

**5.5.3 Constants for the Server Application**

```
#define UPDATE_MAXCALLS 5
#define DEFAULT_ADMIN_FILENAME "admin.up"
```

**5.5.4 Define IS\_COMM\_ERR Function**

```
#define IS_COMM_ERR(s) (s == rpc_s_ok ? 0 : 1$
```

**5.5.5 Typical Serviceability Initialization**

```
char *whoami = "upserver";

dce_svc_handle_t    udt_svc_handle;

void initialize_svc()
{
    error_status_t st;

    udt_svc_handle = dce_svc_register(udt_svc_table, (idl_char *)"udt", &st);
    if (st != svc_s_ok)
        fprintf(stderr, "Cannot register svc 0x%lx0, st);
    dce_svc_routing((unsigned char *) "NOTICE:STDOUT:--", &st);
    if (st != svc_s_ok)
        fprintf(stderr, "Cannot route NOTICE messages 0x%lx0, st);
    dce_svc_routing((unsigned char *) "WARNING:STDOUT:--", &st);
    if (st != svc_s_ok)
        fprintf(stderr, "Cannot route WARNING messages 0x%lx0, st);
    dce_svc_routing((unsigned char *) "ERROR:STDERR:--", &st);
    if (st != svc_s_ok)
        fprintf(stderr, "Cannot route ERROR messages 0x%lx0, st);

    dfs_define_udt_message_tables();
}
```

**5.5.6 Skeleton Server Main**

```
int
main(argc, argv)
    int    argc;
    char  *argv[];
{

    extern char *dfs_dceErrTxt();
    extern UPDATE_v4_0_epv_t UPDATE_v4_0_manager_epv;

    osi_setlocale(LC_ALL, "");

    ...

    initialize_svc();

    initialize_upd_error_table();
```

```

...

us_argSetup();
code = cmd_Dispatch(argc, argv);
if ( code )
    exit(1);

}

initUpserverLog();
updatePthreadInit();

```

### Shutting Down the Duplicate Server

```

code = compat_ShutdownDuplicateServer (
    (rpc_if_handle_t)UPDATE_v4_0_s_ifspec,
    (uuid_t *)NULL,
    0);
if ( code) {
    LogError(udt_s_dup_server_shutdown_failed, code);
    goto exit;
}

```

### Recording the Shutdown Registered Interface

The following example shows how to register an interface with the versioning mechanism.

```

/* recording the above registered interface */
dfs_installInterfaceDescription ((rpc_if_handle_t)UPDATE_v4_0_s_ifspec,
    (rpc_if_handle_t)UPDATE_v4_0_s_ifspec,
    1, /* provider version number */
    (unsigned_char_t *)"Transarc Update Server Interface",
    &st);

if (IS_COMM_ERR(st))
{
    LogError(udt_s_install_reg_intf_failed, st);
    code = st;
    goto exit;
}

uuid_create_nil(&nil_uuid, &st);
if (st != uuid_s_ok) {
    LogError(udt_s_nil_UUID_create_failed, st);
    goto exit;
}

```

**Registering the Interface with the Runtime**

The following is an example of registering an interface with the RPC Runtime.

```

rpc_register_dfs_server ((rpc_if_handle_t)UPDATE_v4_0_s_ifspec,
    (rpc_mgr_epv_t)&UPDATE_v4_0_manager_epv,
    (uuid_t *)&nil_uuid, (uuid_t *)&nil_uuid,
    UPDATE_MAXCALLS, adminFilename,
    "DFS upserver", &st);
if (IS_COMM_ERR(st))
{
    LogError(udt_s_srvr_regist_failed, st);
    code = st;
    goto exit;
}

Log(udt_s_upsrvr_started);

rpc_mgmt_set_server_com_timeout(rpc_c_binding_default_timeout+2, &st);
rpc_server_listen (UPDATE_MAXCALLS, &st);
if (IS_COMM_ERR(st))
{
    LogError(udt_s_srvr_listen_failed, st);
    code = st;
    goto exit;
}

code = compat_UnregisterServer ((rpc_if_handle_t)UPDATE_v4_0_s_ifspec,
    (uuid_t *)NULL);
if (code) {
    LogError(udt_s_unregister_intf_failed, code);
    goto exit;
}
else exit(0);

```

**Sample Error Exit for Preceding Functions**

```

exit:
    exit(2);
}

```



## 5.6 Example Manager Application

The manager implements all RPCs. An additional RPC that it should support is the `RPC_GetServerInterfaces` call which is called by the client to retrieve the server interfaces.

### 5.6.1 Typical MANAGER Headers

```
#include <dcedfs/param.h>
#include <dcedfs/stds.h>
#include <dcedfs/osi.h>

#include <netdb.h>
#include <netinet/in.h>

#include <dirent.h>

#ifdef AFS_AIX_ENV
    /*
     * in AIX valuable fields such as st_blksize are not in
     * the stat structure
     */
#include <sys/statfs.h>
#endif /* AFS_AIX_ENV */

#include <dce/rpc.h>
#include <pthread.h>

#include <dcedfs/compat.h>
#include <dcedfs/dfsauth.h>

#include <update.h>
#include <upcommon.h>
```

### 5.6.2 Headers for Serviceability

The following headers get included for serviceability. The last three are specific to this subcomponent. They would need to be defined for each subcomponent that has implemented serviceability, and the names would be chosen to be specific to the function. The first three are included as is regardless of what subcomponent they are in, as they are needed for general serviceability and are not specific to the subcomponent.

The explanation of serviceability is beyond the scope of this document and this example.

```
/* For serviceability */

#include <dce/dce.h>
#include <dce/dce_msg.h>
#include <dce/dcesvcmsg.h>
#include <dfsudtmac.h>
#include <dfsudtsvc.h>
#include <dfsudtmsg.h>
```

### 5.6.3 Typical Initialization

```
extern char *whoami;
extern pthread_mutex_t updateMutex;
```

### 5.6.4 ANSI C Declaration

```
/* Forward declarations for ANSI C */
long update_SendFile();
long update_SendDirInfo();
```

### 5.6.5 Example Mutex Lock

The following macro is used as an example to demonstrate locking that could be used.

**Note:** Locking is done in this example per exported procedure, rather than per subroutine call. While it serializes operations more than per subroutine locking, it is used here to demonstrate a consistent manner of locking.

```
#define LOCK_UPDATE_MUTEX \
    if (pthread_mutex_lock (&updateMutex) != 0) \
    { \
        Log(udt_s_global_lock_locking_error, \
            whoami, errno); \
        exit(1); \
    }
```

### 5.6.6 Example Mutex Unlock

```
#define UNLOCK_UPDATE_MUTEX \
    if (pthread_mutex_unlock (&updateMutex) != 0) \
    { \
        Log(udt_s_global_lock_unlock_error, \
            whoami, errno); \
        exit(1); \
    }
```

### 5.6.7 Example Error Exit

The following example defines the common error exit that is used in the Manager's example exported procedures.

```
/* used for premature routine exit */
#define ERROR(n) \
{ \
    code = n; \
    goto error_exit; \
}
```

### 5.6.8 UPDATE\_GetServerInterfaces

This example exported procedure sends a list of registered upserver interfaces to the caller.

```

error_status_t UPDATE_GetServerInterfaces (h, serverInterfaces)
    /* [in] */   rpc_binding_handle_t      h;
    /* [out] */  dfs_interfaceList *serverInterfaces;
{
    long code = 0, code2 = 0;
    static char  routineName[] = "UPDATE_FetchFile";
    int savedCancelState;

    DFS_DISABLE_CANCEL(&savedCancelState, code);
    if (code) {
        LogError(udt_s_disable_cancel_failed, code, whoami, routineName);
        serverInterfaces->dfs_interfaceList_len = 0;
        return(code);
    }

    code = dfsauth_server_CheckAuthorization(&h,
        (dacl_permset_t *)NULL,
        (char *)NULL,
        (uid_t *)NULL);
    if ( code )
    {
        if ( code != DAUT_ERROR_ACCESS_DENIED) {
            LogError(udt_s_auth_check_failed, code, whoami, routineName);
        }
        serverInterfaces->dfs_interfaceList_len = 0;
        ERROR(code);
    }
    dfs_GetServerInterfaces (UPDATE_v4_0_s_ifspec, serverInterfaces,
        (unsigned32 *)&code);
    if (code) {
        LogError(udt_s_getserver_intfcs_failed, code, whoami, routineName);
        ERROR(code);
    }

error_exit:
    DFS_ENABLE_CANCEL(savedCancelState, code2);
    if (code2) {
        LogError(udt_s_enable_cancel_failed, code2, whoami, routineName);
        serverInterfaces->dfs_interfaceList_len = 0;
    }
    if (code)
        return code;
    else return code2;
}

```

### 5.6.9 UPDATE\_FetchFile

This example exported procedure sends the contents of a file to the caller.

```

error_status_t UPDATE_FetchFile(h, FileName, pipeStream)
/* [in] */   rpc_binding_handle_t h;
/* [in] */   NameString_t         FileName;
/* [out] */  pipe_t                *pipeStream;
{
    int      fd = -1;
    int      lockHeld = 0;
    struct   stat status;
    char     *pipeBuffP = (char *)0;
    long     code = 0, code2= 0;
    static char  routineName[] = "UPDATE_FetchFile";
    int      savedCancelState;

    /* place a TRY/FINALLY around critical region */
    TRY
    {
        DFS_DISABLE_CANCEL(&savedCancelState, code);
        if (code) {
            /* terminate the pipe */
            pipeStream->push(pipeStream->state, (unsigned char *)pipeBuffP, 0);

            LogError(udt_s_pipe_disable_cancel_failed, code, whoami,
                    routineName);
            ERROR(code);
        }

        LOCK_UPDATE_MUTEX;
        lockHeld = 1;
        code = dfsauth_server_CheckAuthorization(&h,
            (dacl_permset_t *)NULL,
            (char *)NULL,
            (uid_t *)NULL);
        if ( code )
        {
            if ( code != DAUT_ERROR_ACCESS_DENIED) {
                LogError(udt_s_daut_error_auth_failed, code, whoami, routineName);
            }

            /* terminate the pipe */
            pipeStream->push(pipeStream->state, (unsigned char *)pipeBuffP, 0);
            ERROR(code);
        }
    }
    if ( canExportObject(FileName) == 0 )
    {
        /* terminate the pipe */
        pipeStream->push(pipeStream->state, (unsigned char *)pipeBuffP, 0);

        Log(udt_s_req_not_in_list, FileName);
        ERROR(UP_NOT_EXPORTABLE);
    }
}

```

```

    }

    fd = open(FileName, O_RDONLY, 0);          /* Open the target file */
    if (fd < 0 || fstat(fd, &status) < 0) {
        /* terminate the pipe */
        pipeStream->push(pipeStream->state, (unsigned char *)pipeBuffP, 0);

        Log(udt_s_open_file_failed, FileName);
        ERROR(UP_SOFT_ERROR);
    }
    code = update_SendFile(fd, pipeStream, &status);
    if (fd >=0)
        close(fd);

    error_exit;;
}
FINALLY
{
    if (lockHeld)
        UNLOCK_UPDATE_MUTEX;
    DFS_ENABLE_CANCEL(savedCancelState, code2);
}
ENDTRY

if (code2)
    LogError(udt_s_enable_cancel_failed2, code2, whoami, routineName);
if (code)
    return code;
else
    return code2;
}

```

### 5.6.10 UPDATE\_FetchInfo

This example exported procedure fetches directory information about directory name and sends it to remote client.

```

error_status_t UPDATE_FetchInfo(h, DirName, pipeStream)
    /* [in] */   rpc_binding_handle_t h;
    /* [in] */   NameString_t         DirName;
    /* [out] */  pipe_t                *pipeStream;
{
    int fd = -1;
    int lockHeld = 0;
    struct stat status;
    char *pipeBuffP = (char *)0;
    long code = 0, code2 = 0;
    static char routineName[] = "UPDATE_FetchInfo";
    int savedCancelState;

    TRY
    {
        DFS_DISABLE_CANCEL(&savedCancelState, code);
    }
}

```

```

if (code) {
    /* terminate the pipe */
    pipeStream->push(pipeStream->state, (unsigned char *)pipeBuffP, 0);

    LogError(udt_s_disable_cancel_failed2, code, whoami, routineName);
    ERROR(code);
}

LOCK_UPDATE_MUTEX;
lockHeld = 1;
code = dfsauth_server_CheckAuthorization(&h,
    (dacl_permset_t *)NULL,
    (char *)NULL,
    (uid_t *)NULL);
if ( code )
{
    if (code != DAUT_ERROR_ACCESS_DENIED) {
        LogError(udt_s_daut_error_auth_failed2, code, whoami, routineName);
    }

    /* terminate the pipe */
    pipeStream->push(pipeStream->state, (unsigned char *)pipeBuffP, 0);
    ERROR(code);
}
if ( canExportObject(DirName) == 0 )
{
    /* terminate the pipe */
    pipeStream->push(pipeStream->state, (unsigned char *)pipeBuffP, 0);

    Log(udt_s_req_not_in_list2, DirName);
    ERROR(UP_NOT_EXPORTABLE);
}

fd = open(DirName, O_RDONLY, 0);
if (fd < 0 || fstat(fd, &status) < 0) {
    /* terminate the pipe */
    pipeStream->push(pipeStream->state, (unsigned char *)pipeBuffP, 0);

    Log(udt_s_open_dir_failed, DirName);
    ERROR(UP_SOFT_ERROR);
}

if((status.st_mode & S_IFMT) != S_IFDIR){
    /* terminate the pipe */
    pipeStream->push(pipeStream->state, (unsigned char *)pipeBuffP, 0);

    Log(udt_s_obj_not_dir,DirName);
    ERROR(ENOENT);
}
code = update_SendDirInfo(DirName, pipeStream, &status);

error_exit;;

```

```

    }
    FINALLY
    {
    if (fd >= 0)
        close(fd);
    if (lockHeld)
        UNLOCK_UPDATE_MUTEX;

    DFS_ENABLE_CANCEL(savedCancelState, code2);
    }
    ENENTRY

    if (code2)
        LogError(udt_s_enable_cancel_failed3, code2, whoami, routineName);
    if (code)
        return code;
    else
        return code2;
}

```

### 5.6.11 UPDATE\_FetchObjectInfo

This example exported procedure returns a stat of the named object. Selected items from the stat are returned to the caller.

```

error_status_t
UPDATE_FetchObjectInfo(h, objectName, fileStatP)
    /* [in] */  rpc_binding_handle_t h;
    /* [in] */  NameString_t      objectName;
    /* [out]*/  updateFileStatP   fileStatP;
{
    struct stat statBuf;
    long code = 0, code2 = 0;
    int savedCancelState;
    static char routineName[] = "UPDATE_FetchObjectInfo";

    DFS_DISABLE_CANCEL(&savedCancelState, code);
    if (code) {
        LogError(udt_s_disable_cancel_failed4, code, whoami, routineName);
        return(code);
    }

    LOCK_UPDATE_MUTEX;
#ifdef notdef
    code = dfsauth_server_CheckAuthorization(&h,
        (dacl_permset_t *)NULL,
        (char *)NULL,
        (uid_t *)NULL);

    if ( code )
    {
        if (code != DAUT_ERROR_ACCESS_DENIED) {
            LogError(udt_s_daut_error_auth_failed4, code, whoami,

```

```

        routineName);
    }
    ERROR(code);
}
#endif
if ( canExportObject(objectName) == 0 )
{
    Log(udt_s_obj_export_failed, objectName);
    ERROR(UP_NOT_EXPORTABLE);
}
code = stat((char *)objectName, &statBuf);
if ( code )
{
    Log(udt_s_stat_obj_failed, objectName, errno);
    ERROR(UP_SOFT_ERROR);
}

/* copy the information back to the caller's structure */
hset32(fileStatP->fileLength, statBuf.st_size);
fileStatP->mode =    statBuf.st_mode;
fileStatP->uid =    statBuf.st_uid;
fileStatP->gid =    statBuf.st_gid;
fileStatP->mtime =    statBuf.st_mtime;
fileStatP->atime =    statBuf.st_atime;

error_exit:
    UNLOCK_UPDATE_MUTEX;

    DFS_ENABLE_CANCEL(savedCancelState, code2);
    if (code2)
        LogError(udt_s_enable_cancel_failed4, code2, whoami, routineName);
    if (code)
        return code;
    else
        return code2;
}

```

### 5.6.12 UPDATE\_v4\_0\_manager\_epv

This example declares and initializes the example manager's entry point vector.

```

UPDATE_v4_0_epv_t UPDATE_v4_0_manager_epv = {
    UPDATE_GetServerInterfaces,
    UPDATE_FetchInfo,
    UPDATE_FetchFile,
    UPDATE_FetchObjectInfo};

```



## 5.7 Extending the DFS Interface

The owner of an interface implies it is the entity that created the interface with a UUID, major and minor version numbers. Other entities can extend the interface. When an entity extends an interface not owned by it, it needs to create a new interface (UUID, major and minor version numbers) and add new operations in addition to the operations from the interface that is being extended. The DFS interfaces are owned by Transarc Corporation and others would need to create a new interface in order to extend any of the original DFS interfaces.

There are two cases involved in extending an interface. One case is when an owner of an interface extends it; the other is when an entity that is not the owner of an interface extends it. These two cases involve the following changes:

### 5.7.1 Case 1 — An Owner Extends a DFS Interface

When an owner extends a DFS interface, the changes needed are the following:

1. Deciding either to bump the minor version number or the provider version number. This decision is made based upon whether the changes are extensive enough to warrant it's being a new version or not. If extensive, then the provider version number would be increased by one which connotes the meaning of being bumped. Otherwise, increasing or bumping the minor version number by one would be sufficient.
2. If the minor version number is bumped, the change is recorded in the IDL file header. Also a change is to be made in the *dfs\_selectInterface()* function in the client, to select the appropriate interface(stub) to call.
3. If the provider version number is bumped, then the IDL file header should have the new provider version number, and each added RPC should have a comment indicating its provider version level. The server should indicate its provider level when calling *dfs\_installInterfaceDescription()*. Servers must maintain compatibility with clients running at lower provider version levels. Clients must not rely on the availability of any RPCs other than those in provider version 1. When a client makes an RPC to a server running at a lower provider version level, the request is handled by the RPC runtime environment, and the client gets an [rpc\_s\_op\_rng\_error]. Clients should be written so that they can recover from this error by retrying the operation using the functionality available in provider version 1.

Clients should keep an estimate of the provider version level of each server. Initially, this estimate should be the provider version level desired by the client. If an [rpc\_s\_op\_rng\_error] is received for an RPC at provider version N, the server's provider version estimate should be set to N-1. The client should not make any more RPC calls to that server at a provider version higher than N-1. The server's provider version estimate should be reset whenever the server's binding handle is reset.

As an example, the Ubik client interface includes support for server provider versioning, and maintains an estimate of each server's provider version level. Ubik clients making RPCs for provider versions other than provider version 1 should use *ubik\_PVCall()* instead of *ubik\_Call()*. *ubik\_PVCall()* has a parameter for the provider version level of the RPC. The ubik client library keeps an estimate of the provider version level of each server, and returns [UDOWNLVL] when an RPC fails because the server is not at the desired version level. Clients should be written to recover from [UDOWNLVL] errors by retrying the operation using the functionality in provider version 1.

### 5.7.2 Case 2 — A Non-owner Extends a DFS Interface

An entity not owning an interface will herein be called an agent. When an entity extends a DFS interface not owned by it, the changes needed are the following:

1. If the agent modifies an interface not owned (created) by it, then a new Interface UUID (and a new IDL file) should be created. The RPCs in the original interface are supported in the new interface (with a distinguishing prefix) and new RPCs are appended to the older ones.
2. The original DFS interface (Transarc's) and the newly created interface are registered with the RPC runtime and the versioning mechanism. The agent has the option of supporting any intermediate interfaces.
3. A new manager end point vector is created for the new interface with the old entries and the newly added RPCs appended to this. The old entries could refer to functions in any of the managers associated with the interfaces other than the newly added one. This way, the new manager need not duplicate functions already implemented in the previous interfaces.
4. The *dfs\_selectInterface()* function in the client has to be modified to select the new interface. If not available, then the client could use one of the older interfaces supported by the client, and the wrapper for the newly added RPCs could check the interface global to determine if the RPC could be made successfully or not.
5. The wrappers for the new RPCs should check for the new interface. If not supported, then the new RPCs should fail.

# **/** *X/Open Preliminary Specification*

## **Part 4:**

### **DFS Client/Server and Token Manager Interfaces**



## DCE DFS File Exporter Interface

This chapter describes the RPCs that can be issued against the file exporter via a client cache manager. These RPCs allow the cache manager to manipulate files on a file server machine.

The error codes passed back from the DFS server to the DFS client do not come through [ERRNO], and are not the values from the server's version of <errno.h>, but are an architected set that the client will then map into local error codes.

The return codes returned by these calls are the union of the set of errors returned by the RPC, the token manager, and an abstraction of the DFS server's vnode layer itself (for example, DFS\_ENOENT for file not found). The full set of these abstractions of the DFS server's error codes is as follows:

```
#define DFS_ESUCCESS      0    /* Successful */
#define DFS_EPERM         1    /* Operation not permitted */
#define DFS_ENOENT        2    /* No such file or directory */
#define DFS_ESRCH         3    /* No such process */
#define DFS_EINTR         4    /* Interrupted system call */
#define DFS_EIO           5    /* I/O error */
#define DFS_ENXIO         6    /* No such device or address */
#define DFS_E2BIG         7    /* Arg list too long */
#define DFS_ENOEXEC       8    /* Exec format error */
#define DFS_EBADF         9    /* Bad file number */
#define DFS_ECHILD        10   /* No children */
#define DFS_EDEADLK       11   /* Operation would cause deadlock */
#define DFS_ENOMEM        12   /* Not enough core */
#define DFS_EACCES        13   /* Permission denied */
#define DFS_EFAULT        14   /* Bad address */
#define DFS_ENOTBLK       15   /* Block device required */
#define DFS_EBUSY         16   /* Mount device busy */
#define DFS_EEXIST        17   /* File exists */
#define DFS_EXDEV         18   /* Cross-device link */
#define DFS_ENODEV        19   /* No such device */
#define DFS_ENOTDIR       20   /* Not a directory */
#define DFS_EISDIR        21   /* Is a directory */
#define DFS_EINVAL        22   /* Invalid argument */
#define DFS_ENFILE        23   /* File table overflow */
#define DFS_EMFILE        24   /* Too many open files */
#define DFS_ENOTTY        25   /* Not a typewriter */
#define DFS_ETXTBSY       26   /* Text file busy */
#define DFS_EFBIG         27   /* File too large */
#define DFS_ENOSPC        28   /* No space left on device */
#define DFS_ESPIPE        29   /* Illegal seek */
#define DFS_EROFS         30   /* Read-only file system */
#define DFS_EMLINK        31   /* Too many links */
#define DFS_EPIPE         32   /* Broken pipe */
#define DFS_EDOM          33   /* Argument too large */
#define DFS_ERANGE        34   /* Result too large */
#define DFS_EWOULDBLOCK   35   /* Operation would block */
#define DFS_EINPROGRESS   36   /* Operation now in progress */
#define DFS_EALREADY      37   /* Operation already in progress */
```

```

/*
 * ipc/network software
 */
#define DFS_ENOTSOCK 38 /* Socket operation on non-socket */
#define DFS_EDESTADDRREQ 39 /* Destination address required */
#define DFS_EMSGSIZE 40 /* Message too long */
#define DFS_EPROTOTYPE 41 /* Protocol wrong type for socket */
#define DFS_ENOPROTOOPT 42 /* Protocol not available */
#define DFS_EPROTONOSUPPORT 43 /* Protocol not supported */
#define DFS_ESOCKTNOSUPPORT 44 /* Socket type not supported */
#define DFS_EOPNOTSUPP 45 /* Operation not supported on socket */
#define DFS_EPFNOSUPPORT 46 /* Protocol family not supported */
#define DFS_EAFNOSUPPORT 47 /* Address family not supported */
#define DFS_EADDRINUSE 48 /* Address already in use */
#define DFS_EADDRNOTAVAIL 49 /* Can't assign requested address */
#define DFS_ENETDOWN 50 /* Network is down */
#define DFS_ENETUNREACH 51 /* Network is unreachable */
#define DFS_ENETRESET 52 /* Network dropped conn on reset */
#define DFS_ECONNABORTED 53 /* Software caused connection abort */
#define DFS_ECONNRESET 54 /* Connection reset by peer */
#define DFS_ENOBUFS 55 /* No buffer space available */
#define DFS_EISCONN 56 /* Socket is already connected */
#define DFS_ENOTCONN 57 /* Socket is not connected */
#define DFS_ESHUTDOWN 58 /* Can't send after socket shutdown */
#define DFS_ETOOMANYREFS 59 /* Too many references: can't splice */
#define DFS_ETIMEOUT 60 /* Connection timed out */
#define DFS_ECONNREFUSED 61 /* Connection refused */
#define DFS_ELOOP 62 /* Too many levels of symbolic links */
#define DFS_ENAMETOOLONG 63 /* File name too long */
#define DFS_EHOSTDOWN 64 /* Host is down */
#define DFS_EHOSTUNREACH 65 /* No route to host */
#define DFS_ENOTEMPTY 66 /* Directory not empty */
#define DFS_EPROCLIM 67 /* Too many processes */
#define DFS_EUSERS 68 /* Too many users */
#define DFS_EDQUOT 69 /* Disc quota exceeded */
/*
 * NFS errors.
 */
#define DFS_ESTALE 70 /* Stale NFS file handle */
#define DFS_EREMOTE 71 /* Too many levels of remote in path */
#define DFS_EBADRPC 72 /* RPC struct is bad */
#define DFS_ERPCMISMATCH 73 /* RPC version wrong */
#define DFS_EPROGUNAVAIL 74 /* RPC prog. not avail */
#define DFS_EPROGMISMATCH 75 /* Program version wrong */
#define DFS_EPROCUNAVAIL 76 /* Bad procedure for program */
/*
 * misc ..
 */
#define DFS_ENOLCK 77 /* No locks available */
#define DFS_ENOSYS 78 /* Function not implemented */
#define DFS_EAGAIN 79 /* Resource temporarily unavailable */
/*

```

## DCE DFS File Exporter Interface

```

    * SYS V IPC errors
    */
#define DFS_ENOMSG          80 /* No msg matches receive request */
#define DFS_EIDRM          81 /* Msg queue id has been removed */
/*
    * STREAMS
    */
#define DFS_ENOSR          82 /* Out of STREAMS resources */
#define DFS_ETIME          83 /* System call timed out */
#define DFS_EBADMSG       84 /* Next message has wrong type */
#define DFS_EPROTO        85 /* STREAMS protocol error */
#define DFS_ENODATA       86 /* No message on stream head read q */
#define DFS_ENOSTR        87 /* fd not associated with a stream */
/*
    * Not visible outside kernel
    */
#define DFS_ECLONEME       88 /* Tells open to clone the device */
/*
    * Filesystem
    */
#define DFS_EDIRTY        89 /* Mounting a dirty fs w/o force */
/*
    * Loader errors
    */
#define DFS_EDUPPKG       90 /* duplicate package name on install */
#define DFS_EVERSION      91 /* version number mismatch */
#define DFS_ENOPKG       92 /* unresolved package name */
#define DFS_ENOSYM       93 /* unresolved symbol name */
/*
    * To be filled
    */
#define DFS_SPARE94       94
#define DFS_SPARE95       95
#define DFS_SPARE96       96
#define DFS_SPARE97       97
#define DFS_SPARE98       98
#define DFS_SPARE99       99
#define DFS_SPARE100      100
#define DFS_SPARE101      101
#define DFS_SPARE102      102
#define DFS_SPARE103      103
#define DFS_SPARE104      104
#define DFS_SPARE105      105
#define DFS_SPARE106      106
#define DFS_SPARE107      107
#define DFS_SPARE108      108
#define DFS_SPARE109      109
#define DFS_SPARE110      110
#define DFS_SPARE111      111
/*
    * security
    */
```

```
#define DFS_ENOATTR          112 /* no attribute found */
#define DFS_ESAD             113 /* security authentication denied */
#define DFS_ENOTRUST        114 /* not a trusted program */
/*
 * To be filled
 */
#define DFS_SPARE115         115
#define DFS_SPARE116         116
#define DFS_SPARE117         117
#define DFS_SPARE118         118
#define DFS_SPARE119         119
#define DFS_SPARE120         120
#define DFS_SPARE121         121
#define DFS_SPARE122         122
/*
 * Internal Disk/Block Device error codes
 */
#define DFS_ESOFT            123 /* I/O completed but needs relocation*/
#define DFS_EMEDIA           124 /* media surface error */
#define DFS_ERELOCATED       125 /* a relocation request performed ok */
#define DFS_SPARE126         126
#define DFS_SPARE127         127
/*
 * Errors not defined by local kernel.
 */
#define DFS_ENOTDEFINED      128
```



**NAME**

AFS\_SetContext — Establish client context

**SYNOPSIS**

```

error_status_t AFS_SetContext(
    /* IN */ handle_t    h,
    /* IN */ unsigned32  epochTime,
    /* IN */ afsNetData *callbackAddr,
    /* IN */ unsigned32  Flags,
    /* IN */ afsUUID     *secObjectID,
    /* IN */ unsigned32  clientSizeAttrs,
    /* IN */ unsigned32  parm7
);

```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>epochTime</i>	The restart time of the DFS client.
<i>callbackAddr</i>	The RPC endpoint of the client, for token revocation purposes.
<i>Flags</i>	If 0 then the call only defines a primary interface UUID. If the flag, AFS_FLAG_SEC_SERVICE, is set then a secondary interface is defined and is stored in the parameter, <i>secObjectID</i> .
<i>secObjectID</i>	When the AFS_FLAG_SEC_SERVICE flag is set, this parameter holds the secondary interface UUID.
<i>clientSizeAttrs</i>	For 64/32-bit compatibility. Through DCE1.1, this is a spare parameter, <i>parm6</i> .
<i>parm7</i>	A spare parameter.

**DESCRIPTION**

This function establishes some context information at the file exporter that is required before this client can execute calls at the server. If a server crashes, the client must make a new AFS\_SetContext call before it makes other calls to this server. The UUID's passed via this function must have been obtained via the *dfsuuid\_Create()* function.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

[DFS_EINVAL]	The caller passed in some invalid argument.
[FSHS_ERR*]	File server does not recognize client, or is in some form of recovery.

**NAME**

AFS\_LookupRoot — Look up fileset root on file server

**SYNOPSIS**

```
error_status_t AFS_LookupRoot(
    /* IN */ handle_t    h,
    /* IN */ afsFid     *InFidp,
    /* IN */ afsHyper   *minVVp,
    /* IN */ signed32    Flags,
    /* OUT */ afsFid     *OutFidp,
    /* OUT */ afsFetchStatus *OutFidStatusp,
    /* OUT */ afsToken   *OutTokenp,
    /* OUT */ afsVolSync  *Syncp
);
```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>InFidp</i>	The file identifier specifying the fileset whose root directory will be retrieved from the file server.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier describing the root directory of the selected fileset.
<i>OutFidStatusp</i>	Returns the status of that directory after the current operation.
<i>OutTokenp</i>	The promise the file server returns to the cache manager about the provided data; this is only returned if the file resides in a read/write fileset.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**DESCRIPTION**

This call is used to obtain the description of the root directory of a fileset from the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

[DFS_EINVAL]	The caller passed in some invalid argument.
[DFS_ENOENT]	The file specified by <i>inFidp</i> does not exist.
[DFS_EACCES]	The caller lacks read permission on <i>inFidp</i> .
[FSHS_ERR*]	File server does not recognize client, or is in some form of recovery.
[TKM_ERROR*]	A token management failure occurred.

**NAME**

AFS\_FetchData — Retrieve data from file server

**SYNOPSIS**

```

error_status_t AFS_FetchData(
    /* IN */ handle_t    h,
    /* IN */ afsFid     *Fidp,
    /* IN */ afsHyper   *minVVp,
    /* IN */ afsHyper   *Position,
    /* IN */ signed32   Length,
    /* IN */ unsigned32 Flags,
    /* OUT */ afsFetchStatus *OutStatusp,
    /* OUT */ afsToken   *OutTokenp,
    /* OUT */ afsVolSync *Syncp,
    /* OUT */ pipe_t    *fetchStream
);

```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file whose contents will be retrieved from the file server.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Position</i>	Specifies the first byte to be fetched by this call with zero (0) being the first byte in the file.
<i>Length</i>	Specifies the number of bytes desired with the value 0xFFFFFFFF indicating the entire file contents.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutStatusp</i>	Returns the status of the file after the current operation.
<i>OutTokenp</i>	The promise the file server returns to the cache manager about the provided data; this is only returned if the file resides in a read/write fileset.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.
<i>fetchStream</i>	The character pipe parameter returning the data from the file. (See <b>X/Open DCE: Remote Procedure Call</b> specification for implementation details.)

**DESCRIPTION**

This call is used to obtain the contents of the specified file from the file server.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

- [DFS\_EINVAL] The caller passed in some invalid argument.
- [DFS\_ENOENT] The file specified by *Fidp* does not exist.

- [DFS\_EACCES] The caller lacks read permission on *Fidp*.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*] A token management failure occurred.

**NAME**

AFS\_FetchACL — Retrieve Access Control List

**SYNOPSIS**

```

error_status_t AFS_FetchACL(
    /* IN */ handle_t h,
    /* IN */ afsFid *Fidp,
    /* IN */ unsigned32 aclType,
    /* IN */ afsHyper *minVVp,
    /* IN */ unsigned32 Flags,
    /* OUT */ afsACL *AccessListp,
    /* OUT */ afsFetchStatus *OutStatusp,
    /* OUT */ afsVolSync *Syncp
);

```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file or directory whose access control list (ACL) will be retrieved from the file server.
<i>aclType</i>	The type of the access list being modified. One of VNX_ACL_REGULAR_ACL, VNX_ACL_DEFAULT_ACL or VNX_ACL_INITIAL_ACL.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this ACL.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>AccessListp</i>	The access control list returned by the file server for the specified file or directory.
<i>OutStatusp</i>	Returns the current status of the file or directory.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**DESCRIPTION**

Note that a new token is not returned by this call; however, any existing tokens remain in effect.

See *Access Control List Overview* for a description of access list encoding.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

- [DFS\_EINVAL] The caller passed in some invalid argument.
- [DFS\_ENOENT] The file, dir or symlink specified by *Fidp* does not exist.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*] A token management failure occurred.

**NAME**

AFS\_FetchStatus — Obtain file status information

**SYNOPSIS**

```
error_status_t AFS_FetchStatus(
    /* IN */ handle_t    h,
    /* IN */ afsFid      *Fidp,
    /* IN */ afsHyper    *minVvp,
    /* IN */ unsigned32  Flags,
    /* OUT */ afsFetchStatus *OutStatusp,
    /* OUT */ afsToken    *OutTokenp,
    /* OUT */ afsVolSync  *Syncp
);
```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file or directory whose status information will be retrieved from the file server.
<i>minVvp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutStatusp</i>	Returns the current status of the file or directory.
<i>OutTokenp</i>	Returns a token from the file server reflecting guarantees granted by the file server.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**DESCRIPTION**

This call retrieves the status information associated with the specified file or directory. The file server returns a token pertaining to this file if the fileset upon which it resides is read/write.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

[DFS_EINVAL]	The caller passed in some invalid argument.
[DFS_ENOENT]	The file named by <i>Fidp</i> does not exist.
[FSHS_ERR*]	File server does not recognize client, or is in some form of recovery.
[TKM_ERROR*]	A token management failure occurred.

**NAME**

AFS\_StoreData — Write data to file

**SYNOPSIS**

```

error_status_t AFS_StoreData(
    /* IN */ handle_t    h,
    /* IN */ afsFid      *Fidp,
    /* IN */ afsFetchStatus *InStatusp,
    /* IN */ afsHyper     *Position,
    /* IN */ signed32     Length,
    /* IN */ afsHyper     *minVVp,
    /* IN */ unsigned32   Flags,
    /* IN */ pipe_t       *storeStream,
    /* OUT */ afsStoreStatus *OutStatusp,
    /* OUT */ afsVolSync  *Syncp
);

```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file whose status information will be updated from the file server.
<i>InStatusp</i>	The new status information that should be recorded for this file.
<i>Position</i>	Represents the position of the first byte of the data block.
<i>Length</i>	Represents the total length of the transferred data block.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>storeStream</i>	The data stream containing the file updates. (See the <b>X/Open DCE: Remote Procedure Call</b> specification for implementation details.)
<i>OutStatusp</i>	Returns the status of the file.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**DESCRIPTION**

This call writes the specified data back to the file server for the specified file.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

- [DFS\_EINVAL] The caller passed in some invalid argument.
- [DFS\_ENOENT] The file specified by *Fidp* does not exist.
- [DFS\_EACCES] The caller lacks write permission to perform this operation.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.

[TKM\_ERROR\*] A token managment failure occurred.



**NAME**

AFS\_StoreACL — Update access control information for a file

**SYNOPSIS**

```
error_status_t AFS_StoreACL(
    /* IN */ handle_t h,
    /* IN */ afsFid *Fidp,
    /* IN */ afsACL *AccessListp,
    /* IN */ unsigned32 aclType,
    /* IN */ afsFid *aclFidp,
    /* IN */ afsHyper *minVvp,
    /* IN */ unsigned32 Flags,
    /* OUT */ afsFetchStatus *OutStatusp,
    /* OUT */ afsVolSync *Syncp
);
```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file or directory whose access control information will be updated from the file server.
<i>AccessListp</i>	The access control list sent to the file server for the specified file or directory.
<i>aclType</i>	The type of access control list being modified (VNX_ACL_REGULAR_ACL, VNX_ACL_DEFAULT_ACL or VNX_ACL_INITIAL_ACL, in the low-order 8 bits. In the next-higher-order 8 bits, the type of access control list being copied from the file described by <i>aclFidp</i> ).
<i>aclFidp</i>	The file identifier specifying the file or directory whose access control information will be copied.
<i>minVvp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutStatusp</i>	Returns the updated file or directory status.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**DESCRIPTION**

This call updates the access control information associated with a file or directory. It can either set an ACL to a specific value via *AccessListp* or copy an ACL from that possessed by another file via *aclFidp*.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

- [DFS\_EINVAL] The caller passed in some invalid argument.
- [DFS\_ENOENT] The file specified by *Fidp* does not exist.

[DFS\_EACCES] The caller lacks the required access to this file.

[FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.

[TKM\_ERROR\*] A token management failure occurred.

**NAME**

AFS\_StoreStatus — Update status information

**SYNOPSIS**

```

error_status_t AFS_StoreStatus(
    /* IN */ handle_t h,
    /* IN */ afsFid *Fidp,
    /* IN */ afsStoreStatus *InStatusp,
    /* IN */ afsHyper *minVVp,
    /* IN */ unsigned32 Flags,
    /* OUT */ afsFetchStatus *OutStatusp,
    /* OUT */ afsVolSync *Syncp
);

```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	The file identifier specifying the file or directory whose access control information will be updated from the file server.
<i>InStatusp</i>	Contains the new status information for the specified file or directory.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutStatusp</i>	Contains the updated status information for the file or directory.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**DESCRIPTION**

This call updates the status information associated with the file or directory.

It returns the new status information because this operation may change some status fields.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

[DFS_EINVAL]	The caller passed in some invalid argument.
[DFS_ENOENT]	The entry named by <i>Fidp</i> does not exist.
[DFS_EACCES]	The caller lacks the required access rights.
[FSHS_ERR*]	File server does not recognize client, or is in some form of recovery.
[TKM_ERROR*]	A token management failure occurred.

## NAME

AFS\_RemoveFile — Remove a file or symbolic link

## SYNOPSIS

```
error_status_t AFS_RemoveFile(
    /* IN */ handle_t    h,
    /* IN */ afsFid      *DirFidp,
    /* IN */ afsFidTaggedName *Namep,
    /* IN */ afsHyper     *returnTokenIDp,
    /* IN */ afsHyper     *minVVp,
    /* IN */ unsigned32   Flags,
    /* OUT */ afsFetchStatus *OutDirStatusp,
    /* OUT */ afsFetchStatus *OutFileStatusp,
    /* OUT */ afsFid        *OutFileFidp,
    /* OUT */ afsVolSync    *Syncp
);
```

## ARGUMENTS

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory from which to remove the file.
<i>Namep</i>	The complex name of the file to delete.
<i>returnTokenIDp</i>	A token ID being returned, if any.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutDirStatusp</i>	Contains the updated directory status information.
<i>OutFileStatusp</i>	Contains the updated file status information.
<i>OutFileFidp</i>	The file id of the file which was actually deleted.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## DESCRIPTION

This call removes a file or symbolic link, but not a directory, from the file system. The complex name **afsFidTaggedName** *must* include the string name component - it should be filled in with both a **fid** and an **afsTaggedName**. The **afsTaggedName** should have a *codesetTag* (*tn\_tag*) value of 0 as nothing else is presently defined, and a valid *tn\_length* field (if something other than zero was specified for *tn\_tag*).

The cache manager is responsible for decrementing the link count in the file's associated cached status by 1.

## RETURN VALUE

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

## ERRORS

This function fails if:

[DFS\_EINVAL] The caller passed in some invalid argument.

- [DFS\_ENOENT] The entry named by *DirFidp* does not exist, or the entry named by *Namep* does not exist.
- [DFS\_EACCES] The caller lacks the required access rights.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*] A token management failure occurred.

**NAME**

AFS\_Lookup — Obtain directory entry for a file

**SYNOPSIS**

```
error_status_t AFS_Lookup(
    /* IN */ handle_t    h,
    /* IN */ afsFid     *DirFidp,
    /* IN */ afsTaggedName *Namep,
    /* IN */ afsHyper    *minVVp,
    /* IN */ unsigned32  Flags,
    /* OUT */ afsFid     *OutFidp,
    /* OUT */ afsFetchStatus *OutFidStatusp,
    /* OUT */ afsFetchStatus *OutDirStatusp,
    /* OUT */ afsToken     *OutTokenp,
    /* OUT */ afsVolSync   *Syncp
);
```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory from which to obtain the directory information.
<i>Namep</i>	The character string name of the file for which information is being requested.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier of the requested file.
<i>OutFidStatusp</i>	Status information for the specified file.
<i>OutDirStatusp</i>	Status information of the specified directory.
<i>OutTokenp</i>	The token against the directory (allowing the directory entry to be cached).
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**DESCRIPTION**

This call takes a directory identifier and a file name, and returns the file ID of the named file. It returns the status of the directory at the time the operation was performed, as well as a callback on that directory, so that this mapping can be cached.

A non-existent entry is denoted *not* by a special error code, but instead by a special file ID (represented as all zeroes), so that a callback can be returned for entries that don't exist. Caching the non-existence of entries may turn out to be important in some cases.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

[DFS\_EINVAL] The caller passed in some invalid argument.

- [DFS\_ENOENT] The entry named by *DirFidp* does not exist (note that entries that don't exist are handled by success returns).
- [DFS\_EACCES] The caller lacks the required access rights to the directory.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*] A token management failure occurred.

## NAME

AFS\_CreateFile — Create a file

## SYNOPSIS

```
error_status_t AFS_CreateFile(
    /* IN */ handle_t h,
    /* IN */ afsFid *DirFidp,
    /* IN */ afsTaggedName *Namep,
    /* IN */ afsStoreStatus *InStatusp,
    /* IN */ afsHyper *minVVp,
    /* IN */ unsigned32 Flags,
    /* OUT */ afsFid *OutFidp,
    /* OUT */ afsFetchStatus *OutFidStatusp,
    /* OUT */ afsFetchStatus *OutDirStatusp,
    /* OUT */ afsToken *OutTokenp,
    /* OUT */ afsVolSync *Syncp
);
```

## ARGUMENTS

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which to create the requested file.
<i>Namep</i>	The character string name of the file to create.
<i>InStatusp</i>	Specifies the initial status fields for the new file.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier of the newly created file.
<i>OutFidStatusp</i>	The status fields for the newly created file.
<i>OutDirStatusp</i>	The status information of the specified directory.
<i>OutTokenp</i>	A new token granted against the new file.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## DESCRIPTION

This call is used to create a file, but not a symbolic link or a directory.

If the call succeeds, it is the cache manager's responsibility to either create an entry locally in the directory specified by *DirFidp*, or to invalidate this directory's cache entry.

## RETURN VALUE

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

## ERRORS

This function fails if:

- [DFS\_ENOENT] The entry named by *DirFidp* does not exist.
- [DFS\_EACCES] The caller lacks the required access rights.



- [DFS\_EEXIST] An entry named *Namep* already exists in this dir.
- [DFS\_EINVAL] The entry is named “.” or “..”, or the caller passed in some invalid argument.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*] A token managment failure occurred.

## NAME

AFS\_Rename — Rename a file, directory or symbolic link

## SYNOPSIS

```
error_status_t AFS_Rename(
    /* IN */ handle_t h,
    /* IN */ afsFid *OldDirFidp,
    /* IN */ afsFidTaggedName *OldNamep,
    /* IN */ afsFid *NewDirFidp,
    /* IN */ afsFidTaggedName *NewNamep,
    /* IN */ afsHyper *returnTokenIDp,
    /* IN */ afsHyper *minVvp,
    /* IN */ unsigned32 Flags,
    /* OUT */ afsFetchStatus *OutOldDirStatusp,
    /* OUT */ afsFetchStatus *OutNewDirStatusp,
    /* OUT */ afsFid *OutOldFileFidp,
    /* OUT */ afsFetchStatus *OutOldFileStatusp,
    /* OUT */ afsFid *OutNewFileFidp,
    /* OUT */ afsFetchStatus *OutNewFileStatusp,
    /* OUT */ afsVolSync *Syncp
);
```

## ARGUMENTS

*h* The RPC binding handle.

*OldDirFidp* The file identifier specifying the directory in which the file is currently located.

*OldNamep* The complex name of the file or directory to rename.

*NewDirFidp* The file identifier specifying the directory into which the file is to be moved.

*NewNamep* The complex name of the file or directory after it is moved.

*returnTokenIDP* A token ID being returned, if any.

*minVvp* The minimum-acceptable version number on the fileset containing this file or directory.

*Flags* Any of the AFS\_FLAG\_\* flags defined in the <afs4int.h> file.

*OutOldDirStatusp* The status information of the old directory (the one from which the file is being moved) upon termination of the call.

*OutNewDirStatusp* The status information of the new directory (the one to which the file is being moved) upon termination of the call.

*OutOldFileFidp* The file identifier for the file which was moved.

*OutOldFileStatusp* The status information of the file which was moved.

*OutNewFileFidp* The file identifier for the file to which the file identified in *OutOldFileFidp* was, in fact, moved.

*OutNewFileStatusp* The status information of the file identified by *OutNewFileFidp*.

*Syncp*            The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

#### DESCRIPTION

This call renames the specified file or directory, located in the specified directory, to be the new file or directory located in the newly specified directory. The updated directory status for both directories is returned. If the two directories are the same, identical status information is returned in each status output parameter.

The rename must not result in hard links existing to the same object from two different filesets, or the error code EXDEV will be returned.

If a directory is moved from one directory to another, the cache manager must either update the cached copy of the moved directory in order to update its '..' entry, or the cache manager must invalidate the cache entry for the moved directory. The directory link counts will be updated by the server in the returned directory status blocks.

#### RETURN VALUE

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

#### ERRORS

This function fails if:

- [DFS\_ENOENT]    The entry named by *Fidp* does not exist.
- [DFS\_EACCES]    The caller lacks the required access rights.
- [DFS\_ENOTEMPTY]            The target specifies a non-empty directory.
- [DFS\_ENOTDIR]    The source specifies a directory but the target is not a directory.
- [DFS\_EISDIR]     The source specifies a file but the target specifies a directory.
- [DFS\_EXDEV]      Attempting to rename across filesets.
- [DFS\_EINVAL]     The entry is named "." or "..", or the caller passed in some invalid argument.
- [FSHS\_ERR\*]      File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*]    A token management failure occurred.

**NAME**

AFS\_Symlink — Create a symbolic link

**SYNOPSIS**

```

error_status_t AFS_Symlink(
    /* IN */ handle_t    h,
    /* IN */ afsFid      *DirFidp,
    /* IN */ afsTaggedName *Namep,
    /* IN */ afsTaggedPath *LinkContentsp,
    /* IN */ afsStoreStatus *InStatusp,
    /* IN */ afsHyper      *minVVp,
    /* IN */ unsigned32    Flags,
    /* OUT */ afsFid      *OutFidp,
    /* OUT */ afsFetchStatus *OutFidStatusp,
    /* OUT */ afsFetchStatus *OutDirStatusp,
    /* OUT */ afsToken      *OutTokenp,
    /* OUT */ afsVolSync    *Syncp
);

```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the symbolic link is to be created.
<i>Namep</i>	The name of the link to create.
<i>LinkContentsp</i>	The target of the new symbolic link.
<i>InStatusp</i>	This specifies the clientModTime field and unixModeBits of the new link.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier for the newly created symbolic link.
<i>OutFidStatusp</i>	The status information of the newly created symbolic link upon termination of the call.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the symbolic link was created) upon termination of the call.
<i>OutTokenp</i>	The token returned against the specified directory.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**DESCRIPTION**

This call is used to create a symbolic link.

No token is granted on the new symbolic link because symbolic links cannot change, they can only be deleted. It is recommended that the cache manager make use of this fact.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

- [DFS\_ENOENT] The entry named by *DirFidp* does not exist.
- [DFS\_EACCES] The caller lacks the required access rights.
- [DFS\_EINVAL] The entry is named “.” or “..”, or the caller passed in some invalid argument.
- [DFS\_EEXIST] The directory already contains an element named *Namep*.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*] A token management failure occurred.

## NAME

AFS\_MakeMountPoint — Make an AFS mount point

## SYNOPSIS

```
error_status_t AFS_MakeMountPoint(
    /* IN */ handle_t    h,
    /* IN */ afsFid      *DirFidp,
    /* IN */ afsTaggedName *Namep,
    /* IN */ afsTaggedName *cellNamep,
    /* IN */ afsFStype    Type,
    /* IN */ afsTaggedName *volumeNamep,
    /* IN */ afsStoreStatus *InStatusp,
    /* IN */ afsHyper      *minVVp,
    /* IN */ unsigned32    Flags,
    /* OUT */ afsFid      *OutFidp,
    /* OUT */ afsFetchStatus *OutFidStatusp,
    /* OUT */ afsFetchStatus *OutDirStatusp,
    /* OUT */ afsVolSync   *Syncp
);
```

## ARGUMENTS

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the new mount point is to be created.
<i>Namep</i>	The name of the mount point.
<i>cellNamep</i>	The string name of the cell in which the mount point is being created.
<i>Type</i>	The type of the file system mount point.
<i>volumeNamep</i>	The name of the fileset to be mounted on the newly created mount point.
<i>InStatusp</i>	This specifies the clientModTime field and unixModeBits of the new mount point.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier for the newly created mount point.
<i>OutFidStatusp</i>	The status information of the newly created mount point upon termination of the call.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the mount point was created) upon termination of the call.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## DESCRIPTION

This function call is obsolete. This call was used to create a new mount point for the specified fileset.

Note that this call does *not* require that the underlying file system support a separate file type for this type of object; for instance, in previous versions of AFS, these semantics were supported by

using a symbolic link with certain special attributes. This interface simply isolates the client Cache Manager from the server's implementation of **DCE DFS** Ename mount points.

**RETURN VALUE**

This function always returns the error [DFS\_ESRCH].

**ERRORS**

This function fails if:

[DFS\_ESRCH] Operation not supported in this DFS version.

## NAME

AFS\_HardLink — Create a hard link

## SYNOPSIS

```
error_status_t AFS_HardLink(
    /* IN */ handle_t    h,
    /* IN */ afsFid     *DirFidp,
    /* IN */ afsTaggedName *Namep,
    /* IN */ afsFid     *ExistingFidp,
    /* IN */ afsHyper    *minVVp,
    /* IN */ unsigned32  Flags,
    /* OUT */ afsFetchStatus *OutFidStatusp,
    /* OUT */ afsFetchStatus *OutDirStatusp,
    /* OUT */ afsVolSync  *Syncp
);
```

## ARGUMENTS

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the file is currently located.
<i>Namep</i>	The name to use for the new hard link.
<i>ExistingFidp</i>	The file identifier specifying the file identifier of the file to which the hard link should be made.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidStatusp</i>	The status information of the newly created hard link upon termination of the call.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the hard link was created) upon termination of the call.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## DESCRIPTION

This call creates a hard link to the file specified by *existingFid* to the new name *nameP* in the specified directory. The file specified by *existingFid* must not be a directory and must be on the same DCE DFS file system fileset.

## RETURN VALUE

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

## ERRORS

This function fails if:

- [DFS\_ENOENT] The entry named by *DirFidp* does not exist.
- [DFS\_EACCES] The caller lacks the required access rights.
- [DFS\_EXDEV] The target is in a different fileset than the source directory.



- [DFS\_EEXIST] The target name *nameP* already exists.
- [DFS\_EINVAL] The entry is named “.” or “..”, or the caller passed in some invalid argument.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*] A token managment failure occurred.

## NAME

AFS\_MakeDir — Create a directory

## SYNOPSIS

```
error_status_t AFS_MakeDir(
    /* IN */ handle_t h,
    /* IN */ afsFid *DirFidp,
    /* IN */ afsTaggedName *Namep,
    /* IN */ afsStoreStatus *InStatusp,
    /* IN */ afsHyper *minVVp,
    /* IN */ unsigned32 Flags,
    /* OUT */ afsFid *OutFidp,
    /* OUT */ afsFetchStatus *OutFidStatusp,
    /* OUT */ afsFetchStatus *OutDirStatusp,
    /* OUT */ afsToken *OutTokenp,
    /* OUT */ afsVolSync *Syncp
);
```

## ARGUMENTS

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the new directory is to be created.
<i>Namep</i>	The name of the directory to create.
<i>InStatusp</i>	This specifies the new status information, including clientModTime field and unixModeBits, of the new directory point.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutFidp</i>	The file identifier for the newly created directory.
<i>OutFidStatusp</i>	The status information of the newly created directory upon termination of the call.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the directory was created) upon termination of the call.
<i>OutTokenp</i>	A new token granted against the newly created directory.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## DESCRIPTION

This call creates a new directory; a token is granted against the newly created directory.

It is the cache manager's responsibility to either create an entry locally in the directory specified by *DirFidp*, or to invalidate this directory's cache entry.

## RETURN VALUE

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

- [DFS\_ENOENT] The entry named by *DirFidp* does not exist.
- [DFS\_EACCES] The caller lacks the required access rights.
- [DFS\_EEXIST] The target name *nameP* already exists.
- [DFS\_EINVAL] The entry is named “.” or “..”, or the caller passed in some invalid argument.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*] A token management failure occurred.

## NAME

AFS\_RemoveDir — Remove a directory

## SYNOPSIS

```
error_status_t AFS_RemoveDir(
    /* IN */ handle_t h,
    /* IN */ afsFid *DirFidp,
    /* IN */ afsFidTaggedName *Namep,
    /* IN */ afsHyper *returnTokenIDp,
    /* IN */ afsHyper *minVVp,
    /* IN */ unsigned32 Flags,
    /* OUT */ afsFetchStatus *OutDirStatusp,
    /* OUT */ afsFid *OutFidp,
    /* OUT */ afsFetchStatus *OutDelStatusp,
    /* OUT */ afsVolSync *Syncp
);
```

## ARGUMENTS

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file identifier specifying the directory in which the directory to be deleted is located.
<i>Namep</i>	The complex name of the directory to delete.
<i>returnTokenIDp</i>	A token ID being returned, if any.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutDirStatusp</i>	The status information of the directory (the one in which the directory was deleted) upon termination of the call.
<i>OutFidp</i>	
<i>OutDelStatusp</i>	
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## DESCRIPTION

This call removes a directory from the file system.

The directory must be empty (containing only entries for '.' and '..'), otherwise this call will fail. It is the cache manager's responsibility to either create an entry locally in the directory specified by *DirFid*, or to invalidate this directory's cache entry.

## RETURN VALUE

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

## ERRORS

This function fails if:

[DFS\_ENOENT] The entry named by *DirFidp* does not exist, or an entry named *nameP* does not exist in this directory.

- [DFS\_EACCES] The caller lacks the required access rights.
- [DFS\_EINVAL] The entry is named “.” or “..”, or the caller passed in some invalid argument.
- [DFS\_ENOTEMPTY]  
The directory is not empty.
- [DFS\_EBUSY] On some server systems, a process has this directory as its working directory.
- [FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.
- [TKM\_ERROR\*] A token management failure occurred.

**NAME**

AFS\_Readdir — Read an entry from a directory

**SYNOPSIS**

```
error_status_t AFS_Readdir(
    /* IN */ handle_t h,
    /* IN */ afsFid *DirFidp,
    /* IN */ afsHyper *Offsetp,
    /* IN */ unsigned32 Size,
    /* IN */ afsHyper *minVvp,
    /* IN */ unsigned32 Flags,
    /* OUT */ afsHyper *NextOffsetp,
    /* OUT */ afsFetchStatus *OutDirStatusp,
    /* OUT */ afsToken *OutTokenp,
    /* OUT */ afsVolSync *Syncp,
    /* OUT */ pipe_t *dirStream
);
```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>DirFidp</i>	The file descriptor for the specified directory.
<i>Offsetp</i>	The offset into the directory for this entry.
<i>Size</i>	Number of bytes to read.
<i>minVvp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>NextOffsetp</i>	The offset into the directory for the following entry.
<i>OutDirStatusp</i>	Status information for the directory pointed to by <i>DirFidp</i> .
<i>OutTokenp</i>	Token granted against <i>DirFidp</i> .
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.
<i>dirStream</i>	The array of bytes making up the external representation of this part of the directory.

**DESCRIPTION**

This call reads the specified directory entry, returning the requested entry and information to obtain the following entry.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

- [DFS\_EINVAL] The caller passed in some invalid argument.
- [DFS\_ENOENT] The entry named by *Fidp* does not exist.

[DFS\_EACCES] The caller lacks the required access rights.

[DFS\_ENOTDIR] The entry is not a directory.

[FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.

[TKM\_ERROR\*] A token management failure occurred.

## NAME

AFS\_GetToken — Obtain a token

## SYNOPSIS

```
error_status_t AFS_GetToken(
    /* IN */ handle_t    h,
    /* IN */ afsFid      *Fidp,
    /* IN */ afsToken    *MinTokenp,
    /* IN */ afsHyper    *minVVp,
    /* IN */ unsigned32  Flags,
    /* OUT */ afsToken    *OutTokenp,
    /* OUT */ afsRecordLock *OutBlockerp,
    /* OUT */ afsFetchStatus *OutStatusp,
    /* OUT */ afsVolSync  *Syncp
);
```

## ARGUMENTS

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	File identifier of the file to obtain a token against.
<i>MinTokenp</i>	Specification of the minimum requested token.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>OutTokenp</i>	Actual token granted.
<i>OutBlockerp</i>	Information about the possessor of the token that prevents the granting of the requested token (valid only for lock-family tokens).
<i>OutStatusp</i>	Status information on the file specified by <i>Fidp</i> .
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

## DESCRIPTION

This routine is called to obtain a token from the file exporter's token management layer.

## RETURN VALUE

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

## ERRORS

This function fails if:

[DFS_ENOENT]	The entry named by <i>Fidp</i> does not exist.
[DFS_EACCES]	The caller lacks the required access rights.
[TKM_*]	The token manager failed to grant the token.
[DFS_EINVAL]	The file specified by <i>Fidp</i> does not exist, or the caller passed in some invalid argument.
[FSHS_ERR*]	File server does not recognize client, or is in some form of recovery.



[TKM\_ERROR\*] A token managment failure occurred.

## NAME

AFS\_GetStatistics — Obtain file server statistics

## SYNOPSIS

```
error_status_t AFS_GetStatistics(  
    /* IN */ handle_t h,  
    /* OUT */ afsStatistics *Statisticsp  
);
```

## ARGUMENTS

*h*                    The RPC binding handle.  
*Statisticsp*        File server statistics.

## DESCRIPTION

This call returns statistics concerning file server throughput, resource usage and disk storage at the file server to which the call is directed. It is used for status monitoring purposes only. It always returns the value 0 indicating that the call was successful.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function always returns success.

## ERRORS

None.

This function succeeds if:

[error\_status\_ok] This function always returns success.

**NAME**

AFS\_ReleaseTokens — Release tokens granted by file server

**SYNOPSIS**

```
error_status_t AFS_ReleaseTokens(  
    /* IN */ handle_t h,  
    /* IN */ afsReturns *Tokens_Arrayp,  
    /* IN */ unsigned32 Flags  
);
```

**ARGUMENTS**

*h* The RPC binding handle.

*Tokens\_Arrayp* Tokens granted to this file server which are to be relinquished.

*Flags* Any of the AFS\_FLAG\_\* flags defined in the <afs4int.h> file.

**DESCRIPTION**

This call releases an array of tokens to a given file server. The tokens must have been granted originally by the file server to which this call is directed. The cache manager should not attempt to give up more than {MAX\_TOKEN\_RELEASE} tokens in any one call to a file server.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

[DFS\_EINVAL] Caller has not registered a token revocation service.

[FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.

[TKM\_ERROR\*] A token management failure occurred.

## NAME

AFS\_GetTime — Returns the time of day

## SYNOPSIS

```
error_status_t AFS_GetTime(  
    /* IN */ handle_t h,  
    /* OUT */ unsigned32 *Secondsp,  
    /* OUT */ unsigned32 *USecondsp,  
    /* OUT */ unsigned32 *SyncDistance,  
    /* OUT */ unsigned32 *SyncDispersion  
);
```

## ARGUMENTS

<i>h</i>	The RPC binding handle.
<i>Secondsp</i>	Number of seconds since January 1, 1970 UTC.
<i>USecondsp</i>	Number of microseconds into the current segment.
<i>SyncDistance</i>	Estimated path length to source of reliable time.
<i>SyncDispersion</i>	Measure of SyncDistance's variance.

## DESCRIPTION

This call returns the time of day in seconds and microseconds in the same format as returned by the `gettimeofday(system())` call

## RETURN VALUE

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

## ERRORS

This function fails if:

[FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.

**NAME**

AFS\_BulkFetchVV — Obtain many volume version numbers

**SYNOPSIS**

```

error_status_t AFS_BulkFetchVV(
    /* IN */ handle_t h,
    /* IN */ afsHyper *cellIdp,
    /* IN */ afsBulkVolIDs *VolIDsp,
    /* IN */ unsigned32 NumVols,
    /* IN */ unsigned32 Flags,
    /* IN */ unsigned32 spare1,
    /* IN */ unsigned32 spare2,
    /* OUT */ afsBulkVVs *VolVVsp,
    /* OUT */ unsigned32 *spare4
);

```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>cellIdp</i>	The cell identifier for the filesets whose Volume Version numbers are desired.
<i>VolIDsp</i>	The fileset identifiers (within <i>cellIdp</i> ) whose Volume Version numbers are desired.
<i>NumVols</i>	(Redundant) number of identifiers in <i>VolIDsp</i> .
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>spare1, spare2</i>	Spare IN parameters.
<i>VolVVsp</i>	The Volume Version numbers for the filesets whose identifiers were passed in.
<i>spare4</i>	Spare OUT parameter.

**DESCRIPTION**

This call is used to retrieve *en masse* the Volume Version number of a collection of filesets on the same cell. This allows for optimization for maintaining groups of filesets, replicated or otherwise, as might be done by a cache manager or a replication server.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

[DFS_ENOENT]	The fileset specified does not exist.
[DFS_EACCES]	The caller lacks the required access rights.
[DFS_EINVAL]	The value for <i>NumVols</i> is too small .
[FSHS_ERR*]	File server does not recognize client, or is in some form of recovery.

**NAME**

AFS\_ProcessQuota — Get or set quota information

**SYNOPSIS**

```
error_status_t  AFS_ProcessQuota(
    /* IN */ handle_t    h,
    /* IN */ afsFid      *Fidp,
    /* IN */ afsHyper    *minVVp,
    /* IN */ unsigned32  Flags,
    /* INOUT */ afsQuota  *quotaListp,
    /* OUT */ afsFetchStatus *OutStatusp,
    /* OUT */ afsVolSync  *Syncp
);
```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Fidp</i>	File ID of file or directory whose quota is being changed.
<i>minVVp</i>	The minimum-acceptable version number on the fileset containing this file or directory.
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>quotaListp</i>	The quota information being provided or obtained.
<i>OutStatusp</i>	The updated file status after the operation.
<i>Syncp</i>	The current synchronization information about the fileset containing this file; this allows for detection of changes in the fileset containing the specified file or directory.

**DESCRIPTION**

This function is obsolete. It does not perform any useful function. It was used to get or set quota information; the same call was used for both operations. The quota information is interpreted in a file system-specific manner. For example, the quota set with a UFS file system is a Berkeley quota, having per user components, while with a **DCE LFS** file system, the relevant quota is a **DCE LFS** per-fileset quota field.

**RETURN VALUE**

This function always fails and returns the error [DFS\_EINVAL]. A value of -1 is returned and *errno* is set to indicate the error.

**ERRORS**

This function fails if:

[DFS\_EINVAL] The quota structure is illegal. This function always returns this error.

**NAME**

AFS\_BulkKeepAlive — Keep multiple tokens active

**SYNOPSIS**

```

error_status_t AFS_BulkKeepAlive(
    /* IN */ handle_t h,
    /* IN */ afsBulkFEX *KAFEXp,
    /* IN */ unsigned32 numExecFids,
    /* IN */ unsigned32 Flags,
    /* IN */ unsigned32 spare1,
    /* IN */ unsigned32 spare2,
    /* OUT */ unsigned32 *spare4
);

```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>KAFEXp</i>	Collection of afsFid structures.
<i>numExecFids</i>	(Redundant) the number of afsFid structures in <i>KAFEXp</i> .
<i>Flags</i>	Any of the AFS_FLAG_* flags defined in the <afs4int.h> file.
<i>spare1, spare2</i>	Spare IN parameters.
<i>spare4</i>	Spare OUT parameter.

**DESCRIPTION**

This call is used to cause the file exporter to maintain open-for-preserve tokens on the indicated files until some characteristic interval passes. This call is expected to be useful in preventing storage for the lazy replicas of those files from being deleted before all cache managers are finished with it. The characteristic interval is a property of each read-only fileset.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

[DFS\_ENOENT] An entry in the *KAFEXp* structure represents a non-existent file.

**NAME**

AFS\_SetParams — Set cache manager parameters

**SYNOPSIS**

```
error_status_t AFS_SetParams(  
    /* IN */ handle_t h,  
    /* IN */ unsigned32 Flags,  
    /*INOUT*/ struct afsConnParams *paramsP  
);
```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Flags</i>	Either AFS_PARAM_TSR_COMPLETE or AFS_PARAM_RESET_CONN.
<i>paramsP</i>	A block of connection-oriented parameters.

**DESCRIPTION**

This call allows the Cache Manager to request that the File Server alter its connection parameters, in particular the lifetimes and timeouts associated with token state maintenance. Parameter suggestions are supplied by the Cache Manager as the input values in *\*paramsP*, and the values selected by the File Server in response are returned to the Cache Manager as the output values returned via *\*paramsP*.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a non-zero value will indicate the error.

**ERRORS**

This function fails if:

[FSHS\_ERR\*] File server does not recognize client, or is in some form of recovery.



**NAME**

AFS\_GetServerInterfaces — Obtain file server's supported interfaces

**SYNOPSIS**

```
error_status_t AFS_GetServerInterfaces(  
    /* IN */ handle_t h,  
    /*INOUT*/ dfs_interfaceList *serverInterfacesP  
);
```

**ARGUMENTS**

*h* The RPC binding handle.

*serverInterfacesP* Will contain the interface information when the call returns. Currently there is only one interface defined.

**DESCRIPTION**

This call returns information about the RPC interfaces that the file server is exporting.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success.

**ERRORS**

None.

This function succeeds if:

[error\_status\_ok] This function always returns success.



## *Cache Manager Service Interface*

The cache manager interface describes the interface between the client's cache manager and the server's file exporter (where the *server* initiates the operation). This section details the services provided by the cache manager; those provided by the file exporter (file server) are described in the topic, *DCE DFS File Exporter Interface*, found in Chapter 6 on page 99.

All of these calls are expected to return success ([*error\_status\_ok*]) unless noted otherwise.

**NAME**

TKN\_InitTokenState — Initialize token state

**SYNOPSIS**

```

error_status_t TKN_InitTokenState(
    /* IN */ handle_t h,
    /* IN */ unsigned32 Flags,
    /* IN */ unsigned32 hostLifeGuarantee,
    /* IN */ unsigned32 hostRPCGuarantee,
    /* IN */ unsigned32 deadServerTimeout,
    /* IN */ unsigned32 serverRestartEpoch,
    /* IN */ unsigned32 spare1,
    /* IN */ unsigned32 spare2,
    /* IN */ unsigned32 spare3,
    /* OUT */ unsigned32 *spare4,
    /* OUT */ unsigned32 *spare5,
    /* OUT */ unsigned32 *spare6
);

```

**ARGUMENTS**

*h* The RPC binding handle.

*Flags* Flags (TKN\_FLAG\*) indicating the server's likely token state.

*hostLifeGuarantee* Duration of the host lifetime on the server.

*hostRPCGuarantee* Duration of the promise of RPC attempts from the server.

*deadServerTimeout* Suggested maximum interval between probes from Cache Managers to unresponsive servers.

*serverRestartEpoch* Numeric tag for this run of the File Server, to allow Cache Managers to distinguish network partitions from server crashes.

*spare1 .. spare6* Spares.

**DESCRIPTION**

Initialize token state at the workstation, with respect to the calling server. All tokens from the calling server should be discarded or renegotiated.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success.

**ERRORS**

None.

This function succeeds if:

[error\_status\_ok] This function always returns success.

**NAME**

TKN\_Probe — Probe state of cache manager

**SYNOPSIS**

```
error_status_t TKN_Probe(  
    /* IN */ handle_t h  
);
```

**ARGUMENTS**

*h*                    The RPC binding handle.

**DESCRIPTION**

This call is a noop used by the server to check that a cache manager is still running. By periodically polling apparently inactive cache managers, a file server can timeout down workstations before actually trying to revoke a token at that workstation, allowing the operation eventually triggering a token revocation operation to complete faster.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success.

**ERRORS**

None.

This function succeeds if:

[error\_status\_ok] This function always returns success.

**NAME**

TKN-TokenRevoke — Request revocation of a previously granted token

**SYNOPSIS**

```
error_status_t TKN-TokenRevoke(  
    /* IN */ handle_t h,  
    /* INOUT */ afsRevokes *revokeDescp  
);
```

**ARGUMENTS**

*h*                    The RPC binding handle.  
*revokeDescp*        The tokens which are being revoked.

**DESCRIPTION**

This call is used by the file exporter requesting that a set of previously granted tokens be relinquished. If possible, the client's cache manager is to release these tokens.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. Otherwise, a value of zero is returned and the revoke descriptor is set to indicate exactly which tokens were not revoked. This function always returns success.

**ERRORS**

None.

This function succeeds if:

[error\_status\_ok] This function always returns success. Even when a token is not returned, this call returns success but indicates in the revoke descriptor exactly which tokens were not revoked.

**NAME**

TKN\_GetCellName — Request client cell name

**SYNOPSIS**

```
error_status_t TKN_GetCellName(  
    /* IN */ handle_t h,  
    /* IN */ afsHyper *Cellp,  
    /* OUT */ NameString_t CellNamep  
);
```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Cellp</i>	The ID for a cell.
<i>CellNamep</i>	The string name corresponding to <b>cellp</b> .

**DESCRIPTION**

This routine is used to perform cell id to cell name translation.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success.

**ERRORS**

None.

This function succeeds if:

[error\_status\_ok] This function always returns success.

**NAME**

TKN\_SetParams — Request client change parameters

**SYNOPSIS**

```
error_status_t TKN_SetParams(  
    /* IN */ handle_t h,  
    /* IN */ unsigned32 Flags,  
    /*INOUT*/ afsConnParams *paramsP  
);
```

**ARGUMENTS**

<i>h</i>	The RPC binding handle.
<i>Flags</i>	For future expansion.
<i>paramsP</i>	A block of connection-oriented parameters.

**DESCRIPTION**

This call allows the File Server to request that a cache manager alter its connection parameters, in particular the lifetimes and timeouts associated with token state maintenance. Parameter suggestions are supplied by the file server as the input values in *\*paramsP*, and manager in response are returned to the file server as the output values returned via *\*paramsP*.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success.

**ERRORS**

None.

This function succeeds if:

[error\_status\_ok] This function always returns success.



**NAME**

TKN\_GetServerInterfaces — Determine interfaces supported by client

**SYNOPSIS**

```
error_status_t TKN_GetServerInterfaces(  
    /* IN */ handle_t h,  
    /*INOUT*/ dfs_interfaceList *serverInterfacesP  
);
```

**ARGUMENTS**

*h* The RPC binding handle.  
*serverInterfacesP* List of supported interfaces.

**DESCRIPTION**

This call could be used by the file server to determine what interfaces a cache manager client supports. This would permit compatibility between older/newer versions of file exporters and cache managers.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success.

**ERRORS**

None.

This function succeeds if:

[error\_status\_ok] This function always returns success.

## NAME

TKN\_AsyncGrant — Notify client of token availability

## SYNOPSIS

```
error_status_t TKN_AsyncGrant(  
    /* IN */ handle_t    h,  
    /* IN */ afsFid     *grantedFileIDp,  
    /* IN */ afsToken   *grantedTokenP,  
    /* IN */ long       grantedRequestID  
);
```

## ARGUMENTS

*h* The RPC binding handle.

*grantedFileIDp* File identifier of the file or directory to which the granted token applies.

*grantedTokenP* The token that was granted.

*grantedRequestID* Identifier of the original async grant request.

## DESCRIPTION

This call allows the file server to notify a cache manager client that an async grant request for a token has succeeded.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function always returns success.

## ERRORS

None.

This function succeeds if:

[error\_status\_ok] This function always returns success.

# ***X/Open Preliminary Specification***

## **Part 5:**

### **Access Control Lists (ACLs) in DFS**

This part of the **DCE DFS** provides an overview of Access Control Lists (ACLs) in DFS, their formats and the interface provided to manipulate them.



## Access Control List Overview

This chapter describes the use of DCE Access Control Lists (ACLs) with DFS. An ACL lists the authorized users of an object and their allowed accesses. In DFS, ACLs control access to filesystem objects (files and directories) and server processes. These filesystem objects are stored in the DCE Local File System (LFS) as filesets. A fileset as used in this document is a hierarchical grouping of files that is managed as a single unit or entity. These filesets can vary in size but are usually smaller than a standard UNIX disk partition. In DFS, multiple DCE LFS filesets can be stored on a single UNIX disk partition. A non-LFS partition (for example, a UNIX Partition) in DFS can support (contain) only one fileset.

The ACLs used in DCE LFS filesets and discussed in this document in *Part 5, Access Control Lists (ACLs) in DFS*, differ from those used in other DCE Components. Refer to the *Security Service* portion of the **OSF DCE Administration Guide — Core Components** for details concerning DCE ACLs if further information is desired. This chapter also discusses *groups*. Further information about them can also be found in the *Security Service* portion of the **OSF DCE Administration Guide — Core Components** as well. Information on DFS usage of *groups* can also be found in the **OSF DCE Administration Guide — Core Components for OSF DCE 1.1**.

Associated with most DFS server processes is an administrative list that defines the users and server machines (called principals) and also the groups that can perform operations that affect the (server) process. Furthermore, server processes on different machines can have different lists (or each process may use a copy of the same list). Additionally, processes can share the same administrative list even though they may be of different types. Refer to Chapter 4 of the **OSF DCE DFS Administration Guide and Reference for OSF DCE 1.1** for information on using administrative lists in DFS.

### 8.1 ACL Types

There are three types of ACLs. They are:

#### *Object ACL*

Controls access to the object itself. Both filesystem objects and administrative lists have this ACL.

#### *Initial Object Creation ACL*

Possessed only by a directory. It is the default ACL inherited by a file created in the directory.

#### *Initial Container Creation ACL*

Possessed only by a directory. It is the default ACL inherited by a sub-directory created in the directory.

The existence of any of these three types of ACLs is mutually exclusive of any other. They can exist independently of one another or may not exist at all. Or, a given directory can have all, some, or none of these ACLs. The **OSF DCE DFS Administration Guide and Reference** (for **OSF DCE 1.1**) discusses ACLs and their inheritance in detail and explains the correspondence between their existence or absence and their contents.

## 8.2 ACL Entry Types

An **ACL** consists of a number of entries, each of which typically specify the access allowed to a user, group or organization. Typically, an **ACL** defines the operations that a user or group (or organization) can perform on an object (file or directory) in the form of permissions granted to various users or groups. However, these permissions can also apply to an entity called a principal, such as a server, that might be invoked on behalf of a user in order to complete a function requested by the user.

There is the concept of a default cell that is discussed in this section. A default cell of an **ACL** is the cell with respect to which the entries for the **ACL** are defined. A *user* or *group* named as an entry type in an **ACL** is defined to be from the default cell unless the entry type explicitly names a different cell.

In this document, a local user is one whose local cell is the same as the default cell of an **ACL**. Similarly, a foreign user is one whose default cell is different from the default cell of an **ACL**. Also, a foreign cell is different from the default cell of the **ACL**. A user's local cell is the cell in whose registry database the user's principal and account are defined.

The various types of entries are as follows:

<i>user_obj</i>	Specifies access rights for the owner principal of the object.
<i>group_obj</i>	Specifies access rights for the member principals of the group that owns the object.
<i>user</i>	Specifies access rights for a specific principal in the default cell of the <b>ACL</b> .
<i>group</i>	Specifies access rights for the member principals of a specific group in the default cell of the <b>ACL</b> .
<i>other_obj</i>	Specifies access rights for all other principals in the default cell of the <b>ACL</b> .
<i>foreign_user</i>	Specifies access rights for a specific principal in a specific foreign cell.
<i>foreign_group</i>	Specifies access rights for the member principals of a specific group in a specific foreign cell.
<i>foreign_other</i>	Specifies access rights for all principals in a specific foreign cell who are not covered by the <i>foreign_user</i> and <i>foreign_group</i> entries.
<i>any_other</i>	Specifies access rights for all other principals in all other foreign cells and unauthenticated users.
<i>extended</i>	A provision for future types.

### 8.2.1 Rules Governing Entries for Filesystem ACLs

- A filesystem **ACL** must have the *user\_obj*, *group\_obj* and *other\_obj* entries. All others are optional.
- An **ACL** can only have 1 entry each of types *user\_obj*, *group\_obj*, *other\_obj* and *any\_other*.
- If an entry other than *user\_obj*, *group\_obj* or *other\_obj* exists on an **ACL**, then the optional entry, *mask\_obj*, defined in Section 8.2.2 on page 157 must also exist.
- An **ACL** that has an *authenticated* entry, discussed in Section 8.2.3 on page 157, cannot be modified (until and unless this entry is removed).
- The *user\_obj* entry must always explicitly retain the *c* permission in order that the owner not be denied access to an owned object.

### 8.2.2 Optional ACL Entry Mask Types for Authenticated Users

The following optional entries are masks that serve to limit the maximum access allowed. They apply to authenticated users. An authenticated user is one whose DCE identity has been verified by the DCE Security Service and whose credentials have not expired.

*mask\_obj* Specifies maximum access allowed to all principals except principals covered by *user\_obj* and *other\_obj*. Only the access found in both the entry and the *mask\_obj* are granted. In absence of this mask, the default action is to place no limitations at all on types of access allowed.

This optional entry applies only to **DCE LFS** objects. It does not apply to **ACLs** on objects associated with other DCE components.

#### Authenticated Access Determination

In order to determine access permitted for an object for a user that is authenticated, the **ACL** is evaluated according to the steps in Section 8.7 on page 161.

### 8.2.3 ACL Entry Types for Unauthenticated Users

An unauthenticated user is one that has logged in to a machine but has not logged into DCE. Thus, the user's DCE identity has not been verified by the DCE Security Service. Once authenticated, a user is provided a set of credentials by the DCE Security Service. If the DCE credentials of such a user have expired, the user is also considered to be an unauthenticated user.

DCE **ACLs** used with objects for other DCE components include an additional *unauthenticated* entry type that masks the access that is permitted for unauthenticated users. There are no entry types on an **ACL** for unauthenticated users for **DCE LFS** objects.

#### Unauthenticated Access Determination

For users that are unauthenticated, access permissions are determined thusly rather than by Section 8.7 on page 161:

- A. For access to an object in a **DCE LFS** fileset:
  - The user is assigned the identity **nobody** which is treated as an authenticated user from an unknown foreign cell<sup>1</sup>. The primary group of the identity **nobody** is the group **nogroup**.
  - The user permissions are associated with the *any\_other* entry<sup>2</sup>.

**Note:** Prior to DCE 1.1, it was possible for *unauthenticated* entries to exist on **ACLs** of **DCE LFS** objects; however, they were ignored when determining the access permissions allowed for the user on the objects. If an existing object contains such an entry, it should be removed<sup>3</sup>. For DCE 1.1 and newer versions, *unauthenticated* entries for an **ACL** continue to be ignored when determining the access permissions.

1. Since the cell is unknown, its default cell is also unknown. Thus, no entries pertaining to the unknown (foreign) cell can be created on any **ACL**.

2. If an *any\_other* entry is not present on the **ACL**, the user has no permissions. If unauthenticated users are to be prevented from acquiring permissions for an object, this can be accomplished by not having an *any\_other* entry on the object's **ACL**.

- B. For access to an object in a **non-LFS** fileset:
- The user is assigned the identity **nobody** which is treated as an authenticated user regardless of their cell<sup>4</sup>.
  - The user permissions are associated with the *other* UNIX mode bits<sup>5</sup>.

### 8.2.4 Simple and Complex Entry Types

*Simple* ACL entries refer to entries that do not need to store a principal id in them. These are the *user\_obj*, *group\_obj*, *other\_obj*, *any\_other* and *mask\_obj* entries. *Complex* entries refer to all other entry types and store a principal id in the entry.

## 8.3 Delegates and DCE LFS Objects

A user initially requesting some operation on an **DCE LFS** object is referred to as the initiator of the operation. Typically the request is made for direct access to such an object. There are instances when a server may perform the work involved in a request on behalf of the initiator, such as a request to print a file. In such instances, the server principal that has been delegated the operation on the initiator's behalf is referred to as the delegate. The delegate performs operations on an object on behalf of the initiator. The operation performed in this manner is referred to as a delegation operation.

For an operation that does not involve delegation, only the initiator needs to have the necessary access permissions to perform a requested operation. However, for a delegation operation, both the initiator and the delegate must have the access permissions necessary to perform the operation.

DCE ACLS permit the granting of access permissions to a principal that apply only when the principal is a delegate on behalf of the operation. The next section describes the ACL entry types that are provided for delegation operations.

- 
3. Moving or restoring a **DCE LFS** fileset to a fileserver that is running DCE 1.1 or a newer version will automatically result in having any *unauthenticated* ACL entries removed from all objects in the fileset.
  4. This includes all *foreign* users (authenticated and unauthenticated alike).
  5. For clarity, as previously stated, authenticated users from foreign cells are granted the permissions associated with their authenticated *foreign* identities when they access objects in **DCE LFS** filesets.



## 8.4 ACL Entry Types for Delegation

There are entries for delegation corresponding to most of the entries specified in that can be used for delegation with **DCE LFS** objects. Since **DCE LFS** does not examine the entry types associated with the following:

- *user\_obj*
- *group\_obj*
- *other\_obj*.

for delegation, entries for these types shown in Section 8.2 on page 156 are omitted from this list of entry types for delegation. The valid entry types for delegation are the following:

<i>user_delegate</i>	Specifies access rights for a specific principal from the default cell of the <b>ACL</b> acting as a delegate.
<i>group_delegate</i>	Specifies access rights for the members of a specific group from the default cell of the <b>ACL</b> acting as a delegate.
<i>foreign_user_delegate</i>	Specifies access rights for a specific principal from a specific foreign cell acting as a delegate.
<i>foreign_group_delegate</i>	Specifies access rights for the member principals of a specific group in a specific foreign cell.
<i>foreign_other_delegate</i>	Specifies access rights for all principals in a specific foreign cell who do not match any of the preceding delegation entries acting as delegates.
<i>any_other_delegate</i>	Specifies access rights for all other principals in all other foreign cells who do not match any of the preceding delegation entries acting as delegates.

### 8.4.1 ACL Access Types for Delegation

Each of the delegation **ACL** entry types in the previous section can grant any of the access permissions available for **DCE LFS** objects. Further, each permission has the same meaning for a delegation entry that it has for a non-delegation entry. Refer to Section 8.5 for the list of the permissions and their meaning.

## 8.5 Access Types

A principal may access an object in a number of ways which define the operations the principal may perform on that object. These types of access are frequently termed permissions. If the object is a directory, certain permissions which apply only to it further define the types of access permitted on the directory object. These are the *insert* and *delete* permissions shown below.

The types of access that a principal can be granted to an object are defined to be the following:

<i>read</i>	Rights grant the user the ability to read and/or stat the contents of the object in question. It gives the possessor the ability to read the names and file IDs stored within a directory, as well as the ability to stat the directory. The read right for a file gives the possessor the ability to <b>stat</b> the file as well as read the data stored within the file.
<i>write</i>	Rights allow the principal to update the contents of a file. This right, for directories, is split into two other rights, <i>insert</i> and <i>delete</i> rights.

<i>insert</i>	Rights allow the principal to add entries to a directory. This includes the right to perform <b>renames</b> where the target is located within the directory. Note that if the target already exists, one will also require <i>delete</i> rights to remove the target of the <b>rename</b> system call.
<i>delete</i>	Rights allow the principal to remove entries from a directory. This includes the right to <b>rename</b> an object from one directory, as well as the right to obliterate an entry by renaming another entry “on top of it”.
<i>execute/search</i>	Rights allow the principal to execute a program (for files), or examine individual directory entries (for directories). Examining individual entries differs from enumerating the contents of a directory in that the former allows the retrieval of a piece of information only if the entire correct string name of the relevant file is presented to the file system.
<i>control</i>	Rights allow the principal to perform those operations normally reserved to a file or directory’s owner. This includes executing the <b>chmod</b> and <b>utimes</b> system calls, as well as modifying the <b>ACL</b> on the object.

## 8.6 Interaction of Filesystem ACLs with UNIX Permission Bits

In the UNIX operating system, file system protection for files and directories (henceforth called file and directory objects in this document) is provided by a set of bits called mode bits. The bits govern the operations a user can perform, and are called *read*, *write*, and *execute* mode bits. These bits determine access permissions on file and directory objects for three classifications of users:

- the user who owns the object
- members of the group that owns the object
- all other users.

File and directory objects in DCE filesystems, frequently referred to as **DCE LFS** filesets, also have mode bits. In addition, the protection of such files and directories can be augmented with **DCE ACLs**.

**Note:** DCE ACLs are used only with objects in **DCE LFS** filesets. For most non-LFS filesets, UNIX mode bits are the only form of protection.

The interaction of filesystem **ACLs** with UNIX permission mode bits depends in part upon the characteristics of the parent directory, for both directories and files within the filesystem. It also depends upon whether the file or directory is created by a foreign user. The governance is as follows:

- If an object’s (file or directory) parent directory has an appropriate *Initial Creation ACL* (object or container), the object has inherited an *Object ACL* as its form of protection. The object has UNIX mode bits but these are also augmented by the inherited *Object ACL*.
- If an object’s parent directory does *not* have an appropriate *Initial Creation ACL*, the object relies entirely upon mode bits as its only form of protection.

There is an exception to this rule in that if a file or directory is created by a foreign user, an *Object ACL* is always created for protection even if the parent directory does not have the appropriate *Initial Creation ACL*.

In DFS, an object’s UNIX mode bits and its **ACL** permissions are synchronized. The UNIX *owner* permission bits always correspond to *user\_obj* **ACL** entry permissions. The UNIX *other* permission bits always correspond to *other\_obj* **ACL** entry permissions. The correspondence for

UNIX *group* permission bits depends upon the presence or absence of the *mask\_obj* ACL entry. If the ACL contains a *mask\_obj* entry, the UNIX group permission bits correspond to the *mask\_obj* entry permissions, else the UNIX group permission bits correspond to *group\_obj* ACL entry permissions.

Modifications to the simple ACL entries *user\_obj*, *group\_obj* and *other\_obj* result in modification of the corresponding UNIX permission bits and vice versa. Appendix A on page 221, *Mapping DFS ACLs to UNIX mode bits*, discusses in more detail the relationship between ACLs and UNIX mode bits and provides a number of examples demonstrating the synchronization provided by DFS as changes are made.

## 8.7 Access Check Algorithm

To determine the access allowed to a principal on an object, the entry in the object's ACL matching the principal is located. This determination is made by examining the object's ACL entries using the following sequence of steps. Each step is only executed if the previous step did not result in a match.

1. The principal id is compared with the id specified in the *user\_obj*, *user* and *foreign\_user* ACL entries in sequence. The first matching entry is used. If the matching entry is not *user\_obj*, the *mask\_obj* permissions mask is applied to the permissions allowed by the matching entry. Otherwise (the principal owns the object), the access is *not* filtered through the *mask\_obj* entry.
2. The principal id is matched against the *group\_obj* entry. If a match is found the *mask\_obj* entry is used, as specified earlier, to filter the access allowed.
3. The principal id is matched against all the *group* and *foreign\_group* entries. The sequence in which these entries are matched is not relevant. The set-union of permissions allowed by all matched entries, is filtered by both the *mask\_obj* entry (if it exists) to determine the access permitted.
4. If the principal is a member of the default cell, the permissions specified in the *other\_obj* entry determine the access allowed.  
**Note:** The permissions are *not* filtered through the *mask\_obj* entry.
5. If the principal is from a foreign cell, the principal is matched with the *foreign\_other* entries. If a match is found, the permissions allowed in the matching entry are filtered by the *mask\_obj* entry to determine the access allowed.
6. If the principal is from a foreign cell and if that cell does not have a *foreign\_other* entry, the permissions allowed by the *any\_other* entry are filtered through the *mask\_obj* to determine the access allowed.
7. If a match is not found in any of the previous steps, the principal matches no entry and is denied access to the object.

## 8.8 Access Check Algorithm for Delegation

A principal that is initiating an operation is granted access permission from the **ACL** entry types described in Section 8.2 on page 156. A principal acquires permissions from a delegation entry only when acting as a delegate. Additionally, for delegation, **DCE LFS** filters all access permissions to be granted via delegation entries through the *mask\_obj* entry of the **ACL**. The delegation entry types are described in Section 8.4 on page 159.

**Note:** All delegation entry types in an **ACL** are optional. No delegation is permitted for an **ACL** with no delegation entry types.

In the **DCE LFS**, an operation requested by a delegate is performed only if the initiator and the delegate have the access permissions necessary to perform the operation. An **DCE LFS** object may have more than one delegate that has been specified. In this instance, if the initiator or one of the delegates does not have the required access, **DCE LFS** will refuse to perform the request.

The following sequence is used in determining the access allowed to a principal that is acting as a delegate. The sequence involves both delegation and non-delegation **ACL** entry types. Each step is only executed if the previous step did not result in a match.

1. The delegate owns the object. The delegate is granted the access permission from the *user\_obj* entry.
2. If one of the following entries exists, listed in checking priority order, the delegate is granted access permissions from the first of these entries that the delegate matches:
  - *user*
  - *user\_delegate*
  - *foreign\_user*
  - *foreign\_user\_delegate*.
3. One or the other of the following is true for the delegate:
  - A. The delegate is a member of the group that owns the object. In this instance, the permissions are granted via the *group\_obj* entry.
  - B. The delegate is a member of a group for which one of the following entries listed in checking priority order, exist. In this instance, the permissions from *all* the entries that the delegate matches are granted to the delegate. The entries are:
    - *group*
    - *group\_delegate*
    - *foreign\_group*
    - *foreign\_group\_delegate*.
4. The delegate is from the default cell. In this instance, the delegate is granted the access permissions from the *other\_obj* entry.
5. The delegate is from a foreign cell for which either of the following entries listed in checking priority order, exist. In this instance, the delegate is granted the access permissions associated with the first of these entries that the delegate matches. They are:
  - *foreign\_other*
  - *foreign\_other\_delegate*.

6. The delegate is from a foreign cell and either of the following entries listed in checking priority order, exist. In this instance, the delegate is granted the access permissions associated with the first of these entries that the delegate matches. They are:
  - *any\_other*
  - *any\_other\_delegate*.
7. The delegate matches no entry. In this instance, the delegate is denied access to the **DCE LFS** object.

### 8.8.1 Delegation and non-DCE LFS Objects

Since non-LFS objects do not have **DCE ACLs**, the access permissions required for delegation involving a non-LFS object are based solely upon the identity and access permissions of the last delegate in the chain. This is to say, the last delegate in the chain acquires the access permissions on a non-LFS object via the *user*, *group* or *other* (UNIX) mode bits.

Similarly, for administrative lists, delegation is not considered when determining administrative privileges. Thus, the last delegate in the chain must be included in the appropriate administrative list in order to be able to perform a privileged operation.

Also, since delegation is first available with DCE 1.1, the following restrictions apply. It is not possible to do the following for a fileset having **ACLs** containing one or more objects having delegation entries.

- Move the fileset to:
- Restore the fileset to:
- Add a replication site for the fileset on:

a File Server that uses a version of DCE earlier than DCE1.1.

Refer to the **OSF DCE DFS Administration Guide and Reference for OSF DCE 1.1**, for information relative to restrictions for delegation. Chapter 3 of that document provides information of this nature. Chapter 4 of that document provides information about administration considerations.



# ACL Storage Format

This chapter describes the structured formats used for ACLs in memory. ACLs are stored in different formats in memory (or, in core) than on disk. In-core ACLs use the structured format described in this chapter, but disk ACLs are stored in an unstructured byte-stream format. In-core ACLs are packed into a byte-stream before being stored on disk and similarly when ACLs are read from disk, they are unpacked into the in-core format.

This in-core format can seemingly differ depending upon whether the ACL is being manipulated in the kernel or in user space. The next section, Section 9.1, shows how the differences are accounted for. The in-user-space ACLs are referred to as administrative list ACLs. Chapter 8 on page 155, *Access Control List Overview*, provides additional information about administrative lists and their usage in DCE DFS.

ACLs specify access permissions for the users and groups that can perform the specified types of operations (read, or write, for instance) on files and directories. Administrative lists specify the users and groups that can perform specified actions affecting specified server processes. For a DFS server process, groups can be specified on the administrative list associated with it. Then if certain users need access to the process, they are included in the group. Similarly, if a group of users needs access to more than one server process, they can be specified in the Registry Database as a group. Thus, they can be added or removed from the appropriate administrative list(s) as a group. A system administration team is a good example of such a group.

## 9.1 Principal ID Formats

Administrative lists always store principal ids in 128 bits format on disk and in memory. In filesystem ACLs, principal ids represent users or groups, which have only 32 significant bits. Hence, a school of thought proposed storing filesystem ACLs on disk using only 32 bits instead of 128 bits to prevent wasting disk space. Though this idea has been refuted, the DFS ACL interface does provide such flexibility and this point is important to understand the interface. In-core filesystem ACLs use the full 128 bits always.

The next section, Section 9.2, shows the two formats defined as either the full 128 bits or the 32 bits format. The two formats are manipulated by the functions in Chapter 10 on page 171, *Access Control List Interface Functions*.

## 9.2 Principal Identifier Format

In-core ACLs store principal ids in the full 128 bit UUID<sup>6</sup> format.

```
typedef struct epi_uuid {
    u_int32      longField1;
    u_int16      shortField1;
    u_int16      shortField2;
    unsigned char miscChars[8];
};
```

---

6. UUID is a universally unique identifier for an object.

```

} epi_uid_t;

#if defined(EPI_USE_FULL_ID)
typedef epi_uid_t epi_principal_id_t;
#else /* defined(EPI_USE_FULL_ID) */
typedef u_int32 epi_principal_id_t;
#endif /* defined(EPI_USE_FULL_ID) */

```

### 9.3 Foreign Cell Principal Identifier Format

A principal from a foreign cell is represented by type **epi\_uid\_foreign\_t**.

```

typedef struct epi_uid_foreign {
    epi_uid_t    id;
    epi_uid_t    realm;
} epi_uid_foreign_t;

```

### 9.4 Access Type Format

The access permissions are represented with a bitset type **dacl\_permset\_t**. The set also includes provision for future additional types.

```

/* the permission sets stored in an acl */

```

```

typedef u_int32 dacl_permset_t;

```

#### 9.4.1 Access Type Definitions

```

/* the various permissions that may be stored in such a set */
/* first the POSIX ones */
#define dacl_perm_read      0x01
#define dacl_perm_write    0x02
#define dacl_perm_execute  0x04
#define dacl_perm_control  0x08
#define dacl_perm_insert   0x10
#define dacl_perm_delete   0x20
#define DACL_USERLIST_PERMS (dacl_perm_read | dacl_perm_write | \
                             dacl_perm_control)
#define DACL_USERLIST_PERMSTRING "rwc"

#define ntoh_dacl_entry_permset_t(permsetP) *entryTypeP = ntohl(*permsetP)
#define hton_dacl_entry_permset_t(permsetP) *entryTypeP = htonl(*permsetP)

```



## 9.5 ACL Entry Type Format

The ACL entry types are specified by the type `dacl_entry_type_t`.

```
typedef u_int32 dacl_entry_type_t;
#define dacl_entry_type_user_obj      0
#define dacl_entry_type_group_obj     1
#define dacl_entry_type_other_obj     2
#define dacl_entry_type_user         3
#define dacl_entry_type_group         4
#define dacl_entry_type_mask_obj     5
/*
#define dacl_entry_type_class_group   5
#define dacl_entry_type_class_owner   6
#define dacl_entry_type_class_other   7
*/
#define dacl_entry_type_foreign_user  8
#define dacl_entry_type_foreign_group 9
#define dacl_entry_type_foreign_other 10
#define dacl_entry_type_unauth        11
#define dacl_entry_type_extended      12
#define dacl_entry_type_anyother      13

#define dacl_entry_type_user_obj_delegate 14
#define dacl_entry_type_group_obj_delegate 15
#define dacl_entry_type_other_obj_delegate 16
#define dacl_entry_type_user_delegate     17
#define dacl_entry_type_group_delegate    18
#define dacl_entry_type_foreign_other_delegate 19
#define dacl_entry_type_foreign_user_delegate 20
#define dacl_entry_type_foreign_group_delegate 21
#define dacl_entry_type_any_other_delegate 22

#define ntoh_dacl_entry_type_t(entryTypeP) *entryTypeP = ntohl(*entryTypeP)
#define hton_dacl_entry_type_t(entryTypeP) *entryTypeP = htonl(*entryTypeP)
```

**Note:** The entry type *unauth* is no longer valid, for DCE 1.1 and newer versions of DFS. See the note in Section 8.2.3 on page 157 of the chapter, *Access Control List Overview*.

## 9.6 Simple ACL Entry Format

Each of the simple ACL entries is represented by type **dacl\_simple\_entry\_t** defined as:

```
typedef struct dacl_simple_entry {
    u_int32      is_entry_good;
    dacl_permset_t perms; /* permissions allowed */
} dacl_simple_entry_t;
```

The simple entries are stored in the ACL as an array indexed by the entry type **dacl\_simple\_entry\_type\_t**.

```
/* the following enumerated type is used only in memory */
typedef enum dacl_simple_entry_type {
    dacl_simple_entry_type_userobj,
    dacl_simple_entry_type_groupobj,
    dacl_simple_entry_type_otherobj,
    dacl_simple_entry_type_maskobj,
    dacl_simple_entry_type_another,
    dacl_simple_entry_type_userobj_delegate,
    dacl_simple_entry_type_groupobj_delegate,
    dacl_simple_entry_type_otherobj_delegate,
    dacl_simple_entry_type_another_delegate,
    dacl_simple_entry_type_unauthmask /* this must remain as the last
                                        of the enums */
} dacl_simple_entry_type_t;
```

**Note:** The entry, *unauthmask*, must remain as the last of the enums. See Section 9.8 on page 169 for usage information.

## 9.7 Complex ACL Entry Format

Each of the complex ACL entries is represented by type **dacl\_complex\_entry\_t** defined as follows. This structure holds entries that have data associated with them.

```
typedef struct dacl_complex_entry {
    u_int32      num_entries;
    u_int32      entries_allocated;
    dacl_entry_t* complex_entries;
} dacl_complex_entry_t;

typedef struct dacl_entry {
    dacl_permset_t perms;
    dacl_entry_type_t entry_type;
    union {
        epi_uid_t id; /* for locally keyed entries */
        epi_uid_foreign_t foreign_id; /* for foreign keyed entries */
        dacl_extended_info_t extended_info; /* for uninterpreted data */
    } entry_info;
} dacl_entry_t;
```

The complex entries are stored in the ACL as an array indexed by the entry type **dacl\_complex\_entry\_type\_t**. The *user* and *group* entries include both *foreign* and *local* entries. The *other* entry includes both *foreign other* and *extended* entries.

```

typedef enum dacl_complex_entry_type {
    dacl_complex_entry_type_user, /* includes both user and foreign_user */
    dacl_complex_entry_type_group, /* includes both group and foreign_group */
    dacl_complex_entry_type_user_delegate,
    dacl_complex_entry_type_group_delegate,
    dacl_complex_entry_type_for_other_delegate,
    dacl_complex_entry_type_other /* includes both foreign_other
                                   & extended */
} dacl_complex_entry_type_t;

```

**Note:** The entry, *other*, must remain as the last of the enums. See Section 9.8 for usage information.

### 9.7.1 Extended Complex Entry Format

The extended complex ACL entry type `dacl_extended_info_t` is described as:

```

typedef struct dacl_format_label {
    u_int16    miscShorts[3];
    char      charField;
} dacl_format_label_t;

typedef struct dacl_extended_info {
    epi_uid_t    extensionType;
    dacl_format_label_t formatLabel;
    u_int32     numberBytes;
    char *      infoP;
} dacl_extended_info_t;

```

## 9.8 ACL Structure

The in-core ACL is described by the type `dacl_t`.

```

typedef struct dacl {
    epi_uid_t    mgr_type_tag; /* ACL Manager type */
    epi_uid_t    default_realm; /* Default cell UUID */
    u_int32     num_entries; /* the total number of entries
                              stored on disk */
    dacl_simple_entry_t simple_entries[MAX_SIMPLE_ENTRIES];
    dacl_complex_entry_t complex_entries[MAX_COMPLEX_ENTRIES];
} dacl_t;

```

**Note:** `MAX_SIMPLE_ENTRIES = ((int)dacl_simple_entry_type_unauthmask) + 1`  
`MAX_COMPLEX_ENTRIES = ((int)dacl_complex_entry_type_other) + 1`

## 9.9 The Structures for Reading Lists of `sec_id_t`

These are used to avoid copying **uuids** out of lists in the **pac**s. Their usage can be seen in Chapter 10 on page 171, *Access Control List Interface Functions*. A **pac** contains certified data that describes a principal's privilege attributes including principal name, groups of which the principal is a member and the native cell of the principal.

### 9.9.1 Define the `epi_sec_id` Structure

```
typedef struct epi_sec_id {
    epi_uuid_t    uuid;
    char *       name;
} epi_sec_id_t;
```

### 9.9.2 Define the `epi_sec_id_foreign` Structure

```
typedef struct epi_sec_id_foreign {
    epi_sec_id_t  id;
    epi_sec_id_t  realm;
} epi_sec_id_foreign_t;
```

## *ACL Interface Functions*

This chapter describes the functions available in the **ACL** interface by which the **ACLs** can be manipulated. The functions are used for both in-kernel filesystem **ACLs** and in-user-space administrative list **ACLs**.

Some functions have another interface specialized for filesystem **ACLs**. These specialized functions have a **dacl\_** prefix typically and are indicated in the interface description.

Appendix B on page 229 contains a list of the return values, along with their explanations, that can be returned by the functions in this chapter. Refer to it for those instances where the function lists a return value without a reason for it being returned.

## NAME

dacl\_CheckAccessId — perform authorization check

## SYNOPSIS

```
long dacl_CheckAccessId(
    dacl_t *          aclP,
    u_int32 *         permBitsP,
    dacl_permset_t * accessRequestedP,
    epi_uid_t *       realmIdP,
    epi_uid_t *       userIdP,
    epi_uid_t *       groupIdP,
    epi_sec_id_t *    groupIdListP,
    unsigned int      numGroups,
    epi_sec_id_foreign_t * foreignGroupIdListP,
    unsigned int      numForeignGroups,
    epi_principal_id_t * userObjP,
    epi_principal_id_t * groupObjP,
    int               isAuthn
);
```

## ARGUMENTS

<i>aclP</i>	Acl for the object to which access is desired.
<i>permBitsP</i>	UNIX permission bits for the object, if applicable.
<i>accessRequested</i>	Access desired.
<i>realmIdP</i>	Cell to which the specified principal belongs. A NULL value or all zero UID specifies the default cell (the cell to which the object belongs) .
<i>userIdP</i>	Principal desiring access to the object.
<i>groupIdP</i>	Primary group of principal .
<i>groupIdListP</i>	Supplementary local groups of which the principal is a member.
<i>numGroups</i>	Number of local groups.
<i>foreignGroupIdListP</i>	Supplementary groups in foreign cells of which the principal is a member.
<i>numForeignGroups</i>	Num of supplementary foreign groups.
<i>userObjP</i>	Principal who owns the object.
<i>groupObjP</i>	Group that owns the object.
<i>isAuthn</i>	Non-zero if the principal desiring access is authenticated.

## DESCRIPTION

This function performs the authorization check to determine if the access desired by the specified principal to the object, corresponding to the specified **ACL**, is allowed. If *permBitsP* is non-NULL, the UNIX permission bits are taken into account while performing the authorization check.

This function has a filesystem specific version called *dacl\_epi\_CheckAccessId()* with the same interface.

## RETURN VALUE

[**DACL\_ERROR\_ACCESS\_DENIED**], if the access desired is not available in the **ACL**,

[DACL\_ERROR\_ACCESS\_EXPLICITLY\_DENIED] if there was a user match but the access requested is not permitted.

[DACL\_ERROR\_MGR\_PARAMETER\_ERROR] if the manager type check found a null principal id or permission bits.

[DACL\_ERROR\_PARAMETER\_ERROR] if *accessRequestedP* was null.

Otherwise, if the access requested is granted, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_DetermineAccessAllowed()* which is invoked by this routine as part of its determination.

## NAME

dacl\_DetermineAccessAllowed — determine access allowed

## SYNOPSIS

```
long dacl_DetermineAccessAllowed(
    dacl_t *          aclP,
    u_int32 *         permBitsP,
    dacl_permset_t * accessRequestedP,
    epi_uid_t *       realmIdP,
    epi_uid_t *       userIdP,
    epi_uid_t *       groupIdP,
    epi_sec_id_t *    groupIdListP,
    unsigned int      numGroups,
    epi_sec_id_foreign_t * foreignGroupIdListP,
    unsigned int      numForeignGroups,
    epi_principal_id_t * userObjP,
    epi_principal_id_t * groupObjP,
    int               isAuthn,
    dacl_permset_t * accessAllowedP
);
```

## ARGUMENTS

<i>aclP</i>	Acl for the object to which access is desired.
<i>permBitsP</i>	UNIX permission bits for the object, if applicable.
<i>accessRequested</i>	Access desired.
<i>realmIdP</i>	Cell to which the specified principal belongs. A NULL value or all zero UID specifies the default cell (the cell to which the object belongs).
<i>userIdP</i>	Principal desiring access to the object.
<i>groupIdP</i>	Primary group of principal.
<i>groupIdListP</i>	Supplementary local groups of which the principal is a member.
<i>numGroups</i>	Number of local groups.
<i>foreignGroupIdListP</i>	Supplementary groups in foreign cells of which the principal is a member.
<i>numForeignGroups</i>	Num of supplementary foreign groups.
<i>userObjP</i>	Principal who owns the object.
<i>groupObjP</i>	Group that owns the object.
<i>isAuthn</i>	Non-zero if the principal desiring access is authenticated.
<i>accessAllowedP</i>	If non-NULL, set to access allowed.

## DESCRIPTION

This function determines the set of access allowed to the specified principal to the object corresponding to the specified **ACL**. If *permBitsP* is non-NULL, the UNIX permission bits are taken into account while determining the access allowed.

This function has a filesystem specific version called *dacl\_epi\_DetermineAccessAllowed()* with the same interface.



**RETURN VALUE**

[DACL\_ERROR\_ACCESS\_DENIED], if the access desired is not available in the **ACL**,

[DACL\_ERROR\_ACCESS\_EXPLICITLY\_DENIED] if there was a user match but the access requested is not permitted.

[DACL\_ERROR\_MGR\_PARAMETER\_ERROR] if the manager type check found a null principal id or permission bits.

[DACL\_ERROR\_PARAMETER\_ERROR] if *accessRequestedP* was null.

Otherwise, if the access requested is granted, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_CheckAccessId()* which invokes this routine as part of its determination.

**NAME**

dacl\_CheckAccessPac — perform authorization check

**SYNOPSIS**

```
long dacl_CheckAccessPac(
    dacl_t *          aclP,
    u_int32 *         permBitsP,
    dacl_permset_t * accessRequestedP,
    sec_id_pac_t *    pacP,
    epi_principal_id_t * userObjP,
    epi_principal_id_t * groupObjP
);
```

**ARGUMENTS**

<i>aclP</i>	<b>ACL</b> for the object to which access is desired.
<i>permBitsP</i>	UNIX permission bits for the object, if applicable.
<i>accessRequested</i>	Access desired.
<i>pacP</i>	The PAC <sup>7</sup> for the principal desiring access to the object .
<i>userObjP</i>	Principal who owns the object.
<i>groupObjP</i>	Group that owns the object.

**DESCRIPTION**

This function performs the authorization check to determine if the access desired by the principal, corresponding to the specified PAC, to the object, corresponding to the specified **ACL**, is allowed. If *permBitsP* is non-NULL, the UNIX permission bits are taken into account while performing the authorization check.

This function has a file system specific version called *dacl\_epi\_CheckAccessPac()* with the same interface through DCE 1.1. For versions of DCE newer than 1.1, the two interfaces are no longer the same.

**RETURN VALUE**

[DACL\_ERROR\_ACCESS\_DENIED], if the access desired is not available in the **ACL**,  
 Otherwise, if the access requested is granted, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_CheckAccessAllowedPac()* which is invoked by this routine as part of its determination.

---

7. A PAC is an acronym of *Privilege Attribute Certificate*. It contains certified data that describes a principal's privilege attributes including principal name, groups of which the principal is a member and the native cell of the principal.

**NAME**

dacl\_epi\_CheckAccessPac — perform authorization check

**SYNOPSIS**

```
long dacl_epi_CheckAccessPac(
    dacl_t *          aclP,
    u_int32 *         permBitsP,
    dacl_permset_t * accessRequestedP,
    sec_id_pac_t *    pacP,
    epi_principal_id_t * userObjP,
    epi_principal_id_t * groupObjP,
    int               networkRootCheck
);
```

**ARGUMENTS**

*aclP*                    **ACL** for the object to which access is desired.

*permBitsP*             UNIX permission bits for the object, if applicable.

*accessRequested*      Access desired.

*pacP*                    The PAC<sup>8</sup> for the principal desiring access to the object.

*userObjP*                Principal who owns the object.

*groupObjP*               Group that owns the object.

*networkRootCheck*      This is a new argument for versions of DCE greater than 1.1. If true, a check is to be made for the network root ID which has full access. If true, and *pacP* principal UUID is the network root ID, this principal has all rights and the **ACL** is ignored.

**DESCRIPTION**

This function performs the authorization check to determine if the access desired by the principal, corresponding to the specified PAC, to the object, corresponding to the specified **ACL**, is allowed. If *permBitsP* is non-NULL, the UNIX permission bits are taken into account while performing the authorization check.

This function has a non-file system specific version called *dacl\_CheckAccessPac()* with the same interface for versions of DCE prior to DCE 1.2. In DCE 1.2, the argument *networkRootCheck* has been added.

**RETURN VALUE**

[DACL\_ERROR\_ACCESS\_DENIED], if the access desired is not available in the **ACL**,  
Otherwise, if the access requested is granted, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_epi\_CheckAccessAllowedPac()* which is invoked by this routine as part of its determination.

---

8. A PAC is an acronym of *Privilege Attribute Certificate*. It contains certified data that describes a principal's privilege attributes including principal name, groups of which the principal is a member and the native cell of the principal.

## NAME

dacl\_CheckAccessAllowedPac — determine access allowed

## SYNOPSIS

```
long dacl_CheckAccessAllowedPac(
    dacl_t *          aclP,
    u_int32 *         permBitsP,
    dacl_permset_t * accessRequestedP,
    sec_id_pac_t *    pacP,
    epi_principal_id_t * userObjP,
    epi_principal_id_t * groupObjP,
    dacl_permset_t * accessAllowedP
);
```

## ARGUMENTS

*aclP*                    **ACL** for the object to which access is desired.

*permBitsP*            UNIX permission bits for the object, if applicable.

*accessRequested*    Access desired.

*pacP*                    The PAC for the principal desiring access to the object.

*userObjP*              Principal who owns the object.

*groupObjP*             Group that owns the object.

*accessAllowedP*      If non-NULL, set to access allowed.

## DESCRIPTION

This function determines the set of access allowed to the principal, corresponding to the specified PAC, to the object, corresponding to the specified **ACL**. If *permBitsP* is non-NULL, the UNIX permission bits are taken into account while determining the access allowed.

This function has a filesystem specific version called *dacl\_epi\_CheckAccessAllowedPac()* with the same interface if the version of DCE is not newer than 1.1. For DCE 1.2, the interfaces are no longer the same, as a new argument, *networkRootCheck*, has been added to *dacl\_epi\_CheckAccessAllowedPac()* after DCE 1.1.

## RETURN VALUE

[DACL\_ERROR\_ACCESS\_DENIED], if the access desired is not available in the **ACL**,  
 Otherwise, if the access requested is granted, [DACL\_SUCCESS] is returned.

## ERRORS

None.

## SEE ALSO

*dacl\_CheckAccessPac()* which invokes this routine as part of its determination.

**NAME**

dacl\_epi\_CheckAccessAllowedPac — determine access allowed

**SYNOPSIS**

```
long dacl_epi_CheckAccessAllowedPac(
    dacl_t *          aclP,
    u_int32 *         permBitsP,
    dacl_permset_t * accessRequestedP,
    sec_id_pac_t *    pacP,
    epi_principal_id_t * userObjP,
    epi_principal_id_t * groupObjP,
    int               networkRootCheck,
    dacl_permset_t * accessAllowedP
);
```

**ARGUMENTS**

*aclP*                    **ACL** for the object to which access is desired.

*permBitsP*            UNIX permission bits for the object, if applicable.

*accessRequested*    Access desired.

*pacP*                    The PAC for the principal desiring access to the object.

*userObjP*              Principal who owns the object.

*groupObjP*             Group that owns the object.

*networkRootCheck* This is a new argument for versions of DCE greater than 1.1. If true, a check is to be made for the network root ID which has full access. If true, and *pacP* principal UUID is the network root ID, this principal has all rights and the **ACL** is ignored.

*accessAllowedP*    If non-NULL, set to access allowed.

**DESCRIPTION**

This function determines the set of access allowed to the principal, corresponding to the specified PAC, to the object, corresponding to the specified **ACL**. If *permBitsP* is non-NULL, the UNIX permission bits are taken into account while determining the access allowed.

This function has a non-filesystem specific version called *dacl\_CheckAccessAllowedPac()* with the same interface through DCE 1.1. After that, the interfaces are not the same, as this function then checks for the network root ID and grants it certain priviledges.

**RETURN VALUE**

[DACL\_ERROR\_ACCESS\_DENIED], if the access desired is not available in the **ACL**,  
 Otherwise, if the access requested is granted, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_epi\_CheckAccessPac()* which invokes this routine as part of its determination.

**NAME**

dacl\_PacFromUcred — construct a PAC from UNIX credentials

**SYNOPSIS**

```
void dacl_PacFromUcred(  
    sec_id_pac_t *    pacP,  
    struct ucred *    ucredP  
);
```

**ARGUMENTS**

<i>pacP</i>	The PAC.
<i>ucredP</i>	The UNIX credentials.

**DESCRIPTION**

This function constructs a PAC corresponding to the specified UNIX credentials of a principal. If the local group pointer in the specified PAC is non-NULL, the memory associated with the local group list is reused to fill it with group information from the specified UNIX credentials. This assumes that the group list memory reused is big enough to hold all the groups from the ucred.

If the local group pointer is NULL, this routine will allocate a block of memory big enough to hold all the groups in the ucred. This memory must be released by the caller by calling `osi_Free(pacP->groups, pacP->num_groups)`.

The foreign group list of the *pacP* is preserved. If there are any foreign groups associated with the *pacP*, *pacP->foreign\_groups* and *pacP->num\_foreign\_groups* must be set appropriately. All other fields in the *pacP* will be modified, as appropriate.

**Note:** The caller must ensure that the group list is big enough and that the memory is released appropriately.

**RETURN VALUE**

None.

**ERRORS**

None.

**NAME**dacl\_FlattenAclWithModeBits — pack an **ACL****SYNOPSIS**

```
long dacl_FlattenAclWithModeBits(
    dacl_t *      aclP,
    char **      fileBufferPP,
    unsigned int * bytesInBufferP,
    u_int32      modeBits,
    int          useModeBits,
    int          flattenForDisk
);
```

**ARGUMENTS**

<i>aclP</i>	The in-core <b>ACL</b> .
<i>fileBufferPP</i>	The buffer to hold the packed <b>ACL</b> .
<i>bytesInBufferP</i>	Size of the buffer.
<i>modeBits</i>	UNIX permission bits for the object, if applicable.
<i>useModeBits</i>	Flag that whether the UNIX permission bits should be used.
<i>flattenForDisk</i>	Essentially specifies principal id format to be used in the packed <b>ACL</b> <sup>9</sup> .

**DESCRIPTION**

This function converts an in-core **ACL** to its byte-stream representation. If the *useModeBits* flag is set, the packed **ACL** is modified to correspond to the UNIX permission bits.

**RETURN VALUE**

[**DACL\_ERROR\_TOO\_FEW\_BYTES**], if *bytesInBufferP* does not contain enough space. In this case, no processing is done.

[**DACL\_ERROR\_PARAMETER\_ERROR**], if *aclP* or *fileBufferPP* or *bytesInBufferP* is not valid (null). In this instance the appropriate error(s) is(are) logged.

Otherwise, [**DACL\_SUCCESS**] is returned to indicate a successful conversion.

**ERRORS**

None.

**SEE ALSO**

This function also uses the following routines — *Epi\_PrinId\_FromUuid()*, *hton\_epi\_principal\_id()*, *hton\_epi\_uuid()*.

---

9. This *flattenForDisk* parameter provides the support for storing principal ids in the 32 bits format for filesystem **ACLs**. When the corresponding argument is zero, the flattened **ACL** always has 128 bit principal ids. If the argument is non-zero, then the code performs a runtime check to determine the size of principal ids in the flat **ACL**. If other appropriate flags have been defined, this would result in 32 bits principal ids being used in the flattened **ACL**. However, it was decided not to use 32 bits format and currently, the runtime check should always result in 128 bits format being used. Also note, that in-core **ACLs** always store principal ids in the 128 bits format.

**NAME**

`dacl_FlattenAcl` — pack an **ACL**

**SYNOPSIS**

```
long dacl_FlattenAcl(
    dacl_t *      aclP,
    char **       fileBufferPP,
    unsigned int * bytesInBufferP
);
```

**ARGUMENTS**

*aclP*                The in-core **ACL**.

*fileBufferPP*        The buffer to hold the packed **ACL**.

*bytesInBufferP*      Size of the buffer.

**DESCRIPTION**

This function converts an in-core **ACL** to its byte-stream representation. There are no UNIX permission bits considered in This function. The size of principal ids is determined at run time. This is usually used when a filesystem **ACL** is flattened to be stored to disk.

This function uses *dacl\_FlattenAclWithModeBits()* with *modeBits* set to `DACL_NO_MODE_BITS` (not applicable), *useModeBits* set to 0 (don't use mode bits), and *flattenForDisk* set to 1 (flatten for disk).

This function has a filesystem specific version called *dacl\_epi\_FlattenAcl* with a different interface described below.

**RETURN VALUE**

[`DACL_ERROR_TOO_FEW_BYTES`], if *bytesInBufferP* does not contain enough space. In this case, no processing is done.

[`DACL_ERROR_PARAMETER_ERROR`], if *aclP* or *fileBufferPP* or *bytesInBufferP* is not valid (null). In this instance the appropriate error(s) is(are) logged.

Otherwise, [`DACL_SUCCESS`] is returned to indicate a successful conversion.

**ERRORS**

None.

**SEE ALSO**

*dacl\_FlattenAclWithModeBits()*. This function also uses the following routines — *Epi PrinId FromUuid()*, *hton\_epi\_principal\_id()*, *hton\_epi\_uuid()*.



**NAME**

dacl\_epi\_FlattenAcl — pack an ACL

**SYNOPSIS**

```
long dacl_epi_FlattenAcl(
    dacl_t *      aclP,
    char **      fileBufferPP,
    unsigned int * bytesInBufferP,
    u_int32      modeBits
);
```

**ARGUMENTS**

*aclP*                    The in-core ACL.

*fileBufferPP*           The buffer to hold the packed ACL.

*bytesInBufferP*        Size of the buffer.

*modeBits*                UNIX permission bits for the object, if applicable.

**DESCRIPTION**

This function converts an in-core ACL to its byte-stream representation. The packed ACL is modified to correspond to the specified UNIX permission bits. The 128 bit principal id format is used in the flattened ACL.

This function uses *dacl\_FlattenAclWithModeBits()* with *modeBits* set to the applicable mode bits, *useModeBits* set to 1 (use mode bits), and *flattenForDisk* set to 0 (not for disk).

**RETURN VALUE**

[DACL\_ERROR\_TOO\_FEW\_BYTES], if *bytesInBufferP* does not contain enough space. In this case, no processing is done.

[DACL\_ERROR\_PARAMETER\_ERROR], if *aclP* or *fileBufferPP* or *bytesInBufferP* is not valid (null). In this instance the appropriate error(s) is(are) logged.

Otherwise, [DACL\_SUCCESS] is returned to indicate a successful conversion.

**ERRORS**

None.

**SEE ALSO**

*dacl\_FlattenAclWithModeBits()*. This function also uses the following routines — *Epi PrinId FromUuid()*, *hton\_epi\_principal\_id()*, *hton\_epi\_uuid()*.

## NAME

dacl\_ParseAclDiskOption — unpack a flattened ACL

## SYNOPSIS

```
long dacl_ParseAclDiskOption(
    char *      fileBufferP,
    int         bytesInBuffer,
    dacl_t *    aclBufferP,
    epi_uuid_t * mgrUuidP,
    int         parseFromDisk
);
```

## ARGUMENTS

*fileBufferPP*      The buffer holding the packed ACL.

*bytesInBufferP*    Size of the packed ACL.

*aclBufferP*        The in-core ACL.

*mgrUuidP*         The UUID for the ACL manager.

*parseFromDisk*    Essentially specifies principal id format in the specified packed ACL<sup>10</sup>

## DESCRIPTION

This function parses the given packed byte-stream format ACL and constructs an in-core ACL.

It uses the following functions in the parse operation - *ntoh\_epi\_uuid()* and *Epi\_PrinId\_toUuid()*. If a parameter error is found, it is logged.

## RETURN VALUE

[DACL\_ERROR\_BUFFER\_ALLOCATION], [DACL\_ERROR\_INCORRECT\_MGR\_UUID],  
 [DACL\_ERROR\_TOO\_FEW\_BYTES], [DACL\_ERROR\_TOO\_MANY\_BYTES],  
 [DACL\_ERROR\_DUPLICATE\_ENTRY\_FOUND],  
 [DACL\_ERROR\_UNRECOGNIZED\_ENTRY\_TYPE],  
 [DACL\_ERROR\_ENTRY\_TYPE\_TOO\_LARGE],  
 [DACL\_ERROR\_REQUIRED\_ENTRY\_MISSING] or [DACL\_ERROR\_PARAMETER\_ERROR].

Otherwise, [DACL\_SUCCESS] is returned.

## ERRORS

None.

## SEE ALSO

*dacl\_FlattenAclWithModeBits()*. This function uses the following routines — *ntoh\_epi\_uuid()* and *Epi\_PrinId\_toUuid()*. Also, see Appendix B on page 229 for a list of return values and their meanings.

---

10. Refer to earlier footnote in *dacl\_FlattenAclWithModeBits()* about *flattenForDisk* parameter.

**NAME**

dacl\_ParseSyscallAcl — unpack a flattened ACL

**SYNOPSIS**

```
long dacl_ParseSyscallAcl(
    char *      fileBufferP,
    int         bytesInBuffer,
    dacl_t *    aclBufferP,
    epi_uid_t * mgrUuidP
);
```

**ARGUMENTS**

*fileBufferP*      The buffer holding the packed ACL.

*bytesInBufferP*    Size of the packed ACL.

*aclBufferP*        The in-core ACL.

*mgrUuidP*         The UUID for the ACL manager.

**DESCRIPTION**

This function parses the given packed byte-stream format ACL and constructs an in-core ACL. The packed ACL is assumed to store principal ids in the 128 bits format.

It uses the function, *dacl\_ParseAclDiskOption()* with *parseFromDisk* set to 0 (parse from the syscall buffer).

**RETURN VALUE**

[DACL\_ERROR\_BUFFER\_ALLOCATION], [DACL\_ERROR\_INCORRECT\_MGR\_UUID],  
 [DACL\_ERROR\_TOO\_FEW\_BYTES], [DACL\_ERROR\_TOO\_MANY\_BYTES],  
 [DACL\_ERROR\_DUPLICATE\_ENTRY\_FOUND],  
 [DACL\_ERROR\_UNRECOGNIZED\_ENTRY\_TYPE],  
 [DACL\_ERROR\_ENTRY\_TYPE\_TOO\_LARGE],  
 [DACL\_ERROR\_REQUIRED\_ENTRY\_MISSING], or [DACL\_ERROR\_PARAMETER\_ERROR].

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_FlattenAclWithModeBits()* and *dacl\_ParseAclDiskOption()*. Also, see Appendix B on page 229 for a list of return values and their meanings.

**NAME**

dacl\_ParseAcl — unpack a flattened **ACL**

**SYNOPSIS**

```
long dacl_ParseAcl(
    char *      fileBufferP,
    int         bytesInBuffer,
    dacl_t *    aclBufferP,
    epi_uuid_t * mgrUuidP
);
```

**ARGUMENTS**

*fileBufferP*      The buffer holding the packed **ACL**.

*bytesInBufferP*    Size of the packed **ACL**.

*aclBufferP*        The in-core **ACL**.

*mgrUuidP*         The UUID for the **ACL** manager.

**DESCRIPTION**

This function parses the given packed byte-stream format **ACL** and constructs an in-core **ACL**. The principal id format used in the packed **ACL** is determined at runtime.

It uses the function, *dacl\_ParseAclDiskOption()* with *parseFromDisk* set to 1 (parse from the disk buffer).

**RETURN VALUE**

[DACL\_ERROR\_BUFFER\_ALLOCATION], [DACL\_ERROR\_INCORRECT\_MGR\_UUID],  
[DACL\_ERROR\_TOO\_FEW\_BYTES], [DACL\_ERROR\_TOO\_MANY\_BYTES],  
[DACL\_ERROR\_DUPLICATE\_ENTRY\_FOUND],  
[DACL\_ERROR\_UNRECOGNIZED\_ENTRY\_TYPE],  
[DACL\_ERROR\_ENTRY\_TYPE\_TOO\_LARGE],  
[DACL\_ERROR\_REQUIRED\_ENTRY\_MISSING], or [DACL\_ERROR\_PARAMETER\_ERROR].

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_FlattenAclWithModeBits()* and *dacl\_ParseAclDiskOption()*. Also, see Appendix B on page 229 for a list of return values and their meanings.

**NAME**

dacl\_ExtractPermBits — extract UNIX permission bits from an **ACL**

**SYNOPSIS**

```
long dacl_ExtractPermBits(  
    dacl_t *    daclP,  
    u_int32 *   permBitsP  
);
```

**ARGUMENTS**

*daclP*                    The **ACL**.  
*permBitsP*                The UNIX permission bits corresponding to the **ACL**.

**DESCRIPTION**

This function determines the UNIX permission bits corresponding to the **ACL**. Refer to Section 8.6 on page 160, for the correspondence between UNIX permission bits and the **ACL** access rights.

**RETURN VALUE**

[DACL\_ERROR\_PARAMETER\_ERROR], if a parameter error is found, which is also logged.  
Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

The correspondence between UNIX permission bits and **ACL** access rights can be found in Section 8.6 on page 160 in the chapter, *Access Control List Overview*.

**NAME**

dacl\_ChmodAcl — synchronize **ACL** with UNIX permission bits

**SYNOPSIS**

```
long dacl_ChmodAcl(
    dacl_t * daclP,
    u_int32 permBits,
    int forDirectory
);
```

**ARGUMENTS**

*daclP*            The **ACL** .

*permBits*        The UNIX permission bits.

*forDirectory*    Flag to specify a directory object **ACL**.

**DESCRIPTION**

This function modifies the **ACL** entries to correspond to the specified UNIX permission bits. Refer to Section 8.6 on page 160 in the chapter, *Access Control List Overview*, for the correspondence between UNIX permission bits and the **ACL** access rights. If the *forDirectory* flag is set, UNIX *write* permission results in *write*, *insert* and *delete* access rights being granted in the corresponding **ACL** entry.

**RETURN VALUE**

[DACL\_ERROR\_PARAMETER\_ERROR], if a parameter error is found (no **ACL** is passed), which is also logged.

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**NAME**

dacl\_FreeAclEntries — free memory allocated for **ACL** entries

**SYNOPSIS**

```
long dacl_FreeAclEntries(  
    dacl_t *    theAclP  
);
```

**ARGUMENTS**

*theAclP*            The **ACL**.

**DESCRIPTION**

This function releases the memory allocated for the complex entries in the **ACL**.

**RETURN VALUE**

[DACL\_ERROR\_PARAMETER\_ERROR], if a parameter error is found (no **ACL** is passed), which is also logged.

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

## NAME

dacl\_PrintAclEntry — print an **ACL** entry

## SYNOPSIS

```
long dacl_PrintAclEntry(  
    FILE *          stream,  
    dacl_entry_t *  aclEntryP,  
    epi_uid_t *     defaultRealmUuidP,  
    int             local  
);
```

## ARGUMENTS

*stream*                Output file stream.

*aclEntryP*            The **ACL** entry.

*defaultRealmUuidP*    The native cell UUID.

*local*                Specifies whether ids should be interpreted locally or via DCE security registry.

## DESCRIPTION

This function outputs an ascii representation of the **ACL** entry to the specified file stream. The various principal/group/organization ids present in the **ACL** entry are translated into their symbolic representations. Translations are done either locally or via the security service depending on the *local* flag.

## RETURN VALUE

[DACL\_ERROR\_UNRECOGNIZED\_USER\_OR\_GROUP],  
[DACL\_ERROR\_SEC\_RGY\_PGO\_ERROR] or [DACL\_ERROR\_NONLOCAL\_ENTRY\_TYPE].  
Otherwise, [DACL\_SUCCESS] is returned.

## ERRORS

None.

## SEE ALSO

This function also calls *dacl\_NameAndTypeStringsFromEntry()*, *dacl\_EntryType\_ToString()*, and *dacl\_Permset\_ToString()*.



**NAME**

dacl\_PrintAcl — print an **ACL**

**SYNOPSIS**

```
long dacl_PrintAcl(
    FILE *      stream,
    dacl_t *    aclP,
    int         local
);
```

**ARGUMENTS**

<i>stream</i>	Output file stream.
<i>aclEntryP</i>	The <b>ACL</b> .
<i>local</i>	Specifies whether ids should be interpreted locally or via DCE security registry.

**DESCRIPTION**

This function outputs an ascii representation of the **ACL** to the specified file stream. The various principal/group/organization ids present in the **ACL** are translated into their symbolic representations. Translations are done either locally or via the security service depending on the *local* flag.

**RETURN VALUE**

[DACL\_ERROR\_SEC\_RGY\_PGO\_ERROR], if the cell name cannot be found.

[DACL\_ERROR\_UNRECOGNIZED\_USER\_OR\_GROUP] or  
[DACL\_ERROR\_NONLOCAL\_ENTRY\_TYPE].

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

This function also calls *dacl\_PrintAclEntry()*.

**NAME**

dacl\_WriteToDisk — store an **ACL** into a file

**SYNOPSIS**

```
long dacl_WriteToDisk(  
    dacl_t *  aclP,  
    char *    filenameP  
);
```

**ARGUMENTS**

*aclP*                The **ACL**.  
*filenameP*          The file.

**DESCRIPTION**

This constructs a packed version of the specified in-core **ACL** and stores it in the specified file. There are no UNIX permission bits considered in This function. The principal id format for the packed **ACL** is determined at run time.

If a buffer was allocated for packing and a failure occurs in processing, the buffer is freed.

**RETURN VALUE**

[DACL\_ERROR\_FS\_OPEN], [DACL\_ERROR\_FS\_WRITE], [DACL\_ERROR\_FS\_CLOSE],  
[DACL\_ERROR\_TOO\_FEW\_BYTES] or [DACL\_ERROR\_PARAMETER\_ERROR].

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

This function also calls *dacl\_FlattenAcl()* to flatten for disk, *dacl\_epi\_FlattenAcl()* to flatten into the caller's buffer, *dacl\_ValidateBuffer()* to validate the **ACL** for correctness.

**NAME**

dacl\_CreateAclOnDisk — construct an **ACL** and store it

**SYNOPSIS**

```
long dacl_CreateAclOnDisk(
    char *          filenameP,
    epi_uuid_t *    mgrUuidP,
    epi_uuid_t *    defaultRealmUuidP,
    dacl_permset_t * userObjPermsP,
    dacl_permset_t * groupObjPermsP,
    dacl_permset_t * otherObjPermsP,
    int             useEpisodeFile,
    long            epiSyscallFlags
);
```

**ARGUMENTS**

<i>filenameP</i>	The file for the <b>ACL</b> for non-filesystem <b>ACLs</b> ; otherwise the filesystem object for the <b>ACL</b> .
<i>mgrUuidP</i>	The UUID for the <b>ACL</b> manager.
<i>defaultRealmUuidP</i>	The native cell UUID.
<i>userObjPermsP</i>	The permission set for the <i>user_obj</i> entry.
<i>groupObjPermsP</i>	The permission set for the <i>group_obj</i> entry.
<i>otherObjPermsP</i>	The permission set for the <i>other_obj</i> entry.
<i>useEpisodeFile</i>	Filesystem <b>ACL</b> indicator flag.
<i>epiSyscallFlags</i>	Flags used in a system call if its a filesystem <b>ACL</b> .

**DESCRIPTION**

This function constructs an **ACL** having only *user\_obj*, *group\_obj* and *other\_obj* simple entries and converts it into a packed byte-stream. If *useEpisodeFile* is non-zero, it indicates that the **ACL** is a filesystem **ACL** and *filenameP* specifies the filesystem object for which the **ACL** is being set. Otherwise, its a non-filesystem **ACL** and *filenameP* specifies the file to store the packed **ACL**.

**RETURN VALUE**

[DACL\_ERROR\_ACL\_FILE\_EXISTS] or [DACL\_ERROR\_FS\_STAT].

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

This function also calls *dacl\_WriteToDisk()* if this is a non-filesystem **ACL**, *dacl\_epi\_FlattenAcl()* to flatten for filesystems if this is a filesystem object, *dacl\_FreeAclEntries*, to free memory allocated for **ACL** entries.

**NAME**

`dacl_ReadFromDisk` — read an **ACL** from disk and unpack it

**SYNOPSIS**

```
long dacl_ReadFromDisk(  
    dacl_t *      aclP,  
    char *        filenameP,  
    epi_uuid_t *  mgrUuidP  
);
```

**ARGUMENTS**

*aclP*                    The resultant in-core **ACL**.

*filenameP*              The file storing the packed **ACL**.

*mgrUuidP*                UUID for the **ACL** manager.

**DESCRIPTION**

This function reads the packed **ACL** from the specified file and unpacks it into the in-core **ACL** format. The format of principal id in the packed **ACL** is determined at runtime.

**RETURN VALUE**

[`DACL_ERROR_FS_READ`], if the **ACL** could not be read from disk.

[`DACL_ERROR_BUFFER_ALLOCATION`], [`DACL_ERROR_FS_STAT`], if the **ACL** could not be found.

[`DACL_ERROR_FS_CLOSE`], if the file descriptor for the **ACL** could not be closed,

[`DACL_ERROR_ACLFILE_NOT_FOUND`], [`DACL_ERROR_FS_OPEN`], as well as those listed for `dacl_ParseAcl()`.

Otherwise, [`DACL_SUCCESS`] is returned.

**ERRORS**

None.

**SEE ALSO**

`dacl_ParseAcl()` which unpacks the **ACL** read from disk,

**NAME**

dacl\_AddEntryToAcl — add an entry to an **ACL**

**SYNOPSIS**

```
long dacl_AddEntryToAcl(
    dacl_t *      aclP,
    char *        typeStringP,
    char *        nameStringP,
    char *        permStringP,
    char *        dataStringP,
    epi_uuid_t *  mgrUuidP,
    int           local
);
```

**ARGUMENTS**

<i>aclP</i>	The <b>ACL</b> .
<i>typeStringP</i>	<b>ACL</b> entry type, in symbolic format.
<i>nameStringP</i>	Entry name, in symbolic format for complex <b>ACL</b> entries.
<i>permStringP</i>	<b>ACL</b> entry access permissions, in symbolic format.
<i>dataStringP</i>	Information for extended <b>ACL</b> entry type.
<i>mgrUuidP</i>	The <b>ACL</b> manager UUID.
<i>local</i>	Specifies whether the entry name should be interpreted locally or via DCE security registry.

**DESCRIPTION**

This function constructs an **ACL** entry with the specified characteristics and adds it to the specified **ACL**. The *nameStringP* parameter is only used for a complex **ACL** entry type and the *dataStringP* parameter only for an extended **ACL** entry type.

**RETURN VALUE**

[DACL\_ERROR\_ENTRY\_EXISTS], [DACL\_ERROR\_UNRECOGNIZED\_ENTRY\_TYPE] or [DACL\_ERROR\_ENTRY\_TYPE\_TOO\_LARGE].

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_InitAclEntryFromStrings()*.

## NAME

dacl\_ModifyAclEntry — modify an **ACL** entry

## SYNOPSIS

```
long dacl_ModifyAclEntry(  
    dacl_t * aclP,  
    char *   typeStringP,  
    char *   nameStringP,  
    char *   permStringP,  
    char *   dataStringP,  
    int      local  
);
```

## ARGUMENTS

<i>aclP</i>	The <b>ACL</b> .
<i>typeStringP</i>	<b>ACL</b> entry type, in symbolic format.
<i>nameStringP</i>	Entry name, in symbolic format for complex <b>ACL</b> entries.
<i>permStringP</i>	<b>ACL</b> entry access permissions, in symbolic format.
<i>dataStringP</i>	Information for extended <b>ACL</b> entry type.
<i>local</i>	Specifies whether the entry name should be interpreted locally or via DCE security registry.

## DESCRIPTION

This function modifies the specified **ACL** entry to have the specified characteristics.

## RETURN VALUE

[**DACL\_ERROR\_ENTRY\_NOT\_FOUND**], [**DACL\_ERROR\_UNRECOGNIZED\_ENTRY\_TYPE**], [**DACL\_ERROR\_ILLEGAL\_ENTRY**], [**DACL\_ERROR\_UNRECOGNIZED\_USER\_OR\_GROUP**], [**DACL\_ERROR\_NONLOCAL\_ENTRY\_TYPE**], [**DACL\_ERROR\_SEC\_RGY\_PGO\_ERROR**] or [**DACL\_ERROR\_MISSING\_NAME**].

Otherwise, [**DACL\_SUCCESS**] is returned.

## ERRORS

None.

## SEE ALSO

*dacl\_InitAclEntryFromStrings()*.

**NAME**

`dacl_DeleteAclEntry` — delete an **ACL** entry

**SYNOPSIS**

```
long dacl_DeleteAclEntry(
    dacl_t *      aclP,
    char *        typeStringP,
    char *        nameStringP,
    epi_uuid_t *  mgrUuidP,
    int           local
);
```

**ARGUMENTS**

<i>aclP</i>	The <b>ACL</b> .
<i>typeStringP</i>	<b>ACL</b> entry type, in symbolic format.
<i>nameStringP</i>	Entry name, in symbolic format for complex <b>ACL</b> entries.
<i>mgrUuidP</i>	The <b>ACL</b> manager UUID.
<i>local</i>	Specifies whether the entry name should be interpreted locally or via DCE security registry.

**DESCRIPTION**

This function deletes the specified entry from an **ACL**.

**RETURN VALUE**

[`DACL_ERROR_UNRECOGNIZED_ENTRY_TYPE`], [`DACL_ERROR_ILLEGAL_ENTRY`], [`DACL_ERROR_UNRECOGNIZED_USER_OR_GROUP`], [`DACL_ERROR_NONLOCAL_ENTRY_TYPE`], [`DACL_ERROR_SEC_RGY_PGO_ERROR`], [`DACL_ERROR_MISSING_NAME`], [`DACL_ERROR_ENTRY_NOT_FOUND`] or [`DACL_ERROR_ENTRY_REQUIRED`].

Otherwise, [`DACL_SUCCESS`] is returned.

**ERRORS**

None.

**SEE ALSO**

`dacl_InitAclEntryFromStrings()`.

## NAME

dacl\_DeleteAllEntries — invalidate all **ACL** entries

## SYNOPSIS

```
long dacl_DeleteAllEntries(  
    dacl_t *    aclP  
);
```

## ARGUMENTS

*aclP*            The **ACL**.

## DESCRIPTION

This function invalidates all entries in an **ACL** and releases any allocated memory for complex entries.

## RETURN VALUE

[DACL\_ERROR\_PARAMETER\_ERROR], if a parameter error is found (no **ACL** is passed), which is also logged.

Otherwise, [DACL\_SUCCESS] is returned.

## ERRORS

None.

## SEE ALSO

*dacl\_FreeAclEntries()*.



**NAME**

sec\_acl\_FlattenAcl — convert a DCE ACL into a flat DFS ACL

**SYNOPSIS**

```
long sec_acl_FlattenAcl(
    sec_acl_t * secAclP,
    char ** byteBufferPP,
    unsigned int * bytesInBufferP
);
```

**ARGUMENTS**

*secAclP*            The DCE ACL.

*byteBufferPP*      The buffer to hold the packed DFS ACL.

*bytesInBufferP*    Size of the buffer.

**DESCRIPTION**

This function converts a DCE ACL into a packed DFS ACL. No UNIX permission are used as DCE objects have no such construct. The 128 bits principal id format is used in the packed DFS ACL.

**RETURN VALUE**

[DACL\_ERROR\_BUFFER\_ALLOCATION], [DACL\_ERROR\_DUPLICATE\_ENTRY\_FOUND],  
[DACL\_ERROR\_INCORRECT\_MGR\_UUID] or [DACL\_ERROR\_PARAMETER\_ERROR].

The following can be returned from called functions. [DACL\_ERROR\_TOO\_FEW\_BYTES].

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_FlattenAclWithModeBits()* and *dacl\_FreeAclEntries()*.

**NAME**

`sec_acl_ParseAcl` — convert a flat DFS **ACL** into a DCE **ACL**

**SYNOPSIS**

```
long sec_acl_ParseAcl(  
    sec_acl_t * secAclP,  
    char *      byteBufferP,  
    unsigned int bytesInBuffer  
);
```

**ARGUMENTS**

*secAclP*            The DCE **ACL**.

*byteBufferPP*     The buffer holding the packed DFS **ACL**.

*bytesInBufferP*   Size of the packed DFS **ACL**.

**DESCRIPTION**

This function converts a packed DFS **ACL** into a DCE **ACL**. The packed DFS **ACL** is assumed to store the principal ids in the 128 bit format.

**RETURN VALUE**

[`DACL_ERROR_PARAMETER_ERROR`], if a parameter error is found (no **ACL** is passed), which is also logged.

Otherwise, [`DACL_SUCCESS`] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_FreeAclEntries()*.

**NAME**

`dacl_InitAclEntryFromStrings` — construct an **ACL** entry

**SYNOPSIS**

```
long dacl_InitAclEntryFromStrings(
    dacl_entry_t *   aclEntryP,
    epi_uuid_t *    defaultRealmUuidP,
    char *          typeStringP,
    char *          nameStringP,
    char *          permStringP,
    char *          dataStringP,
    epi_uuid_t *    mgrUuidP,
    int             local
);
```

**ARGUMENTS**

<i>aclEntryP</i>	The resultant <b>ACL</b> entry.
<i>defaultRealmUuidP</i>	The native cell UUID.
<i>typeStringP</i>	<b>ACL</b> entry type, in symbolic format.
<i>nameStringP</i>	Entry name, in symbolic format for complex <b>ACL</b> entries.
<i>permStringP</i>	<b>ACL</b> entry access permissions, in symbolic format.
<i>dataStringP</i>	Information for extended <b>ACL</b> entry type.
<i>mgrUuidP</i>	The <b>ACL</b> manager UUID.
<i>local</i>	Specifies whether the entry name should be interpreted locally or via DCE security registry.

**DESCRIPTION**

This function constructs an **ACL** entry with the specified characteristics.

**RETURN VALUE**

[`DACL_ERROR_UNRECOGNIZED_ENTRY_TYPE`], [`DACL_ERROR_ILLEGAL_ENTRY`], [`DACL_ERROR_UNRECOGNIZED_USER_OR_GROUP`], [`DACL_ERROR_NONLOCAL_ENTRY_TYPE`], [`DACL_ERROR_SEC_RGY_PGO_ERROR`] or [`DACL_ERROR_MISSING_NAME`].

Otherwise, [`DACL_SUCCESS`] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_EntryType\_FromString()*.

**NAME**

`dacl_NameAndTypeStringsFromEntry` — generate ascii representations of complex **ACL** entry name and type

**SYNOPSIS**

```
long dacl_NameAndTypeStringsFromEntry(  
    dacl_entry_t *   aclEntryP,  
    epi_uuid_t *    defaultRealmUuidP,  
    char **          typeStringPP,  
    sec_rgy_name_t  globalName,  
    int              local  
);
```

**ARGUMENTS**

<i>aclEntryP</i>	The <b>ACL</b> entry.
<i>defaultRealmUuidP</i>	The native cell UUID.
<i>typeStringP</i>	Resultant symbolic type of entry type.
<i>globalName</i>	Resultant symbolic name of entry principal.
<i>local</i>	Specifies whether the entry id should be interpreted locally or via DCE security registry.

**DESCRIPTION**

This function generates the ascii string representations of the name and type of the complex **ACL** entry. The various principal/group/organization ids present in the **ACL** entry are translated into their symbolic representations. Translations are done either locally or via the security service depending on the *local* flag.

**RETURN VALUE**

[`DACL_ERROR_UNRECOGNIZED_USER_OR_GROUP`],  
[`DACL_ERROR_SEC_RGY_PGO_ERROR`] or [`DACL_ERROR_NONLOCAL_ENTRY_TYPE`].  
Otherwise, [`DACL_SUCCESS`] is returned.

**ERRORS**

None.

**SEE ALSO**

Also, see Appendix B on page 229 for a list of return values and their meanings.

**NAME**

dacl\_EntryType\_ToString — generate symbolic representation of **ACL** entry type

**SYNOPSIS**

```
char * dacl_EntryType_ToString(  
    dacl_entry_type_t    probeType  
);
```

**ARGUMENTS**

*probeType*            The **ACL** entry type.

**DESCRIPTION**

This function generates an ascii string representation from the specified internal representation of the **ACL** entry type.

**RETURN VALUE**

None.

**ERRORS**

None.

## NAME

dacl\_EntryType\_FromString — generate internal representation of symbolic **ACL** entry type

## SYNOPSIS

```
long dacl_EntryType_FromString(  
    dacl_entry_type_t * probeTypeP,  
    char * probeString  
);
```

## ARGUMENTS

*probeType*            The **ACL** entry type.  
*probeStringP*        Symbolic representation of the **ACL** entry type.

## DESCRIPTION

This function generates an internal representation of an **ACL** entry type from the specified ascii string representation.

## RETURN VALUE

[DACL\_ERROR\_UNRECOGNIZED\_ENTRY\_TYPE] or [DACL\_SUCCESS].

## ERRORS

None.

**NAME**

dacl\_Permset\_ToString — generate a symbolic representation of ACL entry permissions

**SYNOPSIS**

```
void dacl_Permset_ToString(  
    dacl_permset_t * thePermSetP,  
    char *          stringBufferP  
);
```

**ARGUMENTS**

*thePermSetP*      The ACL entry permissions bitset.  
*stringBufferP*    The buffer to hold the symbolic representation.

**DESCRIPTION**

This function generates the symbolic string representation of the specified ACL entry permissions bitset.

**RETURN VALUE**

None.

**ERRORS**

None.

**NAME**

`dacl_Permset_FromString` — generate internal representation of symbolic **ACL** entry permissions

**SYNOPSIS**

```
void dacl_Permset_FromString(  
    dacl_permset_t * thePermSetP,  
    char *          permStringP  
);
```

**ARGUMENTS**

*thePermSetP*      The generated permission bitset.  
*permStringP*      Symbolic **ACL** entry permissions.

**DESCRIPTION**

This function generates the internal representation of **ACL** entry permissions, a bitset, from its symbolic ascii string representation.

**RETURN VALUE**

None.

**ERRORS**

None.



**NAME**

dacl\_ValidateBuffer — parse and validate an ACL

**SYNOPSIS**

```
long dacl_ValidateBuffer(
    char *          byteBufferP,
    unsigned int *  bytesInBufferP,
    epi_uuid_t *    mgrUuidP,
    int             makeMinorRepairs,
    dacl_t **       daclPP
);
```

**ARGUMENTS**

*byteBufferP*            The packed format of the ACL.

*bytesInBufferP*        Size of the packed ACL.

*mgrUuidP*              The UUID of the ACL manager.

*makeMinorRepairs*    Flag that causes correction of minor errors detected during validation.

*daclPP*                The parsed and validate ACL.

**DESCRIPTION**

This function unpacks the given packed ACL into its in-core format. The 128 bit principal ids format is assumed in the packed ACL. The resultant in-core ACL is then validated for correctness. If *makeMinorRepairs* argument is set, any minor errors detected are corrected.

This function has a filesystem specific version called *dacl\_epi\_ValidateBuffer()* with the same interface.

**RETURN VALUE**

[DACL\_ERROR\_PARAMETER\_ERROR], which is also logged, [DACL\_ERROR\_VALIDATION\_FAILURE], if the *user\_obj* entry does not have c permission, or if the *mask\_obj* entry is missing. It can also indicate that an entry specifies *foreign\_user* within the default realm.

The following values can be returned from the called functions, below:

[DACL\_ERROR\_BUFFER\_ALLOCATION], [DACL\_ERROR\_INCORRECT\_MGR\_UUID], [DACL\_ERROR\_TOO\_FEW\_BYTES], [DACL\_ERROR\_TOO\_MANY\_BYTES], [DACL\_ERROR\_DUPLICATE\_ENTRY\_FOUND], [DACL\_ERROR\_UNRECOGNIZED\_ENTRY\_TYPE], [DACL\_ERROR\_ENTRY\_TYPE\_TOO\_LARGE] or [DACL\_ERROR\_REQUIRED\_ENTRY\_MISSING].

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_ParseSyscallAcl()*, *dacl\_AreObjectEntriesRequired()* and *dacl\_FreeAclEntries()*.

**NAME**

dacl\_InitEpiAcl — initialize a filesystem **ACL**

**SYNOPSIS**

```
long dacl_InitEpiAcl(
    dacl_t *    theAclP,
    u_int32 *   modeP,
    int        forDirectory,
    epi_uid_t * aclRealmP
);
```

**ARGUMENTS**

<i>theAclP</i>	The <b>ACL</b> .
<i>modeP</i>	The UNIX permission bits.
<i>forDirectory</i>	Flag to specify a directory object <b>ACL</b> .
<i>aclRealmP</i>	Default realm of <b>ACL</b> (optional).

**DESCRIPTION**

This function initializes the **ACL** to be a filesystem **ACL** in the local cell. The **ACL** only has 3 simple object entries *user\_obj*, *group\_obj* and *other\_obj*. If the UNIX permission bits argument *modeP* is non-NULL, the simple object **ACL** entries permissions are set to have corresponding access rights, otherwise *user\_obj* is set to have only *control* access right and *group\_obj* and *other\_obj* to have none. If the *forDirectory* flag is set, UNIX *write* permission results in *write*, *insert* and *delete* access being granted in the corresponding **ACL** entry.

If the realm (*aclRealmP*) is not null, it is used as the default realm of the **ACL**. Otherwise, for the kernel, the realm of the **ACL** is set to the local cellID; and for a non-kernel user, the realm is set to zero.

**RETURN VALUE**

[DACL\_ERROR\_PARAMETER\_ERROR], if a parameter error is found (no **ACL** is passed), which is also logged.

Otherwise, [DACL\_SUCCESS] is returned.

**ERRORS**

None.

**SEE ALSO**

*dacl\_ChmodAcl()*.

**NAME**

dacl\_AreObjectEntriesRequired — determine if simple object entries are needed in the **ACL**

**SYNOPSIS**

```
int dacl_AreObjectEntriesRequired(
    epi_uuid_t *    mgrUuidP
);
```

**ARGUMENTS**

*mgrUuidP*            The UUID for the **ACL** manager.

**DESCRIPTION**

This function determines if the specified **ACL** requires the *user\_obj*, *group\_obj* and *other\_obj* simple entries.

**RETURN VALUE**

[1], if the query determines the requirement must be met.

Otherwise, [0] is returned.

**ERRORS**

None.

## **dacl\_AreObjectUidsRequiredOnAccessCheck( )** *ACL Interface Functions*

### **NAME**

`dacl_AreObjectUidsRequiredOnAccessCheck` — determine if object UUIDs are required for authorization check

### **SYNOPSIS**

```
int dacl_AreObjectUidsRequiredOnAccessCheck(  
    epi_uuid_t *    mgrUuidP  
);
```

### **ARGUMENTS**

*mgrUuidP*            The UUID for the **ACL** manager.

### **DESCRIPTION**

This function determines if the the UUIDs corresponding to *user\_obj* and *group\_obj* are required during the authorization check procedure.

### **RETURN VALUE**

[1], if the query determines the requirement must be met.

Otherwise, [0] is returned.

### **ERRORS**

None.

**NAME**

dacl\_ArePermBitsRequiredOnAccessCheck — determine if UNIX permission bits are required for authorization check

**SYNOPSIS**

```
int dacl_ArePermBitsRequiredOnAccessCheck(
    epi_uid_t *    mgrUuidP
);
```

**ARGUMENTS**

*mgrUuidP*            The UUID for the **ACL** manager.

**DESCRIPTION**

This function determines if UNIX permission bits are required during the authorization check procedure.

**RETURN VALUE**

[1], if the query determines the requirement must be met.

Otherwise, [0] is returned.

**ERRORS**

None.

**NAME**

dacl\_AclMgrName — determine symbolic name of **ACL** manager

**SYNOPSIS**

```
char * dacl_AclMgrName(  
    epi_uuid_t *    mgrUuidP  
);
```

**ARGUMENTS**

*mgrUuidP*            The UUID for the **ACL** manager.

**DESCRIPTION**

This function determines the symbolic name for the **ACL** manager specified by the given **ACL** manager UUID.

**RETURN VALUE**

The following symbolic names are returned: [epsiode], [bossserver], or [unknown].

**ERRORS**

None.

**NAME**

Epi\_PrinId\_ToUuid — generate an UUID from a principal id

**SYNOPSIS**

```
void Epi_PrinId_ToUuid(  
    epi_principal_id_t * epiPrinIdP,  
    epi_uuid_t *        epiUuidP  
);
```

**ARGUMENTS**

<i>epiPrinIdP</i>	The principal id.
<i>epiUuidP</i>	The UUID.

**DESCRIPTION**

This is a macro that generates an UUID corresponding to the specified principal id.

**RETURN VALUE**

None.

**ERRORS**

None.

## NAME

Epi\_PrinId\_FromUuid — generate principal id from an UUID

## SYNOPSIS

```
void Epi_PrinId_FromUuid(  
    epi_principal_id_t * epiPrinIdP,  
    epi_uuid_t * epiUuidP  
);
```

## ARGUMENTS

<i>epiPrintIdP</i>	The principal id.
<i>epiUuidP</i>	The UUID.

## DESCRIPTION

This is a macro that generates the principal id corresponding to the specified UUID.

## RETURN VALUE

None.

## ERRORS

None.



**NAME**

Epi\_PrinId\_Cmp — compare 2 principal ids or UUIDs

**SYNOPSIS**

```
int  Epi_PrinId_Cmp(
    (epi_principal_id_t|epi_uuid_t) *   idBlob1P,
    (epi_principal_id_t|epi_uuid_t) *   idBlob2P
);
```

**ARGUMENTS**

*idBlob1P*            The first principal id or UUID.  
*idBlob2P*            The second principal id or UUID.

**DESCRIPTION**

This is a macro that compares the 2 specified principal ids or the 2 specified UUIDs for equality. It returns zero if they are the same, else it returns a non-zero value.

**RETURN VALUE**

This routine returns zero if the comparison is the same, otherwise it returns a non-zero value.

**ERRORS**

None.

## NAME

hton\_epi\_uuid — convert an UUID from host to network byte format

## SYNOPSIS

```
void hton_epi_uuid(  
    epi_uuid_t *    epiUuidP  
);
```

## ARGUMENTS

*epiUuidP*            The UUID.

## DESCRIPTION

This function converts an UUID from the host byte order to the network byte order format.

## RETURN VALUE

None.

## ERRORS

None.

**NAME**

ntoh\_epi\_uuid — convert a UUID from network to host byte format

**SYNOPSIS**

```
void ntohs_epi_uuid(  
    epi_uuid_t *    epiUuidP  
);
```

**ARGUMENTS**

*epiUuidP*            The UUID.

**DESCRIPTION**

This function converts an UUID from the network byte order to the host byte order format.

**RETURN VALUE**

None.

**ERRORS**

None.

## NAME

hton\_epi\_principal\_id — convert a principal id from host to network byte format

## SYNOPSIS

```
void hton_epi_principal_id(  
    epi_principal_id_t *    epiPrinIdP  
);
```

## ARGUMENTS

*epiPrinIdP*      The principal id.

## DESCRIPTION

This function converts a principal id from host byte order to network byte order format.

## RETURN VALUE

None.

## ERRORS

None.

**NAME**

ntoh\_epi\_principal\_id — convert a principal id from network to host byte format

**SYNOPSIS**

```
void ntoh_epi_principal_id(  
    epi_principal_id_t *    epiPrinIdP  
);
```

**ARGUMENTS**

*epiPrinIdP*      The principal id.

**DESCRIPTION**

This function converts a principal id from network byte order to host byte order format.

**RETURN VALUE**

None.

**ERRORS**

None.



## Mapping DFS ACLs to UNIX Mode Bits

This appendix discusses the relationship between ACLs and UNIX mode bits and also the file creation mask and how it relates to ACLs. More detailed information on ACLs can be found under the topic, *Access Control List Overview*, in Chapter 8 on page 155, relating to the ACL types discussed in this appendix.

### A.1 Relationship Between ACLs and UNIX Mode Bits

The following summarizes the relationship between ACLs and mode bits:

- *owner* mode bits correspond to ACL *user\_obj*
- *group* mode bits correspond as follows:
  - *group* mode bits correspond to ACL *group\_obj* if no *mask\_obj* is present
  - UNIX *group* mode bits correspond to the *mask\_obj* if the ACL has a *mask\_obj*
- *other* mode bits correspond to ACL *other\_obj*.

This relationship will be demonstrated in the next set of examples. They start out with a basic set of characteristics and then are cumulatively changed. This should be kept in mind as the examples are examined.

#### A.1.1 ACLs and UNIX Mode Bits

##### Example A-1 ACLs and UNIX Mode Bits

In this example, file **foo** is owned by `rajesh@minyan.dce.transarc.com`, **foo**'s owning group is `transarc@transarc.dce.com`. The prompt has been shortened to consist of the text, `[test]`.

The `acl_edit` command is shown being used with the list option (`-l`) to display the file's ACL.

```
[test] pwd
/.../minyan.dce.transarc.com/fs/test
[test] acl_edit foo -l
# SEC_ACL for foo:
# Default cell = /.../minyan.dce.transarc.com
user_obj:rw-c--
group_obj:r-----
other_obj:r-----
[test] ls -l foo
-rw-r--r--  1 rajesh  transarc      0 Sep 15 11:51 foo
```

### A.1.2 Changing other Mode Bits

#### Example A-2 Changing *other* Mode Bits

This example starts with the results demonstrated in Example A-1 on page 221. It demonstrates that *other* mode bits correspond to ACL *other\_obj* according to the relationship stated in Section A.1 on page 221.

```
[test] chmod o+x foo
[test] ls -l foo
-rw-r--r-x  1 rajesh  transarc      0 Sep 15 11:51 foo
[test] acl_edit foo -l
# SEC_ACL for foo:
# Default cell = ../../minyan.dce.transarc.com
user_obj:rw-c--
group_obj:r-----
other_obj:r-x---
```

Notice *other\_obj* now has execute rights.

### A.1.3 Changing group Mode Bits on a File with ACL with No *mask\_obj*

#### Example A-3 Changing *group* Mode Bits on a File with ACL and No *mask\_obj*

According to the relationship stated in Section A.1 on page 221, this example demonstrates that *group* mode bits correspond to the ACL *group\_obj* entry since no *mask\_obj* is present in the object. The base characteristics of the ACL are as they ended in the previous example, Example A-2.

```
[test] chmod g+w foo
[test] ls -l foo
-rw-rw-r-x  1 rajesh  transarc      0 Sep 15 11:51 foo
[test] acl_edit foo -l
# SEC_ACL for foo:
# Default cell = ../../minyan.dce.transarc.com
user_obj:rw-c--
group_obj:rw----
other_obj:r-x---
```

Notice that *group\_obj* now has w rights.

### A.1.4 Changing the *user\_obj* Entry

#### Example A-4 Changing the *user\_obj* Entry

In this example, *owner* mode bits will exhibit the changed *user\_obj*, as they correspond to the ACL *user\_obj* according to Section A.1 on page 221.

```
[test] acl_edit foo -m user_obj:rwxc
[test] acl_edit foo -l
# SEC_ACL for foo:
# Default cell = ../../minyan.dce.transarc.com
user_obj:rwxc--
group_obj:rw----
other_obj:r-x---
[test] ls -l foo
```



```
-rwxrw-r-x  1 rajesh  transarc      0 Sep 15 11:51 foo
  ^
  ^
```

Notice *owner* (mode bits) show *x* rights in the `ls -l` output.

### A.1.5 Changing the `group_obj` Entry

#### Example A-5 Changing the `group_obj` Entry

In this example, since the `ACL` has no `mask_obj`, the change will affect the `group` mode bits. The starting setup for this example is from the ending of Example A-4 on page 222.

```
[test] acl_edit foo -m group_obj:rwx
[test] acl_edit foo -l
ls -l
# SEC_ACL for foo:
# Default cell = /.../minyan.dce.transarc.com
user_obj:rwx--
group_obj:rwx---
other_obj:r-x---
[test] ls -l foo
-rwxrwxr-x  1 rajesh transarc      0 Sep 15 11:51 foo
```

The `group_obj` entry change impacts `group` mode bits.

### A.1.6 Changing the `group_obj` Entry in an `ACL` with a `mask_obj` Entry

#### Example A-6 Changing the `group_obj` Entry in an `ACL` with a `mask_obj` Entry

At this point, we introduce a `mask_obj` entry for the `ACL`. Now let's look at what happens if the `ACL` has an `mask_obj` entry.

**Note:** A `mask_obj` entry in an `ACL` is required if it has entries other than `user_obj`, `group_obj` and `other_obj`.

According to the relationship stated in Section A.1 on page 221, we should expect that a change to the `group_obj` entry in the `ACL` object will not affect the UNIX `group` mode bits.

```
[test] acl_edit foo -m user:ashok:r
[test] acl_edit foo -l
# SEC_ACL for foo:
# Default cell = /.../minyan.dce.transarc.com
mask_obj:rwx---
user_obj:rwx--
user:ashok:r-----
group_obj:rwx---
other_obj:r-x---
[test] ls -l foo
-rwxrwxr-x  1 rajesh  transarc      0 Sep 15 11:51 foo
```

### A.1.7 The Results of `chmod` of the group Bits

#### Example A-7 The Results of `chmod` of the *group* Bits

Starting with where we left off in the previous example, (Example A-6 on page 223,) the following scenario demonstrates the results of `chmod` of the *group* bits. According to the relationship established in Section A.1 on page 221, we should expect to see a corresponding change to the *mask\_obj* in the ACL, and no effect upon the *group\_object* entry.

```
[test] chmod g-x foo
[test] ls -l foo
-rwxrw-r-x  1 rajesh transarc      0 Sep 15 11:51 foo
[test] acl_edit foo -l
# SEC_ACL for foo:
# Default cell = ../../minyan.dce.transarc.com
mask_obj:rw----
user_obj:rwxc--
user:ashok:r-----
group_obj:rwx---      #effective:rw----
other_obj:r-x---
```

Notice that changing the group mode bits affected the *mask\_obj* entry and not the *group\_obj* entry.

### A.1.8 The Results of `acl_edit` of the *mask\_obj* Bits

#### Example A-8 The Results of `acl_edit` of the *mask\_obj* Bits

This example starts with the setup from the previous example, Example A-8, as the base. The following scenario shows the result of modification of the *mask\_obj* entry.

According to Section A.1 on page 221, modifying the *mask\_obj* entry in the ACL object directly affects the UNIX *group* mode bits:

```
[test] acl_edit foo -m mask_obj:r
[test] acl_edit foo -l
# SEC_ACL for foo:
# Default cell = ../../minyan.dce.transarc.com
mask_obj:r-----
user_obj:rwxc--
user:ashok:r-----
group_obj:rwx---      #effective:r-----
other_obj:r-x---
[test] ls -l foo
-rwxr--r-x  1 rajesh  transarc      0 Sep 15 11:51 foo
```

Notice that changing the *mask\_obj* to `r` changed the group bits shown by `ls` to `r`.

### A.1.9 Changing the `group_obj` Entry for an ACL with a `mask_obj`

#### Example A-9 Changing the `group_obj` Entry for an ACL with a `mask_obj`

According to Section A.1 on page 221, for an ACL object with a `mask_obj`, changing the `group_obj` entry for an ACL with `mask_obj` should have no impact on the `group` mode bits. The following scenario demonstrates this.

This example's setup starts with the results of the previous example, Example A-8 on page 224.

```
[test] acl_edit foo -m group_obj:- -l
# SEC_ACL for foo:
# Default cell = ../../minyan.dce.transarc.com
mask_obj:r-----
user_obj:rwxc--
user:ashok:r-----
group_obj:-----
other_obj:r-x---
[test] ls -l foo
-rwxr--r-x  1 rajesh transarc      0 Sep 15 11:51 foo
```

Notice that there is no change in `group` mode bits.

## A.2 File Creation Mask and ACLs

The following rules apply to the file creation mask. The file creation mask refers to the **umask**. In this appendix, an **Initial Object ACL** will be referred to as an **IO ACL**.

- A. In DFS if a directory has an appropriate **Initial Object ACL**, then the *umask* setting is not taken into account when creating a file in that directory. Only the mode bits specified to the **creat** syscall and the **IO ACL** determine the resulting **ACL** on the file.
- B. If the directory does not have an **IO ACL**, then the *umask* setting of the process creating the file is taken into account. In this case the mode bits on the resulting file are determined using the mode bits specified to **creat** syscall and the *umask*.

### A.2.1 File Creation with no Initial Object ACL

#### Example A-10 File Creation with no IO ACL

Let's say we are in a directory with no **IO ACL** called **simple\_dir** and we create a file using **touch**. In this example, **touch** specifies **0666** creation mode bits to the **creat** syscall.

First, let's set a *umask*, and then display it.

```
[simple_dir] umask 77
[simple_dir] umask
77
```

Now, let's *touch* the file, **foo** (and display it):

```
[simple_dir] touch foo
[simple_dir] ls -l foo
-rw-----  1 rajesh  transarc      0 Sep 15 12:08 foo
```

**touch foo** specified **creat(foo, 0666)**. DFS determined that the parent directory did not have an **IO ACL**, so case B applies.

The algorithm to determine the mode bits on the resulting file is the same as in traditional UNIX systems (as in POSIX.1):

```
resulting file mode bits = syscall creation mode bits & ~umask =
```

Thus, in this case:

```
resulting file mode bits = 0666 & ~077
                        = 0600
                        = owner:rw group:- other:-
```

as shown in the **ls -l** output.

### A.2.2 File Creation with an Initial Object ACL

#### Example A-11 File Creation with an IO ACL

Let's say we are in a directory with an **IO ACL** called **acl\_dir**. Case A should apply here on both file and directory creation. First, let's display the **Initial Object ACL**. This is done by using **acl\_edit** with the **-io -l** options:

```
[acl_dir] acl_edit . -io -l
# Initial SEC_ACL for objects created under: .:
# Default cell = ../../minyan.dce.transarc.com
```

```
mask_obj: rwx---
user_obj: rwxc--
user: ashok: rwx---
group_obj: rw----
other_obj: r-----
```

Let's set the **umask** and display it, as in the previous example:

```
[acl_dir] umask 77
[acl_dir] umask
77
```

Now, let's *touch* the file, **bar** (and display it):

```
[acl_dir] touch bar
[acl_dir] ls -l bar
-rw-rw-r-- 1 rajesh transarc 0 Sep 15 15:21 bar
```

Notice that the permissions **664** on the new file, **bar**, mean that the *umask* was not applied, else the *group* and *other* permissions would be zero.

### A.2.3 File ACL Creation with an initial object ACL

#### Example A-12 File ACL Creation with an initial object ACL

Let's see what the **ACL** for **bar** is. This can be done by using **acl\_edit** with the **-l** option, as follows:

```
[acl_dir] acl_edit bar -l
# SEC_ACL for bar:
# Default cell = ../../minyan.dce.transarc.com
mask_obj: rw----
user_obj: rw-c--
user: ashok: rwx--- #effective: rw----
group_obj: rw----
other_obj: r-----
```

Under this case (A), the permissions on the new file's **ACL** are derived by intersecting the parent directory's **IO ACL** permissions with the mode bits specified to the **creat** syscall, as follows.

**Note:** The parent directory's **IO ACL** permissions are those from the **ACL** shown in Example A-11 on page 226.

Also, recall that *touch* uses **creat** mode bits as **0666**. Thus, the respective *owner*, *group*, and *other* creation mode bits are **rw**.

```
user_obj perms = parent dir IO ACL user_obj perm &
                owner creation mode bits
```

Thus, in this case:

```
user_obj perms = rwx & rw
                = rw.
```

Also *user\_obj* is always granted **c** rights.

Similarly for *other\_obj*, intersect **IO ACL** *other\_obj* permissions with *other creat* mode bits:

```
other_obj perms = parent dir IO ACL other_obj perm &
                  other creation mode bits
```

Thus:

```
other_obj perms = r & rw = r
```

Recall from Section A.1 on page 221, that since the **IO ACL** has a *mask\_obj*, the *IO mask\_obj* will be intersected with the group **creat** mode bits to determine the *mask\_obj* on the file **ACL** as follows:

```
mask_obj perms = parent dir IO ACL mask_obj perm &
                  group creation mode bits
```

```
mask_obj perms = rwx & rw = rw.
```

**Note:** If the **IO ACL** did not have *mask\_obj*, then the group **creat** mode bits would be intersected with the IO ACL *group\_obj* permissions to determine the file **ACL**'s *group\_obj* permissions.

## Access Control List Package Error List

The following sections contain the return values that are returned from the *Access Control List Interface Functions* found in Chapter 10 on page 171.

### B.1 Access Control List Return Values

[DACL\_SUCCESS]

The operation was successful or no errors were found.

[DACL\_ERROR\_TOO\_MANY\_BYTES]

Too many bytes in acl byte buffer.

[DACL\_ERROR\_TOO\_FEW\_BYTES]

Too few bytes in acl byte buffer.

[DACL\_ERROR\_ENTRY\_NOT\_FOUND]

ACL entry not found.

[DACL\_ERROR\_ENTRY\_EXISTS]

ACL entry already exists.

[DACL\_ERROR\_UNRECOGNIZED\_ENTRY\_TYPE]

Unrecognized entry type.

[DACL\_ERROR\_UNRECOGNIZED\_ENTRY\_CLASS]

Unrecognized entry class.

[DACL\_ERROR\_UNIMPLEMENTED\_ENTRY\_TYPE]

Unimplemented entry type.

[DACL\_ERROR\_SEC\_RGY\_PGO\_ERROR]

Error returned from security service.

[DACL\_ERROR\_MISSING\_NAME]

Name required for specified entry type.

[DACL\_ERROR\_ENTRY\_REQUIRED]

Attempt to remove required ACL entry.

[DACL\_ERROR\_BUFFER\_ALLOCATION]

Buffer allocation error.

[DACL\_ERROR\_ACCESS\_DENIED]

Requested access (implicitly) denied by ACL.

[DACL\_ERROR\_ACCESS\_EXPLICITLY\_DENIED]

Requested access explicitly denied by ACL.

[DACL\_ERROR\_ACL\_FILE\_EXISTS]

File in which ACL creation was requested already exists.

[DACL\_ERROR\_INCORRECT\_MGR\_UUID]

Attempt to parse ACL by incorrect ACL manager.

[DACL\_ERROR\_ILLEGAL\_ENTRY]

ACL contains an entry not appropriate for the ACL manager type.

- [DACL\_ERROR\_MGR\_PARAMETER\_ERROR]  
ACL manager requires parameters not passed to routine.
- [DACL\_ERROR\_REQUIRED\_ENTRY\_MISSING]  
Required ACL entry missing from ACL.
- [DACL\_ERROR\_PARAMETER\_ERROR]  
Required pointer parameter has NULL value.
- [DACL\_ERROR\_UNRECOGNIZED\_MGR\_TYPE]  
Unrecognized ACL manager type uuid.
- [DACL\_ERROR\_ENTRY\_TYPE\_TOO\_LARGE]  
ACL entry type is too large to be processed by current code.
- [DACL\_ERROR\_DUPLICATE\_ENTRY\_FOUND]  
Duplicate ACL entry found.
- [DACL\_ERROR\_ACLFILE\_NOT\_FOUND]  
Specified ACL file not found.
- [DACL\_ERROR\_VALIDATION\_FAILURE]  
Flat ACL buffer is of incorrect form.
- [DACL\_ERROR\_UNRECOGNIZED\_USER\_OR\_GROUP]  
User or group not recognized by Registry Server.
- [DACL\_ERROR\_NONLOCAL\_ENTRY\_TYPE]  
Non-local entry type encountered running in local mode.

## **B.2 Filesystem Access Control List Return Values**

The following are filesystem return values.

- [DACL\_ERROR\_FS\_OPEN]  
File system open error.
- [DACL\_ERROR\_FS\_CLOSE]  
File system close error.
- [DACL\_ERROR\_FS\_STAT]  
File system stat error.
- [DACL\_ERROR\_FS\_READ]  
File system read error.
- [DACL\_ERROR\_FS\_WRITE]  
File system write error.



# *X/Open Preliminary Specification*

## **Part 6:**

### **The DCE DFS VFS+ Interface Specification**

This part of the **DCE DFS** provides an overview of the VFS extensions provided by DFS, called the *DCE DFS VFS+ Interface*. This interface is an enhancement of the **DCE 1.1: Distributed File Service** specification that permits different file systems to coexist within one UNIX kernel.



# DCE DFS VFS+ Interface Introduction

This chapter provides an overview of the VFS+ interfaces and behaviors that a DFS-compliant physical file system needs to support. Such a file system will be fully interoperable with other implementations of **DCE DFS** file systems from the viewpoint of all network protocols and the vast majority of user-space software for extended filesystem fileset operations, known as **EFS** software.

## 11.1 Definition of Terms

<b>DFS</b>	The DCE Distributed File System.
<b>EFS</b>	The optional portion of DFS which deals with extended fileset operations (cloning, backup, replicas, and so on).
<b>LFS</b>	Local File System: a physical file system, which in the context of this document, provides DFS semantics.  <b>Episode</b> is the name of the LFS that Transarc Corporation provides as part of the DFS EFS package.
<b>UFS</b>	The <b>UNIX</b> File System. The name of the file system originally provided by vendors of <b>UNIX</b> systems.
<b>Aggregate</b>	An expanse of physical disk which is managed by a DFS LFS. It is similar in context to a <b>UNIX</b> partition and is identified by both a device number and a device file. In many cases, the storage provided to an aggregate will be by a Logical Volume Manager (LVM).
<b>Cell</b>	A collection of DCE machines administered as a single entity. In this context, the key point is that a cell is serviced by a single registry with a single, consistent set of user and group identifiers.
<b>Fileset</b>	A set of related files, connected via a sub-tree with a single root point, that is administered as an entity by DFS. Although filesets cannot span aggregates, a single aggregate can hold many filesets. Filesets which have an associated maximum size or quota, grow independently of each other and compete for space within their aggregate. Individual filesets can be mounted locally, backed up and restored, cloned or transparently moved to another aggregate.  The terms fileset and volume are used interchangeably, although the term volume is becoming archaic.
<b>Group</b>	A membership list, identified by a <b>uuid_t</b> , which contains a list of principals (by <b>uuid_t</b> ) which are members of that group. Again, in a non-DCE environment, they are identified by a 32-bit (at best) <b>gid</b> .
<b>PAC</b>	Privilege Attribute Certificate. A data structure, generated by the DCE Privilege Server, that contains the identity of a DCE authenticated principal. To a first approximation, it contains the principal's <b>uuid_t</b> along with the <b>uuid_t</b> of each group of which it is a member.
<b>Principal</b>	An entity that is interacting with DFS. Although it usually refers to a person, it could as easily correspond to a DCE server or a computer within a DCE cell. In the

DCE environment, principals are identified by a **uuid\_t**. In a non-DCE, standard UNIX, environment, they are identified by a 32-bit (at best) **uid**.

**Realm** Equivalent to a Cell. In security DCE security discussions, the term Realm is typically used instead of Cell.

## 11.2 VFS+ Interface Goals and Constraints

The **DCE DFS VFS+** interface is an enhancement of the VFS interface. The VFS (Virtual File System) interface was introduced by Sun Microsystems to make it possible to implement different file systems within one **UNIX** kernel; specifically, to implement both a physical and a network file system. The original VFS interface, or variants of it used by other vendors of **UNIX** systems, appears in most **UNIX** kernels today. The VFS+ interface extends VFS in several ways, serving to generalize VFS and adapt it to fulfill the requirements of operation in a distributed computing environment.

The principal improvements offered by the VFS+ interface are:

- **Generalized credentials.** The credential structure passed to most of the VFS functions has been augmented to accommodate the use of various different authentication mechanisms. The new credential may thus carry the standard **UNIX** information as well as Kerberos tickets, or any other form of authentication information required. This revision is completely backwards-compatible.
- **Synchronization.** Most operations are redefined to use a synchronization package. In this way, protocol exporters can support notification guarantees to their clients (similar to AFS *callbacks*, for those familiar with the concept). When an object in any type of virtual file system is changed, the synchronization system notifies all interested parties. This allows a more general mechanism for providing appropriate semantics in a distributed file system.
- **Fileset/aggregate interface.** The VFS+ interface includes operations on **filesets**, a key functional enhancement and operability feature in the **DCE DFS** distributed environment. Filesets are located within a new form of data container called an **aggregate**, and the VFS+ interface also allows them to be manipulated.

Historically, filesets were called **volumes**, but the name was changed to avoid confusion with features of other file systems.

- **Portability.** Several portability problems are addressed. Synchronization is equally effective for vendor file systems (such as the **UFS** file system in Sun's SunOS and other Berkeley-derived kernels) and for local file systems (such as the **DCE LFS** that is part of **DCE DFS**). Vendor file systems need not be enhanced to support the full functionality of filesets and aggregates, but they can support the full fileset/aggregate interface, with rudimentary functionality. And lastly, local file systems can be adapted to different kernels. It is necessary to do this without modifying the kernels above the VFS layer. A complication is that some vendor kernels, notably IBM's AIX and DEC's Ultrix, do not have simple variants of Sun's VFS interface, but have their own file system interfaces, with differing data structures (both AIX and Ultrix have a **gnode** structure, and Ultrix does not have a **vnode** structure), and different designs for the boundary between the file system and the higher levels of the kernel. Another complication is that the file system is expected to be integrated with the virtual memory system, and the different kernels have completely dissimilar and independent virtual memory systems.

An important motivation for these changes is compatibility among protocol exporters. For example, one protocol exporter (namely, a VFS+ server) may make guarantees about when an

exported file changes. At the same time, other servers (in particular, NFS), as well as the local UNIX kernel, may simultaneously modify these files.

In order to avoid modifying the file-system-independent layer of the UNIX kernel or the NFS server, the stock VFS vnode-level functions are redefined to call on a synchronization package that keeps track of all guarantees made by the various protocol exporter types. The updated function definitions are implemented by writing what is termed “wrapper” functions for an already-existing virtual file system. These wrappers perform the appropriate synchronization calls before and after the original VFS call.

The VFS+ performance issues are fairly straightforward. The principal additional overheads incurred by using the VFS+ interface are checks for concurrent whole-fileset and checks for conflicting references to vnodes. The check for concurrent use of a fileset is done by lookup in a *fileset registry* (see Section 11.6 on page 239). Checks for conflicting references to vnodes are done by calling a file-system-independent layer with the appropriate file ID.

This layer is necessary to ensure that no matter how a vnode is modified or otherwise accessed, all protocol exporters that may have granted incompatible promises to their clients will have an opportunity to revoke those promises. Promises are represented by *tokens*. The VFS+ interface must obtain *tokens* for files that it manipulates; in this respect it is no different from a DCE DFS client.

The *fileset registry* and the *token cache* are hash tables in which normal lookups, using the file ID, are relatively inexpensive. Optimizations allow the fileset synchronization and token obtaining steps to be skipped, wholly or in part, for those filesets that are not exported, or for filesets that are read-only. Thus, those file systems that may require inter-machine synchronization are the only ones that experience this modest performance penalty.

### 11.3 Overview of Interfaces to the LFS

There are three separate subdomains or areas in the VFS+ interface that are required. A fourth may be present:

1. **Vnode-level facilities.** The vnode functions appropriate to the particular kernel, accessible via a **struct vnodeops** function array or equivalent, are exported. Each type of LFS provides a vector of pointers to its vnode facilities, and a pointer to this function array is placed in vnodes initialized by that LFS. Thus, for vendor file systems, the array will generally not point to the original vendor functions, but will point perhaps to wrapper functions, which perform synchronization. The wrappers will need to call non-wrapped functions (which in the case of vendor file systems could be the original vendor functions). They will access the non-wrapped versions via the same function array, which therefore is extended to about twice the original length.

That, however, isn't the only extension of the **struct vnodeops** function array. The methodology for achieving portability in **DCE LFS** has been to write a set of functions appropriate for Sun's original VFS interface, and to use wrapper functions to extend these for VFS-variant interfaces from other vendors (as well as for Sun's more recent VFS extensions). Another layer of wrapper functions means another extension of the function array, which is consequently about three times the original length.

Of course, the kernel accesses this new triple-length array as if it were just a **struct vnodeops**. But components of the server may access other parts of the array. In particular, the Protocol Exporter accesses the non-synchronized non-portable functions. That is because it does its own synchronization, and it does not have to use the same VFS-variant as the kernel.

2. **Aggregate-level facilities.** An *aggregate* is a data container that holds filesets. Aggregates are a generalization of the UNIX concept of partition, while filesets are a generalization of the UNIX concept of a file system. That is, a fileset provides a single connected file hierarchy, while an aggregate provides a flat address space on disk. In conventional UNIX systems these concepts are identified so that each partition can contain exactly one file system, but in the DFS, an aggregate can contain any number of filesets. The VFS+ interface provides facilities for managing aggregates, including creating new filesets within the aggregate and enumerating existing filesets. Aggregate information is made available via a global *Aggregate Registry* table, where each entry describes the data container for filesets and exposes a function array containing the aggregate operations described above. This table is reached by means of an aggregate type, which is a small integer that is provided by the protocol exporter, that is used as an index to this table. (see Chapter 14 on page 303).
3. **Fileset-level facilities.** Structures and operations are provided to access and manage the contents of individual filesets situated within the above aggregates. Each type of LFS provides a vector of pointers to its vnode operations. A pointer to this vector is returned to DFS by the aggregate operation which attaches filesets. The global *Fileset Registry* table exports entries for each fileset hosted by available aggregates. These entries include a function array (vector of pointers) allowing use of the intra-fileset operations, such as creating, deleting, and getting status on files within a fileset (see Chapter 15 on page 325).
4. **LFS-specific facilities.** A set of facilities that are private or specific to a particular LFS may be present that are used to communicate between LFS-specific utilities and the LFS. These may consist of private system calls or other private routines which are peculiar to a specific vendor's LFS. Nothing else will be said about them other than as noted here - that they are not specifically excluded by this specification.

### 11.3.1 Locking

One of the most important VFS+ aspects is the manipulation of locks above the original VFS interface layer.

In the VFS+ interface, each VFS function is redefined to perform the following operations:

1. **Lock operands and declare the operations to be performed.** In this step, the VFS function calls a locking function common to all virtual file system types. Declarations are made of the file IDs that are about to be processed, as well as the operations that will be performed upon them. As part of this function, other virtual file systems may be invoked to revoke promises (such as *callbacks* and *tokens*) that they may have issued pertaining to the same files.
2. **Perform the desired original virtual file system operation.**
3. **Release the synchronization entries locked in the first step above.** This does not cause the revocation of any guarantees made by the call itself, but rather indicates to the local synchronization package that revocation is now permitted.

### 11.3.2 Credentials

The *credential structure* (**struct ucred**) passed to most VFS functions has also been generalized in the VFS+ interface. Each credential is augmented with a *magic cookie* value that is associated with a property list. This property list (and thus the credential containing the appropriate cookie) is then associated with an arbitrary set of user identification information. Examples of entries on this property list are NFS and Kerberos identity structures.

Using this generalized structure, file systems that understand the credential *magic cookie* will obtain the specific information (such as the user's Kerberos identity) that they need to use the relevant protection mechanisms.

File systems that have not been modified to understand the credential *magic cookie* will continue to use the basic UNIX uid and group information in order to make authorization decisions, oblivious of the extra information available.

### 11.3.3 Operations on Filesets and Aggregates

The final components of the new VFS system are two new virtual interfaces for performing operations on filesets and on aggregates. The functions in the fileset interface implement such operations as:

1. **Taking filesets off-line and putting them back on-line.** While a fileset is off-line, vnode operations on files stored in that fileset will either wait or fail with distinctive error codes, depending on the mode in which the fileset was taken off-line. When the fileset operations in question are completed, that fileset is once again functional, and can be accessed via protocol exporters or via system calls from the local UNIX kernel.
2. **Iterating through all files in a fileset.** Many operations require that some operation be performed on every file within that fileset. Examples are cloning a fileset and restoring a fileset from a dump.
3. **Operating on individual files in a fileset.** Cloning of a fileset makes use of an operation that clones an individual file. Restoring a fileset makes use of several operations on individual files: one to write the file's data; one to write its status; possibly one to create a file, and possibly one to delete a file; and others.

The operations used to support dumping and restoring of filesets correspond roughly to various vnode operations used by the UNIX kernel. But there are important differences. The vnode operations implement UNIX semantics, while the fileset operations do not. For instance, the vnode operation **VOP\_RDWR**, which reads or writes file data, will also set the file's atime (for reading) or mtime and ctime (for writing). The corresponding fileset operations **VOL\_READ** and **VOL\_WRITE**, on the other hand, are carefully coded to avoid modifying *any* of the file's meta-data, except its length. Another important difference is that, given a consistent file system on entry, a vnode operation will restore the file system to consistency at exit, while the volume operations do not maintain consistency. For instance, the above mentioned **VOL\_WRITE** operation is used to restore directories as well as files. When it finishes writing a directory, the directory may well contain dangling references to files, which will be corrected only by subsequent volume operations to create those files. For comparison, vnode operations that modify directories, such as **VOP\_CREATE** and **VOP\_REMOVE**, also create or remove files in such a way that the directory hierarchy is correct by the time the operation returns.

4. **Accessing files on behalf of protocol exporters.** In normal operation, a DFS file exporter obtains a vnode for a file by passing the file's ID and uniquifier (generation number) to a fileset operation.

While a fileset is off-line, an operation on a vnode within that fileset should cause the corresponding virtual file system to either return a distinctive error code or simply wait for the fileset to come back on line.

The functions in the aggregate interface implement such operations as creating a new fileset on an aggregate and enumerating all filesets within an aggregate.

## 11.4 Organization of the VFS+ Switch

Examples of the following description can be found in Section 13.5 on page 290. They can be referred to for further information to gain an enhanced understanding of this description.

### 11.4.1 The O-ops

In the VFS+ interface, the *base* procedures (the ones already defined by the vendor-supplied kernel) are implicitly redefined. In addition to performing their basic functionality, they also call the DFS Token Manager in order to revoke any tokens conflicting with the operation the VFS function is about to perform. The token element will remain locked in the token database until the VFS operation completes.

The redefined procedures, as used by the vendor-supplied kernel, are referred to as the *O-ops*.

### 11.4.2 The N-ops

An additional set of procedures are provided, which represent the *extended* portion of the interface. They are similar to the base set, but assume the caller has performed all the required token synchronization. Since the extended procedures are new to the vnode operations vector, the only programs that will invoke them are components of DFS, such as the *O-ops*. Other clients of the VFS+ interface, such as the UNIX kernel's file system calls, do not know of the existence of these extended functions.

These procedures, having the same specifications as the *O-ops* but not performing synchronization, are referred to as the *N-ops*. If the file system is a vendor-supplied file system, these are generally the original vendor vnode ops.

### 11.4.3 The X-ops

The last extension is a set of functions that do not conform to the vendor's VFS-variant specifications, but rather, conform to the original VFS specifications. These are referred to as the *X-ops*. Because the *X-op* interface is constant across platforms, it is convenient for use by the Protocol Exporter and other server components. For portability, the DCE LFS has been written as a set of *X-ops*, and the corresponding *N-ops* are wrappers for these *X-ops*. Vendor-supplied file systems, on the other hand, come as sets of *N-ops*, and the corresponding base *X-ops* are wrappers for the *N-ops*.

The *X-ops* include some extra operations, besides those in Sun's original VFS. Some of these extra operations facilitate integration with vendor virtual memory systems. Others facilitate communication with the Protocol Exporter, or the Fileset Registry, or allow manipulation of Access Control Lists. A complete set of specifications for the *X-ops* may be found in Section 16.7 on page 424 and Section 16.8 on page 460.



## 11.5 Basic Operation

Because of the extension of the vnode operation vector, there are some basic bootstrapping problems in getting a VFS+ interface working in a kernel. It cannot be assumed that any file system's code is written with any knowledge of the VFS+ interface, or even that the VFS+ interface is present when the kernel is first booted or at any time soon after boot. In other words, the VFS+ interface may be loaded as a kernel extension, after a system has been running for some time.

Pathologically, some kernels may rely upon knowledge of the addresses of vnode operation vectors. For instance, a kernel may determine whether a vnode will require certain resources, by comparing the address of its vnode op vector with the known address of the vnode op vector for a particular file system. This sort of comparison will not work if old vnode ops vectors are replaced by new ones. Thus, such file systems cannot be exported. Appendix C on page 245 discusses typical components in a typical VFS+ package in a typical VFS kernel.

## 11.6 The Fileset Registry

In-memory structures describing filesets local to the host are kept in a hash table called the *Fileset Registry*. The **struct volume** fileset descriptor records kept there (see are hashed by volume ID. One fileset descriptor exists for each fileset served by that host. There are various ways to generate references to these fileset descriptor structures. Section 11.7 shows how to determine the proper fileset pointer given an already existing vnode. Section 15.14 on page 375 describes the *vol\_vget()* call, which converts a file ID to a vnode containing the proper fileset descriptor reference.

Protocol exporters will need to have the *Fileset Registry* pre-loaded by an application program which iterates over the filesets on the server's disk(s) and loads fileset descriptors for each.

Each fileset descriptor stored in the *Fileset Registry* contains a pointer to an array of operations (**struct volumeops**) that may be performed on the associated fileset. Thus, given a fileset ID, acquiring the *Fileset Registry* entry for that fileset allows the caller to perform such operations as attach, destroy, and open, as well as operations on the individual files contained within that fileset. The **struct volumeops** operation array is used by such agents as the Fileset Server to perform higher-level activities such as cloning (replicating), moving, and dumping filesets. Refer to Chapter 15 on page 325 (in particular, the overview in Section 15.1 on page 325) for a more detailed examination of the fileset portion of the VFS+ interface.

## 11.7 The Fileset/Vnode Interfaces

These sections describe the interface between the fileset and vnode layers in more detail. They primarily define how, given a vnode, one finds the corresponding fileset.

A new VFS+ function, **VOPX\_GETVOLUMEa**, is defined that takes a vnode pointer and returns the corresponding fileset pointer.

The UFS implementation of this function looks up the VFS pointer of the vnode in a hash table. The table associates VFS pointers with *fileset registry* entry pointers. In a file system for which exportation would make no sense, such as NFS, the implementation of this function would simply return a NULL pointer, indicating that the fileset concept does not apply to this file system type.

In **DCE LFS**, this function obtains a handle for the fileset from the file-system-specific portion of the vnode structure. From this handle, it obtains the fileset's ID, and looks up that ID in the *fileset*

*registry* in the normal way.

## 11.8 References from the LFS back into DFS

There is a set of core functionality that is shared between **DCE DFS** (herein called DFS) and an **DCE LFS**, herein called an LFS. These are provided by a set of DFS functions which LFS calls. This section will list this functionality. In the interest of minimizing the interactions between DFS and an LFS, only those calls which are truly needed (or seem potentially useful) and cannot be implemented by code within the LFS itself are actually listed.

### 11.8.1 LFS Introducing Itself to DFS

An LFS makes itself known to DFS via an *ag\_setops()* call which supplies DFS with an aggregate operations vector to be used for the given aggregate type (AG\_TYPE\_xxx). These are listed in Section 13.2.4 on page 270 and are used in several places such as Section 13.2.1 on page 269.

**Note:** In Section 13.2.4 on page 270, *Valid Aggregate Types*, there is also a size definition. If one were to add a new type (or types), one would also increase the size of the aggregate operations structures by increasing this size appropriately.

Additionally, the aggregate operations vector whose structure is defined in Section 13.2.6 on page 270, *Aggregate Operations Vector*, is supplied. The following signature defines the *ag\_setops()* function call.

```
/* agtype:   AG_TYPE_xxx encoding for the LFS
 * ops:     supplied aggregate operations vector
 */
extern void ag_setops(
    /* IN */ int    agtype,
    /* IN */ struct aggrops    *ops
);
```

### 11.8.2 Operating on DFS Lock Structures

The next set of functions are used by the *vol\_lock()*, *vol\_unlock()*, *ag\_lock()* and *ag\_unlock()* operations to operate on DFS lock structures.

The following locks are defined in the DFS header file **src/file/osi/lock.h**.

```
#define READ_LOCK      1
#define WRITE_LOCK     2
#define SHARED_LOCK   4
```

The following lock structure is also defined in this same header. It is the standard "data lock". All locks wait on *excl\_locked* except for *READ\_LOCK*, which waits on *readers\_reading*

```
struct lock_data {
    unsigned char wait_states;    /* type of lockers waiting */
    unsigned char excl_locked;    /* boosted, shared or write lock? */
    unsigned char readers_reading; /* readers actually with read locks */
    unsigned char num_waiting;    /* probably need this soon */
};
```

The following locking functions are defined in this header as well. As previously mentioned, these functions are used to operate on the DFS lock structure, **lock\_data**, defined above.

```
extern lock_ObtainRead(lock_data);
extern lock_ObtainWrite(lock_data);
extern lock_ObtainShared(lock_data);
extern lock_ReleaseRead(lock_data);
extern lock_ReleaseWrite(lock_data);
extern lock_ReleaseShared(lock_data);
```

### 11.8.3 Allocating and Freeing Memory

The next two calls are used to allocate and free kernel memory. They are specifically needed by the *ag\_attach()* and *ag\_detach()* functions. They are defined in the DFS header file *src/file/osi/osi.h*. See also Appendix F on page 323 for additional information.

```
/*
 * typedef needed for the Alloc and Free functions
 */

typedef unsigned int size_t;

/*
 * Generic allocation routine
 * Returns pointer to allocated storage, else panics ;
 * CAN block
 *
 * asize:    allocation size wanted, in bytes
 */
extern opaque osi_Alloc(
    size_t    asize
);

/*
 * Generic memory deallocation routine
 *
 * x:    pointer to whatever is being freed
 * asize:    ... and its size (in bytes)
 */
extern void osi_Free(
    opaque    x,
    size_t    asize
);
```

### 11.8.4 Fileset Creation Assistance

The following call is needed by the *ag\_volCreate()* operation. Its use is described in the chapter, *Aggregate Operations Interface*, under that function (see Chapter 14 on page 303). There is also more information in Appendix D on page 317. It is defined in the DFS header file *src/file/xvolume/vol\_init.h*.

```
/* Returns 0 on success, else an error code */

extern int vol_Attach(
    afsHyper        *volId,
    struct vol_status *statusp,
    struct aggr      *aggrp,
```

```

    struct volumeops    *volopsp
    );

```

### 11.8.5 Obtaining a Fileset Structure

The following call is used by several **fileset** and **vnode** operations to determine if a given fileset (by ID) is exported to DFS and, if so, to obtain a held fileset (volume) structure for it. It is defined in the DFS header file **src/file/volreg/volreg.h**. Its signature can be found in Section 15.15 on page 399 in *Fileset Registry Array Functions*.

```

/* Returns 0 on success, else an error code
 *
 * fidp:    fileset ID
 * volp:    if non-NULL, pointer to held volume
 *           structure is returned here
 */
extern int volreg_Lookup(
    struct afsFid *fidp,
    struct volume **volp
);

```

### 11.8.6 Releasing a Fileset

The following call is made by the *vol\_rele()* call in Section 15.14 on page 340 when the usage count on a volume structure goes to zero.

It is defined in the DFS header file **src/file/xvolume/vol\_init.h**.

```

/* Always returns 0 (success)
 *
 * volp:    pointer to volume structure whose usage
 *           count has gone to 0.
 */
extern int vol_VolInactive(
    struct volume *volp
);

```

### 11.8.7 Conversion of Operations Vectors

The following calls assist in the "conversion" of **vnode** and **vfs** operations vectors. Consult the portion of this document on **vnode** operations, Chapter 16 on page 415.

These don't appear to exist in any DFS header files. Their source resides in the DFS file **src/xvnode/xvfs\_vnode.c**.

```

/* afuns:    pointer to LFS extended vnode ops vector
 * axfuns:   pointer to enhanced (3-part) ops vector
 *           that is constructed by this call
 */
extern void xvfs_InitFromXOps(
    struct xvfs_xops    *afuns,
    struct xvfs_vnodeops *axfuns
);

/* aofuns:   pointer to LFS vfs operations vector
 * afuns:    pointer to enhanced (2-part) vfs vector

```

```

*           that is constructed by this call
* getvolfn: pointer to supplied vfs_getvolume
*           function
*/
extern void xvfs_InitFromVFSOps(
    struct osi_vfsops  *aofuns,
    struct xvfs_vfsops *afuns,
    int                (*getvolfn)()
);

```

### 11.8.8 Getting the Local Cell ID

The following call returns the DCE local Cell ID (consult Chapter 12 on page 251 for details on its use).

```

extern dacl_GetLocalCellID(
    afsUUID      *localCellID
);

```

### 11.8.9 Getting the Administrative Group ID

The following call returns the uuid of the DFS administrative group. It is the general function provided for determining if the caller is the member of the DFS administrative group.

```

extern dacl_GetSysAdminGroupID(
    afsUUID      *groupID
);

```

### 11.8.10 Obtaining the Identity of the Principal

As described in the *Obtaining a Principal's Identity* section of the ACL document (see Chapter 12 on page 251), the following calls might be necessary.

**Note:** As part of its support, a vendor could instead provide a wrapper function within DFS that performs the full task of extracting a PAC from a set of credentials. Their signatures can be found in Section 17.1 on page 469, *The xcred Package*, in Section 17.3 on page 477, and Section 17.3 on page 480, respectively.

```

/* aucredP:      input unix credentials
* axcredPP:     pointer to returned extended credentials
*/
extern xcred_UCredToXCred(
    struct ucred *aucredP,
    cred_t      **axcredPP
);

/* axcredP:      input extended credentials
* aattributeP:   desired attribute
* aattributeLength: ... and its length
* avaluePP:      returned attribute value
* areaLengthPP  ... and its length
*/
extern xcred_GetProp(
    xcred_t      *axcredP,
    char         *aattributeP,

```

```
long      aattributeLength,  
char      **avaluePP,  
long      *areaLengthP  
);
```

## Components of a Typical VFS+ Package

### C.1 The VFS Vector

The first item to be considered in providing a VFS+ package is providing a new VFS vector for the file system to be exported. Typically, the new vector is like the old vector, but some of the functions are wrappers. In particular, the *vfs\_root* and *vfs\_vget* functions are wrappers; they call the original function to obtain a vnode, and then *convert* that vnode, by replacing its old vnode ops vector with a new one.

Converted (VFS+) vnodes can be distinguished from normal (VFS) vnodes by a status flag in the **v\_flag** field, **V\_CONVERTED**. (In IBM's AIX, the flags field is the **gn\_flags** field of the gnode; in DEC's Ultrix, the flags field is the **g\_flag** field of the gnode.) If this flag is set, the vnode is a VFS+ vnode, and the functions pointed to by the **v\_op** field have their VFS+ definitions. Otherwise, the functions are basic VFS functions and have their usual definition.

The conversion process works by checking for the **V\_CONVERTED** flag. If present, there is nothing to be done, and the conversion process terminates. Otherwise, an updated set of vnode functions is generated from the existing vnode functions, and the updated function array is placed into the vnode's **v\_op** field. Finally, the **V\_CONVERTED** flag is set in the **v\_flag** field in the converted vnode.

When an updated set of vnode functions has been generated for one vnode in a file system, it can be used for all the other vnodes in the same file system. So to save conversion effort, an association list can be maintained, in which each pair consists of an old vnode ops array, and a new one. When a vnode is to be converted, the address of its vnode ops array can be looked up in the list, and if it is found (as an "old" array), the appropriate "new" array replaces it.

### C.2 Naming Conventions

In this document, naming conventions are used for both the macros defined for vnode functions as well as the functions themselves.

- **Vnode functions.** One naming convention identifies the "family" to which a given vnode function belongs. Let us take the example of vnode functions implemented for the UNIX BSD file system, or **UFS**. Furthermore, let *fname* be the suffix portion of the function name pointed to from the **struct vnodeops** array, such as *rename*, *lookup*, *open*, and *mkdir*. In this appendix, *ufs\_fname()* refers to the version of *fname* as it normally appears in the **struct vnodeops** function array straight from the vendor. The *glue* version of the same function, whose address is overwritten into this array (as an *O-op*) and which performs the additional synchronization operations transparently, is referred to as *xglue\_fname()*. Finally, the extended version of the same function, the one that is normally used by the **DFS** class of exporters directly (as an *X-op*) and for which the caller has already performed all the synchronization work is named *xufs\_fname()*. See Section 11.4 on page 238 for more information on *O-ops*, *N-ops* and *X-ops*.
- **Vnode macros.** There are a standard set of vnode macros defined by the file system. As in the previous paragraph, let *fname* be the suffix portion of the function name pointed to from the **struct vnodeops** array. Then, in Sun's original VFS, **VOP\_FNAME** is the macro that indireacts through the **struct vnodeops** array in the vnode descriptor and calls the appropriate

base function there. For clarity in DFS, an equivalent macro, **VOPO\_FNAME**, is defined which accesses the same slot in the array but indicates clearly that an *O-op* is expected to be there. A second set of macros, **VOPX\_FNAME**, used by DFS exporters, is also defined which indirections through the same array but calls the extended version of the same function (that is, the *X-op*). Similarly a third set of macros, **VOPN\_FNAME**, is defined to indirect through the *N-op* slots in the same array.

### C.3 Glue Functions

Most of these *glue functions* are derived from the corresponding vnode functions by writing a new function that calls the old one as part of its operation. Here is pseudo-code for the *xglue\_rename()* function. It is written with an old-style Sun VFS interface in mind; the versions for other VFS-like interfaces, such as that used in IBM's AIX, would be somewhat different.

#### C.3.1 xglue\_rename() Function

##### Example C-1 xglue\_rename() Function

```
xglue_rename(vnode, name, tovnnode, toname, cred)
struct vnode *vnode, *tovnnode;
char *name, *toname;
struct ucred *cred; {
    struct vnode *tvpl, *tvp2;
    struct afs_tokenSet *tset;
    struct volume *vold;

    if (code = ReferenceCorrespondingVolume(vnode, &vold))
        return(code);

    if (NoTokensRequired (vold)) {
        code = VOPN_RENAME(vnode, name, tovnnode, toname, cred);
        return(code);
    }

    VOPN_LOOKUP(vnode, name, &tvpl, cred);
    ConvertVnode(tvpl);
    VOPN_LOOKUP(tovnnode, toname, &tvp2, cred);
    ConvertVnode(tvp2);

    tset = NewTokenSet();
    AddTokenSet(tset, vnode, DATA_WRITE+STATUS_WRITE);
    AddTokenSet(tset, tovnnode, DATA_WRITE+STATUS_WRITE);
    AddTokenSet(tset, tvpl, STATUS_WRITE);
    AddTokenSet(tset, tvp2, STATUS_WRITE);
    ObtainTokenSet(tset);

    code = VOPN_RENAME(vnode, name, tovnnode, toname, cred);

    ReturnTokenSet(tset);
    VN_RELE(tvpl);
    VN_RELE(tvp2);
}
```



```

    ReleaseCorrespondingVolume(vold);
    return(code);
}

```

The **Token Manager** is called to acquire and release tokens on all of the vnode-class operands, whether they are explicit or implicit in the VFS interface. Thus tokens must be obtained for the object being renamed, and the object (if any) that it will replace, though their vnodes are not explicit parameters of the function.

### C.3.2 xglue\_lookup() Function

#### Example C-2 xglue\_lookup() Function

Calls like those done via the *VOP\_LOOKUP* macro that potentially return new vnodes must also call the conversion function on the returned vnode parameters. For instance, the *xglue\_lookup()* call looks like this:

```

xglue_lookup(vp, name, vpp, cred)
struct vnode *vp, **vpp;
char *name;
struct ucred *cred; {
    long code;
    struct afs_tokenSet *tset;
    struct volume *vold;

    if (code = ReferenceCorrespondingVolume(vp, &vold))
        return(code);

    if (NoTokensRequired (vold)) {
        code = VOPN_LOOKUP(vp, name, vpp, cred);
        return(code);
    }

    tset = NewTokenSet();
    AddTokenSet(tset, vp, STATUS_READ);
    ObtainTokenSet(tset);
    code = VOPN_LOOKUP(vp, name, vpp, cred);
    ReleaseTokenSet(tset);
    if (*vpp)
        ConvertVnode(*vpp);    /*Convert to VFS+ vnode*/

    ReleaseCorrespondingVolume(vold);
    return(code);
}

```

Note that the returned vnode (**\*vpp**) is converted into a VFS+ vnode. Ensuring that all vnodes in the system are converted is accomplished by intercepting the *vfs\_vget()* and *vfs\_root()* VFS functions as well as all returns of newly-created vnodes. The result is that after the interception code is installed, all VFS vnodes will be converted to VFS+ vnodes, and will thus call the VFS+ code instead of the normal VFS code from then on.

In other words, after the VFS+ interception code is enabled, all processing that starts with a mount point and continues down the file system tree via *lookup()* and *create()* calls will deal exclusively with VFS+ vnodes, rather than the generic VFS vnodes.

Notice the calls to *ReferenceCorrespondingVolume()* and *ReleaseCorrespondingVolume()* in the code above. These functions are responsible for ensuring that whole-fileset operations synchronize properly with the basic vnode operations. If the fileset concept is meaningless for a particular type of file system, these functions need do nothing but return zero (in which case the vnode operations will never wait), or return *inaccessible* when referring to data on a particular file system.

If there is a meaningful fileset concept for a particular file system, then there is a reference-counted fileset structure in one of three states:

**Normal** File operations may proceed on this fileset.

**Busy** File operations should wait until the fileset is no longer marked as *busy*.

**Error** File operations should fail with the specified error code.

## C.4 Extended Functions

Pointers to the *xglue\_rename()* and *xglue\_lookup()* functions replace those of native (base) functions in the vnode function array. The extended functions, namely the ones that do *not* make calls to the Token Manager, must also be defined. *In most cases, the extended version of the vnode function is identical with the original, vendor-supplied version.* It is assumed for all extended functions that the caller has already dealt with all the synchronization requirements, so it is safe to call the original vendor code. However, all functions in the extended set that return new vnodes must also make sure that these resulting vnodes are the VFS+ versions (that is, they point to the augmented **struct vnodeops** function array). Thus, one may simply use the *ufs\_rename()* function to implement *xufs\_rename()* since it doesn't return any new vnode references. This cannot be said for *xufs\_lookup()*, which must be implemented as follows (old Sun VFS version):

### C.4.1 xufs\_lookup() Function

#### Example C-3 xufs\_lookup() Function

```
xufs_lookup(vp, name, vpp, cred)
struct vnode *vp, **vpp;
char *name;
struct ucred *cred; {
    long code;

    code = VOPN_LOOKUP(vp, name, vpp, cred);
    if (!code)
        ConvertVnode(*vpp);    /*Convert to VFS+ vnode*/
    return(code);
}
```

This function, like all extended vnode functions, does not attempt fileset synchronization or token acquisition, but simply performs the native version of the lookup operation and then converts the resulting vnode, if one was generated.

## **C.5 Extended Vnode Attributes**

The vnode attribute structure (**struct vattr**) used in UNIX kernels is extended for purposes of exporting files. The extended structure (**struct xvfs\_attr**) is given in Section 13.6.3 on page 294. The X-ops, VOPX\_GETATTR and VOPX\_SETATTR, take an extra argument, which is a flag specifying whether or not the attribute structure argument is extended.



# The DCE DFS ACL Model for an LFS

This chapter discusses ACLs, security and protection checking as they relate to a DFS physical file system.

## 12.1 Overview

DCE DFS ACLs are an extension of POSIX ACLs. The specific algorithms presented here are based on draft 13 (/D13) of the POSIX 1003.6.1 spec. This is a recent draft which is simplified in several respects from the earlier ones which are reflected in some of the existing implementations of the DCE DFS. It is assumed that the reader has access to this document.

**Note:** Since it is not clear what POSIX will finally decide on here, this document also describes (clearly annotated) the older behaviors, as implemented currently in Episode.

## 12.2 Definition of Terms

The terms here are defined in Section 11.1 on page 233. They are reproduced here for convenience as they are used extensively in the descriptions that follow.

<b>Cell</b>	A collection of DCE machines administered as a single entity. In this context, the key point is that a cell is serviced by a single registry with a single, consistent set of user and group identifiers.
<b>Group</b>	A membership list, identified by a <b>uuid_t</b> , which contains a list of principals (by <b>uuid_t</b> ) which are members of that group. Again, in a non-DCE environment, they are identified by a 32-bit (at best) <b>gid</b> .
<b>PAC</b>	Privilege Attribute Certificate. A data structure, generated by the DCE Privilege Server, that contains the identity of a DCE authenticated principal. To a first approximation, it contains the principal's <b>uuid_t</b> along with the <b>uuid_t</b> of each group of which it is a member.
<b>Principal</b>	An entity that is interacting with DFS. Although it usually refers to a person, it could as easily correspond to a DCE server or a computer within a DCE cell. In the DCE environment, principals are identified by a <b>uuid_t</b> . In a non-DCE, standard UNIX, environment, they are identified by a 32-bit (at best) <b>uid</b> .
<b>Realm</b>	Equivalent to a Cell. In security DCE security discussions, the term Realm is typically used instead of Cell.

## 12.3 Primitive Data Types

There are two primitive data types used in this chapter, `sec_id_pac_t` and `permissions_t`. Their definitions can be found in Section 13.1 on page 267.

## 12.4 Local Realm

Several of the **ACL** algorithms detailed below require the identity of the local DCE Realm (Cell). During **DCE DFS** and **LFS (DCE LFS)** initialization, this value (a `uuid_t`) is supplied to the system via a mechanism discussed elsewhere.

The simplest mechanism is to use the `dacl_GetLocalCellID()` call detailed in *Getting the Local Cell ID*, Section 11.8.8 on page 243. In this document, this Cell (by `uuid_t`) is referred to as the LOCAL-REALM.

## 12.5 DFS Administrator

DFS has the notion of a special System Administrative group which is granted the right to perform various control functions. (In many regards, this group enjoys the same file system privileges that ROOT on a typical UNIX system possesses.)

The simplest mechanism is to use the `xvfs_IsAdminGroup()` or `dacl_GetSysAdminGroupID()` call detailed in *Getting the Administrative Group ID*, Section 11.8.9 on page 243. If DFS is not active, there is no notion of an administrator.

In this document, this group (by `uuid_t`) is referred to as the DFS-SYSADMIN-GROUP.

## 12.6 DFS vs non-DFS requests

In certain situations, the desired behavior depends upon whether or not the operation being performed originated from DFS. Although a number of short cuts are possible (say, if DFS is not installed at all), the *logical* determination of this fact is described in *Obtaining the Identity of the Principal*, Section 11.8.10 on page 243.

## 12.7 ACL Contents

A DFS **ACL** consists of a certain amount of fixed information along with a variable number of ACL entries. The fixed information is the following:

*Manager Type*     A UUID defining the type of manager that supports this **ACL** type. Alternatively, one can view it as defining an **ACL** type. In order to be interoperable with other physical file system implementations, the constant picked for the Episode LFS should be used.

*Default. Realm*     A UUID identifying the realm that non-qualified (by an explicit realm) entries in this **ACL** are interpreted relative to. Usually, this will be the local cell within which this file system resides.

The definition for **ACL** structure, of type `dacl_t`, can be seen in *ACL Structure*, Section 9.8 on page 169.

The remainder of the **ACL** consists of some number of **ACL** entries. **ACL** entries come in several different types, each of which contains a type field, a set of permission rights and some amount of information identifying the principal(s) that this entry applies to. The different types of **ACL** entries are listed below. They are an amplification from an LFS point of view, of the information found in *ACL Entry Types*, Section 8.2 on page 156.

- user\_obj* Contains: permission set. The rights granted to principals within the default realm that belong to the user file class of the file. (That is, the file *oid* field matches the principal's id field.)
- group\_obj* Contains: permission set. The rights granted to principals within the default realm that belong to the group file class of the file. (That is, the file *gid* field matches either the principal's *gid* field or one of groups in its group list.)
- other\_obj* Contains: permission set. The rights granted to any principals from the default realm not matched by any of the user or group types.
- Note:** A DFS **ACL** is required to contain exactly one of each of the above three entry types (*user\_obj*, *group\_obj* and *other\_obj*). Likewise, the *user\_obj* entry is required to contain the **perm\_control** right. If either of these restrictions are violated, the **ACL** is corrupt and should be rejected by a *vn\_setacl()* operation.
- user* Contains: user UUID, permission set. The rights granted to the specified user in the default realm. No two user entries can contain the same user UUID.
- group* Contains: user UUID, permission set. The rights granted to the specified group in the default realm. No two group entries can contain the same group UUID.
- foreign\_user* Contains: realm UUID, user UUID, permission set. The rights granted to the specified user from the specified realm. No two *foreign\_user* entries can contain the same realm UUID and user UUID pair.
- foreign\_group* Contains: realm UUID, group UUID, permission set. The rights granted to the specified group from the specified realm. No two *foreign\_group* entries can contain the same realm UUID and user UUID pair.
- Note:** The default realm in the above two entries (*foreign\_user* and *foreign\_group*) is required to be different than the **ACL**'s default realm. This requirement might not be necessary or desirable given **ACL** inheritance.
- foreign\_other* Contains: realm UUID. The rights granted to any principals from the specified realm that are not matched by any of the *foreign\_user* or *foreign\_group* entries. No two *foreign\_other* entries can contain the same realm UUID.
- any\_other* Contains: permission set. The rights granted to any principals not matched by any of the other entries. At most one *any\_other* entry can appear in an **ACL**.
- extended* Contains: an uninterpreted string of bytes (and length). Ignored (but preserved) by DFS.
- Note:** This is for future extensions.
- mask\_obj* Contains: permission set. This entry type corresponds to functionality dropped in the POSIX. The safest behavior, from an interoperability viewpoint, is to ignore but preserve this entry type. At most one **mask\_obj** entry can appear in an **ACL**.
- Note:** This is an optional entry. Not all **ACL**s contain this entry.

## 12.8 Define the External ACL Representation

Although there are no requirements on **ACL** representation on disk and in memory, the *external* representation as seen outside the LFS is fixed.

The only limit on **ACL** size is that its external representation must fit within a 8192-byte **dfs\_acl** structure.

```
struct dfs_acl {
    long    acl_len;
    char    acl_val[8188];
};
```

The *acl\_len* field gives the number of bytes in *acl\_val[]* that are valid. For an invalid or empty **ACL**, this will be 0.

The first few components stored in *acl\_val[]* are fixed as follows:

**uuid\_t** A manager UUID which gives type of **ACL** this is (literally, it identifies the manager which interprets it). This field's only use is to be checked for validity when an **ACL** is presented to the file system. If general interoperability is desired, LFS implementations should use the same UUID value that was selected for its Episode LFS. Specifically:

```
uuid_t LFSMgrUuid = { 0xd076c532, 0x0a1d, 0x11ca,
    {0x95, 0x3d, 0x02, 0x60, 0x2e, 0xa9, 0x6e, 0x00}
};
```

**uuid\_t** The default realm that this **ACL** applies to.

**int32** The number of **ACL** entries which follow.

Following this is some number of **ACL** entries in an unspecified order. **ACL** entry types are identified by a type field with the values shown in the following table. See *ACL Entry Type Format*, Section 9.5 on page 167, for their definitions.

ACL Entry Type	Value
<b>user_obj</b>	0
<b>group_obj</b>	1
<b>other_obj</b>	2
<b>user</b>	3
<b>group</b>	4
<b>mask_obj</b>	5
<b>foreign_user</b>	8
<b>foreign_group</b>	9
<b>foreign_other</b>	10
<b>unauth_mask</b>	11
<b>extended</b>	12
<b>anyother</b>	13

Table 12-1 ACL Entry Types



### 12.8.1 Formats of ACL Types

The formats for the various types of **ACL** entries are shown below. These formats are defined by the **ACL** entry formats in *Complex ACL Entry Format*, Section 9.7 on page 168, and *Simple ACL Entry Format*, Section 9.6 on page 168, for those that are not foreign or extended. That is to say, the simple entries contain the permissions, *perms*, and the value from the above table (to determine if the entry is good), while the foreign types also have a realm *uuid*. This distinction can be seen in the definitions shown in Section 9.9.1 on page 170 and Section 9.9.2 on page 170.

- *user\_obj, group\_obj, other\_obj, mask\_obj, any\_other, unauth\_mask*

**permissions\_t** Associated permissions.

**int32** Type value from above table.

- *user, group*

**permissions\_t** Associated permissions.

**int32** Type value from above table.

**uuid\_t** *user* or *group* **uuid**.

- *foreign\_user, foreign\_group*

**permissions\_t** Associated permissions.

**int32** Type value from above table.

**uuid\_t** *user* or *group* **uuid**.

**uuid\_t** Realm **uuid**.

- *foreign\_other*

**permissions\_t** Associated permissions.

**int32** Type value from above table.

**uuid\_t** Realm **uuid**.

- *extended*

**uuid\_t** Extension type.

**u\_int16[3]** Miscellaneous stuff.

**char[2]**

**U\_int32** The number of bytes following this.

**char[]** (Whatever).

**Note:** All **ACL** entries must begin on an integral long (32-bit) boundary in memory. For the purpose of computing the total size of the **ACL** and placing additional entries after it, the size of an extended entry is logically rounded up to a multiple of 4 bytes in size. (Its internal size field can contain an arbitrary value, however.)

In this external representation, all fields are stored in network canonical order. Fields of type *long*, *int32* and **permissions\_t** are run through *htonl()*. Fields of type **short** (int16) are run through *htons()*. Fields within a **uuid\_t** structure are converted as follows:

- Fields with *.time\_low* are run through *ntohl()*.
- Fields with *.time\_mid* are run through *ntohs()*.

- Fields with *.time\_hi\_and\_version* are run through *ntohs()*.

Consult your `<stds.h>` header file for definitions of the **integer**, **long** and **ntoh** functions.

## 12.9 Relationship — ACLs and UNIX Protections

There is a direct equivalence between the *user\_obj*, *group\_obj* and *other\_obj* required entries in an **ACL** and the standard UNIX protection bits (the mode bits). If an object does contain an **ACL** (as always, optional), these must remain consistent. P Although there are any number of implementation possibilities, the most obvious ones are listed below.

**Approach 1** Don't associate the standard UNIX protection bits with a file. Instead, require that all objects have an associated **ACL** with at least the required *user\_obj*, *group\_obj* and *other\_obj* entries.

From a performance and disk utilization point of view, it seems most reasonable to always hold these **ACL** entries in an object's inode (anode).

With this approach, changing the mode bits on a file automatically adjusts the **ACL** in a corresponding fashion.

**Approach 2** Associate both the UNIX protection bits and an optional **ACL** with an object. When the object's mode bits are modified as a result of a **chmod** operation or file creation, the required entries in the **ACL** are changed.

Likewise, when an object's **ACL** is explicitly changed as a result of a set **ACL** operation, its mode bits are changed as well (to agree with the **ACL**).

With this approach, an object's **ACL** is always consistent with its UNIX permission bits. However, it may suffer from performance and disk utilization problems.

**Approach 3** Associate both the UNIX protection bits and an optional **ACL** with an object.

When an object's **ACL** is explicitly changed as a result of a set **ACL** operation, its mode bits are changed as well (to agree with the **ACL**).

When the object's mode bits are modified as a result of a **chmod** operation or file creation, its **ACL** is left unchanged.

When an **ACL** is returned to user space during a get **ACL** operation, the read-write-execute permissions in its *user\_obj*, *group\_obj* and *other\_obj* entries are changed to agree with potentially more up-to-date mode bits on the object.

**Note:** This does not occur for initial-file and initial-directory **ACLs** on a directory.

When this step is performed for a directory's **ACL**, the **perm\_insert** and **perm\_delete** rights will additionally be turned on if the **perm\_write** one ends up being set (regardless of its initial state). Note the assymetry here: although these **perm\_insert** and (or) **perm\_delete** rights may be turned on during this process, they will not be turned off -- even if the **perm\_write** right is removed.

When permission checking is done against an **ACL**, the permissions from the mode bits are used instead of what is found in the *user\_obj*, *group\_obj* and *other\_obj* entries.

This approach has the same advantageous performance and disk utilization behavior that **Approach 1** enjoys.

**Note:** The second approach is not recommended at all. Either of the other two schemes seems reasonable. Episode follows the third approach.

The access check algorithm in *Access Rights Algorithm*, Section 12.16 on page 262, assumes **Approach 3**.

### 12.9.1 Episode Visibility above VFS+

As pointed out in the **fileset** and **vnode** operations chapters, under *vol\_getacl()*, (see Section 15.14 on page 365 and so forth) *vol\_setacl()*, *vol\_copyacl()*, *vn\_setacl()* and *vn\_getacl()* (in Section 16.8 on page 465 and so forth), Episode's choice of **Approach 3** above appears to make itself visible above the VFS+ layer. Specifically:

- [1] Episode follows approach 3. Although an object's mode bits and regular **ACL** can disagree, the mode bits are the **truth** used during **ACL** checking algorithms.
- [2] The *vn\_getacl()* operation coerces (depending upon **ACL** type) the **ACL**'s required entries to agree with the object's mode bits, while *vol\_getacl()* doesn't.
- [3] The *vn\_setacl()* operation coerces (depending upon **ACL** type) the object's mode bits to agree with an **ACL**'s required entries, while *vol\_setacl()* doesn't.
- [4] The *vol\_copyacl()* operation, along with *vn\_setacl()* when invoked in its copy mode, totally ignore an object's mode bits.

Items [2] and [3] don't seem to be that major since *vol\_getacl()* and *vol\_setacl()* are used during backup/restore/move operations. The important point here is that an object backed up by a *vol\_getattr()* and *vol\_getacl()* (in that order, although it shouldn't matter) appear **UNCHANGED** after it is restored by a *vol\_copyacl()* or *vol\_setacl()* and *vol\_setattr()* (in that order). A major concern is guaranteeing interoperability of the dump/restore format with other LFS implementations such as Episode.

Item [4] is more difficult. However, it has the virtue of allowing for maximum **ACL** sharing if the application (**acl\_edit**, for instance) takes care of the mode bits itself.

Taking everything into account, the recommendation is to take the safest course and do as Episode does (**Approach 3**).

## 12.10 ACL Creation From Mode Bits Algorithm

There are several circumstances in which it is necessary to construct an **ACL** from a set of UNIX mode bits. One example is when a get **ACL** operation is performed on an object which does not contain an **ACL**. In such cases, the **ACL** is built as follows.

- a. The manager type field (**uid\_t**) is set to the correct value (LFSMgrUuid from above).
- b. The default realm field (**uid\_t**) is set to the LOCAL-REALM.
- c. The *user\_obj*, *group\_obj* and *other\_obj* entries (the required ones) all have their permissions taken from the mode bits in question. The **perm\_read**, **perm\_write** and **perm\_execute** rights are copied in directly. If the object in question is a directory, **perm\_insert** and **perm\_delete** rights are granted whenever **perm\_write** rights are. For the *user\_obj* entry, the **perm\_control** right is added as well.
- d. The *number-of-entries* field within the **ACL** is set to 3.

## 12.11 Initial ACL and File Creation Algorithm

Any file system object can have an attached **ACL** which governs access to that object. Directories may additionally have two associated optional **ACLs** that are used during object creation.

The *Initial Object (Creation) ACL*, if present, is applied to any files created under that directory. (Initial **ACL**.)

The *Initial Container (Creation) ACL*, if present, is applied to any sub-directories created under that directory.

**Note:** Within this chapter, an unqualified use of the term **ACL** usually refers to the regular **ACL** used for access checking.

When an object is created, an initial **ACL** may be applied to or created for it as follows.

1. When a file is created, any *Initial Object ACL* of its parent directory is given to the object. If there is no *Initial Object ACL*, proceed to step [4].
2. When a directory is created, any *Initial Container (Directory) ACL* of its parent directory is given to the new directory. Additionally, the directory's *Initial Object* and *Initial Container* **ACLs** are inherited from the parent directory as well. If there is no *Initial Container* directory **ACL**, proceed to step [4].
3. The object's mode bits (0777) are replaced by a value computed from the *user\_obj*, *group\_obj* and *other\_obj* entries of its **ACL** (**perm\_read**, **perm\_write**, **perm\_execute**). If any of these entries are invalid, the corresponding mode bits are left 0. (In other words, the initial **ACL** overrides the initial mode bit setting)
4. If the create operation occurs via DFS, the remote principal's PAC indicates the realm from which the request originates (the *.realm* field of the **sec\_id\_pac\_t** type). For a non-DFS operation (there is no PAC), the request is viewed as originating from the LOCAL-REALM.
5. If the originating realm is the LOCAL-REALM and the newly created object has no **ACL** (no inherited **ACL**), proceed to step [6]. If the originating realm agrees with the default realm in the object's **ACL**, proceed to step [6].

Otherwise, the existing **ACL** (if any) must be modified to reflect the realm from which the create request actually originated.

If the object has no **ACL**, one is synthesized for it from the file's mode bits as described in *ACL Creation From Mode Bits Algorithm*, Section 12.10 on page 257.

The resultant **ACL**, however it was obtained, is modified as listed below.

- If the originating request originated from DFS and is not authenticated (PAC *.authenticated* field = 0), an *unauth\_mask* entry is added to the **ACL** with a permissions field of all-ones.
 

**Note:** Currently, PACs presented to the LFS by DFS are always authenticated. Refer to *Obtaining a Principal's Identity*, Section 12.13 on page 259.
- The default realm in the **ACL** is set equal to the originating realm (step [4]).
- Any user or group entries in the **ACL** are changed to *foreign\_user* and *foreign\_group* ones. The realm field in these *foreign\_xxx* entries is set equal to the ORIGINAL realm that was associated with the **ACL**.
- Any *foreign\_user* or *foreign\_group* entries with a realm field equal to the NEW **ACL** realm (the originating realm) are changed to *user* and *group* ones, respectively.

6. If the new object is a directory that possesses an *Initial Object (Creation) ACL*, perform the same operations detailed in step [5] on THAT initial file **ACL** (instead of on the object's regular **ACL**).

If the new object is a directory that possesses an *Initial Container ACL*, perform the same operations detailed in step [5] on THAT initial directory **ACL** (instead of on the object's regular **ACL**).

## 12.12 User and Group Identities

(This is the saga of the **uuid\_t** versus the 32 bit UNIX ID.) When a DCE registry assigns accounts, it guarantees that a 32 bit UNIX **user** or **group** ID agrees with the first 32 bits of the corresponding DCE **uuid\_t**. Although this behavior could in theory change at some time in the future, it is unlikely to as long as the underlying operating systems upon which DCE and DFS are layered deal natively with these 32s-bit values instead of **uuid\_t**'s.

Because of this fact, the access right algorithm given later is only required to take into account the first 32-bits of *user* and *group* **uuid\_t** values. Likewise, only these first 32 bits are taken into account when comparing a **uuid\_t** against 32-bit file **oid** or **gid** values as well as fields from a non-DFS **ucred** authentication structure.

## 12.13 Algorithm for Obtaining a Principal Identity

During a DFS access, a principal's identity is contained in a PAC structure (**sec\_id\_pac\_t**). In an ideal world, each of the extended **vnode** operations would have been enhanced to include an argument of this type which could then, as discussed in a following section, be used during access checks. However (from the point of view of the LFS), these extended **vnode** operations continue to take only the standard **struct ucred** argument from which, in the DFS case, an appropriate PAC must be obtained.

Before going into how this can be achieved, it might be best to discuss what happens when a remote DFS request arrives at a server. Although the actual specifics are not in themselves relevant, a rough understanding of what is happening is useful.

- a. We initially start out with an rpc-handle, an **afsFid** (object identifier) and the call's arguments.
- b. A call to *px\_AdjustCell()* is made to force the *afsFid.Cell* field into a known state (*.high* = 0, *.low* = 1). (Presumably, certain code relies upon this behavior.)
- c. A call to *fshs\_InqContext()* is made in order to obtain the context (identity, and so forth) of the remote caller. This is followed by a call to *fshs\_GetPrincipal()* to obtain information regarding the principal at the remote site, including a PAC (**sec\_id\_pac\_t**) and an appropriate **ucred** structure.
- d. A series of *xcred\_xxx()* calls are made in order to embed this **ucred** and PAC within an extended credentials **xcred** structure and also plant within the **ucred** a PAG (Process Authentication Group) field which can later be used to "track down" this **xcred** and, hence, the PAC.
- e. A call is made to *volreg\_LookupExtended()* with the object **afsFid** in order to locate the exported fileset in question (indicated by the *afsFid.Volume* field). This call returns a pointer to a **struct volume**. After this volume is held (*vol\_hold()*), the **vnode** of the desired file is obtained by a call to *vol\_vget()*.

- f. A call is made into the file system to perform the operation via a *VOPX\_xxx()* call. From within the file system code, the path back to the PAC structure required for access checks is the PAG which was left in the **ucred** structure.
- g. The volume and vnode are released via *vol\_rele()* and *vn\_rele()* calls. (The principal and context structures are released as well.)

From within the LFS (step (f) above), the desired PAC is obtained as follows. This description is presented as actual C code since it is so dependent upon calls into the DFS portion of the system.

**Note:** The precise declarations for these *xcred\_xxx()* calls can be found in *DCE DFS Credential Design*, Section 17.1 on page 469.

```

/* Call xcred_UcredToXcred to use a PAG embedded within
 * the ucred, if any, to locate an extended credential
 * structure (xcred) which will eventually lead to the
 * PAC. It takes a ucred argument and, if successful,
 * returns a pointer to an xcred structure.
 */
rtn_code = xcred_UcredToXcred(&ucred, &xcred_ptr);
if (rtn_code == 0)
    /* Call xcred_GetProp to actually obtain the PAC.
     * It takes an xcred structure argument and, if
     * successful, returns a pointer to a PAC
     * (* sec_id_pac_t) along with the size of the PAC
     * - which we don't care about.
     */
    rtn_code = xcred_GetProp(xcred_ptr, "DCE_PAC", 7,
                            &pac_ptr, &pac_size);

if (rtn_code == 0) {
    /* We have a PAC: execute the access check
     * detailed below with with this.
     */
}
else {
    /* There is no PAC. Presumably, this is a non-DFS or
     * local access. Either execute an access check
     * algorithm optimized to work off a ucred or "create"
     * a PAC from the ucred via the rules given below and
     * execute the PAC-based access check algorithms. In
     * the latter case, the space for such a PAC (a
     * sec_id_pac_t along with a group list) is allocated
     * and freed in a file system dependent manner.
     */
}

```

If the DFS request in question is unauthenticated, the PAC returned by *xcred\_GetProp()* is set to an unauthenticated identity as shown below. Because even these PACs appear to be authenticated, normal **ACL** entries can be established to match unauthenticated users.

<i>.authenticated</i>	non-zero (true)
<i>.principal.uid</i>	-2 (1st 32 bits: used in access check)

<i>.group.uuid</i>	-2 (1st 32 bits: used in access check)
<i>.realm.uuid</i>	-2 (1st 32 bits)
<i>remaining fields</i>	set to 0

## 12.14 Algorithm for Generating PAC from Ucred Structure

As was mentioned above, (Section 12.13 on page 259), it may be necessary during a on-DFS access to synthesize a PAC from the contents of a standard ucred structure. The simplest way to accomplish this is to start out with both a **sec\_id\_pac\_t** PAC structure and a maximum sized (`OSI_MAXGROUPS_1`) **sec\_id\_t** group list (array of **sec\_id\_t**, that is). The PAC is initialized as follows with fields not mentioned set to 0).

<i>.pac_type</i>	= <i>sec_id_pac_format_v1</i> (= 0)
<i>.authenticated</i>	= true (actually, to, -1)
<i>.realm.uuid</i>	= LOCAL-REALM
<i>.principal.uuid.time_low</i>	= <b>uid</b> from <b>ucred</b> (by <i>osi_GetUID()</i> call)
<i>.group.group.uuid.time_low</i>	= <b>gid</b> from <b>ucred</b> (by <i>osi_GetGID()</i> call)
<i>.num_groups</i>	= number of groups in the ucred group list (via <i>osi_GetNGroups()</i> )
<i>.groups</i>	= pointer to the <b>sec_id_t</b> group array

For each of these groups, the *.groups[]* entry is zeroed and then filled in as follows:

<i>.groups[i].uuid.time_low</i>	= <i>group-id</i> from <b>ucred</b> <i>.cr_groups[i]</i>
---------------------------------	--

Note that a 32-bit UNIX ID-to-**uuid\_t** conversion is performed via a copy into the first 32-bits of the **uuid\_t**.

**Note:** Although this initialization could also be accomplished via a call to the DFS function *dacl\_PacFromUcred()*, a manual description is presented above in the interest of minimizing the actual interactions between an LFS and DFS.

## 12.15 perm\_control Access Right

Operations such as *chmod()* which, would in a traditional UFS file system, require the caller to own a file in question, instead require **perm\_control** access. The owner of a file will always have this right.

## 12.16 Access Rights Algorithm

The access check algorithm presented below uses the following parameters.

<b>desired-access</b>	A <b>permset_t</b> detailing the desired access.
<b>file-mode-bits</b>	The object's mode protection bits.
<b>file-oid</b>	The object's owner-ID value (32-bit UNIX <b>uid</b> ).
<b>file-gid</b>	The object's group-ID value (32-bit UNIX <b>gid</b> ).
	When the file-oid or file-gid are compared against a <b>uuid_t</b> below, the comparison actually occurs against the first 32-bits of the <b>uuid_t</b> .
<b>ucred</b>	The <b>ucred</b> structure of the principal on whose behalf the request is being made.
<b>ACL</b>	The object's <b>ACL</b> , if any.
<b>PAC</b>	The PAC identifying the principal on whose behalf the request is being made. For a DFS access, a PAC is always available (as discussed above). For non-DFS access (local, NFS access), the following algorithm assumes that a PAC is built from the available <b>ucred</b> structure as discussed above. Alternatively, one could provide two bodies of code; one working off of a PAC (DFS access) and another working off of a <b>ucred</b> (non-DFS access).  Associated with this PAC is an indication of whether this is a DFS access or a non-DFS access.  <b>Note:</b> A PAC is only available for requests that arrive via DFS. In the situation where a local user has performed a <b>dce_login</b> , there is no PAC that an LFS (or DFS) has access to.

For the purpose of user and group comparisons, only the first 32-bits of **uuid\_t** values need to be considered. This is the size of **file-oid** and **file-id** fields as well as the that of the relevant **ucred** fields that we started with in the event of a non-DFS access. All realm comparisons are performed on the full **uuid\_t**.

**Note:** This is the behavior of Episode LFS.

Several of the steps below test for an authenticated user. Although they are present in Episode, they are never triggered since PACs presented by DFS to the LFS are always (currently, at least) authenticated. Refer to *Algorithm for Obtaining a Principal's Identity*, Section 12.13 on page 259, and also *ACL Entry Types for Unauthenticated Users*, Section 8.2.3 on page 157.

The algorithm shown below can, in general, be executed in one of two situations:

- The set of rights that the caller has to an object is being computed (as returned by *vn\_getattr()* or *vol\_getattr()* in the *.callerAccess* field).
- A determination is being made as to whether the caller has a specific access to an object (as performed by the *vn\_access()*, *vn\_lookup()*, types of vnode operations)

In the Episode implementation, several of the individual steps differ depending upon which of the two cases is being handled. These are called out below.

If there is no object **ACL**, the rights that the principal has to the object are determined as in Section 12.16.1 on page 263.

**Note:** If the implementation were willing to construct an **ACL** in this case, the "there is an **ACL**" algorithm in Section 12.16.2 on page 263 could presumably be followed.



In either case, the desired access can finally be compared against the granted rights computed in Section 12.16.2 to determine if the access is allowed.

### 12.16.1 No Object ACL Exists Rights Algorithm

1. If the principal in question is root, full rights are granted: **perm\_control**, **perm\_read**, **perm\_write** and **perm\_execute**. If the object in question is a directory, **perm\_insert** and **perm\_delete** rights are granted as well. Root access is indicated by an all-zero *.principal.uid* field in an authenticated PAC (in the non-DFS situation, from a **uid** of 0). Continue with step [7].

**Note:** Episode doesn't perform this step in the case of a *vn* (or *vol\_getattr()*) operation obtaining a remote client's allowed access (*.callerAccess*) to an object.

2. If the PAC *.authenticated* field is 0 (an unauthenticated DFS access), no rights are granted. Exit the algorithm.
3. If the *.realm.uid* field in the PAC differs from the LOCAL-REALM (a DFS request), no rights are granted. Exit the algorithm.
4. If the PAC *.principal.uid* field matches the *file-oid* field, the **owner** bits in the file-mode-bits are granted. Additionally, **perm\_control** access is granted as well. Continue at step [7].
5. If the *file-gid* field matches either the *.group.uid* field or one of the *.groups[].uid* fields in the PAC, the *group* bits in the file-mode-bits are granted. Continue at step [7].
6. If all else fails, the *other* bits in the file-mode-bits is granted.

**Note:** In the above UNIX mode-to-**permset\_t** rights conversion, a straight forward **read-write-execute** to **perm\_read-write-execute** mapping is performed.

7. If the object is a directory and the **perm\_write** right has been granted, the **perm\_insert** and **perm\_delete** ones are as well.
8. Continue with step [k] in the ACL case algorithm in Section 12.16.2 (the next section).

### 12.16.2 Object ACL Exists Rights Algorithm

If there is an object ACL, the rights are determined instead of in the previous section, as follows.

- a. If the PAC *.realm.uid* is all zeroes, the steps below behave as if it were equal to the default **realm** field from the **ACL**.
- b. If the principal in question is root, full rights are granted: **perm\_control**, **perm\_read**, **perm\_write** and **perm\_execute**. If the object in question is a directory, **perm\_insert** and **perm\_delete** rights are granted as well. Root access is indicated by an all-zero *.principal.uid* field in an authenticated PAC (in the non-DFS situation, from a **uid** of 0). Continue with step [l] to perform a Read-only (RO) filesset check.
- c.

**Note:** This step starts *User* entry checking. It matches from the **file-uid**

If the PAC *.realm.uid* field agrees with the default **realm** from the **ACL** and the PAC *.principal.uid* field matches the **file-uid** (or **oid**), the rights from the *user\_obj* entry are granted and continue with step [k].

As discussed earlier, an implementation may need to replace the **ACL perm\_read**, **perm\_write** and **perm\_execute** rights with the corresponding bits from the file-mode bits.

- d. If the PAC *.realm.uuid* field agrees with the default realm from the **ACL**:
- The **ACL** is searched for an entry of type *user* whose user **uuid** matches the PAC *.principal.uuid* field.
- Otherwise, the **ACL** is searched for an entry of type *foreign\_user* whose user **uuid** matches the PAC *.principal.uuid* field and whose **realm uuid** matches the PAC *.realm.uuid* field.
- If such an entry is found, the rights in it are granted and continue with step [k].
- e.
- Note:** This step starts *Group* entry checking. It matches from the **file-gid**
- If the PAC *.realm.uuid* field agrees with the default **realm** from the **ACL**:
- If the **file-gid** matches the PAC *.group.uuid* field or one of the PAC *.groups[].id.value* fields, the rights from the *group\_obj* entry are granted and continue with step [k].
- As discussed earlier, an implementation may need to replace the **ACL perm\_read**, **perm\_write** and **perm\_execute** rights with the corresponding bits from the file-mode bits.
- f.
- Note:** This step matches against some **ACL group** entry.
- If the PAC *.realm.uuid* field agrees with the default **realm** from the **ACL**:
- The **ACL** is searched for entries of type *group* whose *.group.uuid* matches either the *.group.uuid* or one of the *.groups[i].id.uuid* fields in the PAC.
- Otherwise, the **ACL** is searched for entries of type *foreign\_group* whose *.realm.uuid* matches the PAC *.realm.uuid* and whose *.group.uuid* matches either the *.group.uuid* or one of the *.groups[i].id.uuid* fields in the PAC.
- The rights from ALL such **ACL** entries found are "or'd" together to form the granted access.
- If at least one such **ACL** entry is found, continue with step [k].
- g.
- Note:** This step matches against the **ACL other** entry.
- If the PAC *.realm.uuid* field agrees with the default **realm** from the **ACL**, the rights from the *other\_obj* entry are granted and continue with step [k].
- As discussed earlier, an implementation may need to replace the **ACL perm\_read**, **perm\_write** and **perm\_execute** rights with the corresponding bits from the file-mode bits.
- h. The **ACL** is searched for an entry of the type *foreign\_other* whose *.realm.uuid* matches the PAC *.realm.uuid*. If such an entry is found, the rights in it are granted and continue with step [k].
- i. If an **ACL** entry of the type *any\_other* exists, the rights in it are granted and continue with step [k].
- j. If this step is reached, no rights are granted thus far.
- k.
- Note:** This step begins finishing up, and is arrived at from several previous places in the processing algorithm.
- If the PAC *.authenticated* field is zero (0), this is an unauthenticated request and the rights computed thus far are masked with the *any\_other* entry if present, or if not, a mask of zero

is used.

**Note:** See *ACL Entry Types for Unauthenticated Users*, Section 8.2.3 on page 157. If there is an *unauth\_mask* entry in the **ACL**, it should be ignored. If the version of the system is DCE 1.1 or newer, these entries should be removed.

- l. If the object resides on a read-only fileset and is not a device file, the **perm\_write** permission is removed if it has been granted thus far.

**Note:** This allows device files to reside on read-only (RO) filesets.

**Note:** Episode doesn't perform this step in the case of a *vn\_getattr()* or *vol\_getattr()* operation obtaining a remote client's access (.callerAccess) to a file.

- m. If the object is a directory:

The **perm\_insert** and **perm\_delete** rights are removed (if present) if **perm\_write** is not present.

Otherwise, the **perm\_insert** and **perm\_delete** rights are not meaningful and can be removed.

- n.

**Note:** This is for DFS-SYSADMIN-GROUP processing.

If the PAC *.realm.uuid* field equals the LOCAL-REALM (the request originates from the local cell) and the DFS\_SYSADMIN\_GROUP is equal to one of the groups indicated in the PAC *.group.uuid* or *.groups[].uuid* fields, the **perm\_control** right is added as well.

## 12.17 mask\_obj ACL Entry Algorithm Impact

As was mentioned earlier, no mention was made of the *mask\_obj* ACL entry in the above algorithms since the functionality provided by it has been dropped from the POSIX ACL proposals. The impacts that this *mask\_obj* field would have on the above algorithms are listed below for a number of reasons.

1. Episode and DCE in general, implements to an older POSIX specification that does contain the *mask\_obj* functionality.
2. POSIX has not as yet accepted the new proposal.
3. Complete interoperability with Episode and the rest of DCE might require that an LFS implement this functionality.

The general idea is that during a *chmod()*, the group mode bits are used as a mask (the *mask\_obj* ACL entry) that weakens all *user* and *group* type **ACL entries**. In order to accommodate this *mask\_obj* functionality, the above algorithms are modified as follows.

- When an **ACL** is being validated during a *vn\_setacl()* operation, an **ACL** is considered corrupt if it there are any entries of type *user*, *foreign\_user*, *group* or *foreign\_group* and there is not a *mask\_obj* entry.

**Note:** Episode also requires a *mask\_obj* entry if there are any *foreign\_other* entries.

- During a *chmod()* operation to an object with an **ACL**, the **read-write-execute** group mode bits being set are copied into the *mask\_obj* entry if it exists. (If there is no *mask\_obj* entry, they are copied into the *group\_obj* entry instead.)

**Note:** As discussed earlier, an implementation may leave the **ACL** alone and instead force this coercion at the time of a latter *vn\_getacl()* or access rights check.

- During a *vn\_setacl()* operation, the **read-write-execute** rights from the **ACL** *mask\_obj* entry are copied into the group file mode bits. (If there is no *mask\_obj* entry, rights from the *group\_obj* entry are instead copied into the group file mode bits.)
- During the access rights checking algorithm presented earlier, the rights in a *mask\_obj* entry (if present) are used to weaken ("and operation") the rights in any **ACL** entries of type *user*, *group\_obj*, *group*, *foreign\_user* and *foreign\_group*.

Additionally, the presence of an *mask\_obj* entry causes step [f] in the previous section, Section 12.16.2 on page 263, to be performed even if step [e] triggers a hit. In this case, the rights constructed by [f] and [e] are "or'd" together to construct the actual rights granted.

## 12.18 Miscellaneous Topics and Suggestions

1. If **ACLs** become prevalent, performance and disk utilization concerns will likely argue for **ACL** sharing (both on disk and in memory). As pointed out, most objects created by principals in a *foreign* realm have an **ACL**.
2. The DFS *vn\_getattr()* and *vol\_getattr()* requests are required to return the **minimum rights** granted by an object's **ACL** to ANY user (authenticated or not). Performance concerns might warrant "pre-computing" this value and storing it in a dedicated **ACL** or anode (inode) field.

## 13.1 Primitive Data Types

The following are the primitive data types needed by a **DCE LFS**.

- afsFid** A DFS file identifier. See Section 4.11.8 on page 60 for its definition.
- afsHyper** A 64-bit identifier for objects such as cells, volumes, and so on. See Section 4.11.7 on page 60 for its definition.
- afsTimeval** A time value. See Section 4.11.6 on page 60 for its definition.
- afsUUID** Equivalent to an `uuid_t`. See Section 4.11 on page 59 of Chapter 4 on page 57 for its definition.
- permissions\_t** DCE permissions are stored within a 32-bit word. Within this long, DFS is cognizant of the rights listed below.

```
#typedef u_int32 permissions_t

#define perm_read      0x1 /* Read rights */
#define perm_write     0x2 /* Write rights */
#define perm_execute   0x4 /* file: Execute rights */
                        /* directory: Search rights */
#define perm_control   0x8 /* the right perform various */
                        /* control operations (for example, */
                        /* changing the mode bits) */
#define perm_insert    0x10 /* right to add dir entries */
#define perm_delete    0x20 /* right to remove dir entries */
```

The remaining bits, although not interpreted by DFS, are fully preserved in **ACLs**.

- sec\_id\_pac\_t** The PAC (Privilege Attribute Certificate) identifying an authenticated principal.

```
struct sec_id_t {
    uuid_t          uuid;
    idl_char        *name;
};
struct sec_id_foreign_t {
    sec_id_t        id;
    sec_id_t        realm;
};

struct sec_id_pac_t {
    sec_id_pac_format_t pac_type;
    long              authenticated;
    sec_id_t          realm;
    sec_id_t          principal;
    sec_id_t          group;
    short             num_groups;
```

```

        short                num_foreign_groups;
        sec_id_t             *groups;
        sec_id_foreign_t     *foreign_groups;
};

```

The following are the definitions of the **sec\_id\_pac\_t** fields.

*.pac\_type*            A format type which can be safely ignored (currently, always *sec\_id\_pac\_format\_v1* which is 0).

*.authenticated*      Non-zero implies that the principal is authenticated.

*.realm*                The administrative realm with which the user and group identities in this PAC are associated.

*.principal*            Principal's identifier.

*.group*                Principal's primary group.

*.num\_groups*          Number of entries in *.groups*.

*.num\_foreign\_groups*    Number of entries in *.foreign\_groups*.

*.groups*                List of groups associated with the default *.realm* field above to which the principal is a member; each entry consists of a group **uuid\_t** and an ignored name.

*.foreign\_groups*      List of groups associated with other realms to which the principal is a member; each entry consists of both a realm (**uuid\_t** and ignored name) and group (**uuid\_t** and ignored name).

## **uuid\_t**

The globally unique name for entities (objects, interfaces, principals, and so on) known by DCE. The DCE variant of the Universal Unique Identifier is listed in Appendix A of the **X/Open DCE: Remote Procedure Call** specification which lists it in the form of table entries. Its code structure is listed here for convenience, and is as follows:

```

typedef struct uuid_t {
    u_int32         time_low;
    u_int16         time_mid;
    u_int16         time_hi_and_version;
    unsigned char  clock_seq_hi_and_reserved;
    unsigned char  clock_seq_low;
    unsigned char  node[6];
};

```

## 13.2 Aggregates and Aggregate Registry Data Types

```
#include <dcedfs/osi.h>
#include <dcedfs/common_def.h>
#include <dcedfs/lock.h>

#define MAX_AGGRNAME 64 /* XXXX */
```

### 13.2.1 Aggregate Static Status

Aggregate's general static status information: the piece that sits in the **struct aggr** defined in Section 13.2.6 on page 270. These fields are all maintained by software above the LFS and therefore do not need to be stored on disk. For instance, many are taken from or are derived from information in the `dfstab` table.

```
struct ag_status_st {
    long aggrId; /* Aggregate Id */
    long nVolumes; /* Number of attached filesets */
    long spare1, spare2, spare3, spare4; /* Some spares */
    dev_t device; /* major/minor aggr device number */
    char aggrName[MAX_AGGRNAME]; /* Aggregate name */
    char devName[32]; /* Aggr device's name */
    /*(One of the AG_TYPE_xxx codes)*/
    char type; /* UFS, PFS, AIX3, etc. */
    char spares[16]; /* More spares */
};
```

### 13.2.2 Aggregate Dynamic Status

The dynamic aggregate status information: the piece that is computed on the fly. Fields within this structure are maintained (or supplied) by the file system dependent code.

The *minFree* field is the number of blocks on the aggregate held in reserve for the exclusive use of root. If this field is 0 (as it is for Episode), then the *totalUsable* field will be equal to the *blocks* field.

```
struct ag_status_dy {
    long blocks; /* # blks (of size fragsize) in aggr */
    long blocksize; /* block size of aggregate */
    long totalUsable; /* total available 1K blocks on aggr */
    /* available to non-root users */
    long realFree; /* free 1K blocks avail to non-root */
    long minFree; /* min free 1K blocks */
    long fragsize; /* fragment size of aggr */
    long spares[7]; /* Some spares */
};
```

### 13.2.3 Aggregate Status

This structure holds the status of an aggregate.

```
struct ag_status {
    struct ag_status_st    ag_st;
    struct ag_status_dy ag_dy;
};
```

### 13.2.4 Valid Aggregate Types

The following are valid aggregate types (`ag_status.type`). See Section 13.2.1 on page 269.

```
#define    AG_TYPE_UNKNOWN    0
#define    AG_TYPE_UFS        1
#define    AG_TYPE_EPI        2
#define    AG_TYPE_AIX3      3
#define    AG_TYPE_VXFS       4

#define    MAX_AG_TYPE        5 /* size of array of agops structures */
```

### 13.2.5 Adding of New Filesystem Types

Macros to make addition of new file system types easier in the future.

```
#define AG_TYPE_SUPPORTS_EFS(ag)  ((ag) == AG_TYPE_EPI) \
    || ((ag) == AG_TYPE_VXFS)

#define AG_TYPE_TO_STR(ag)  ((ag) == AG_TYPE_EPI ? "LFS" : \
    (ag) == AG_TYPE_VXFS ? "VXFS" : \
    "Non-LFS")
```

### 13.2.6 Aggregate Structure

The `aggr` structure, below, is a memory resident structure maintained for each aggregate. In general, it is maintained by code above the LFS and is passed as an argument to each of the aggregate operations.

```
struct aggr {
    struct aggr *a_next;          /* next on linked list */
    struct aggrOps *a_aggrOpsp; /* Ops on aggregates */
    long a_states;                /* Aggregate states */
    long a_refCount;              /* Reference counter */
    osi_dlock_t a_lock;          /* Lock variable for the structure */
    struct ag_status_st a_stat_st; /* static part of aggregate's status */
    struct vnode *devvp;         /* special vnode for device */
    opaque a_fsDatap;           /* Pointer to File-system-dependent */
                                /* private data */
};
```



### 13.2.7 Aggregate Field Definitions as a\_\* Items

By defining all aggr's fields as a\_\* they can freely be moved around without changing the code all over the place where this is done.

```
#define a_aggrId      a_stat_st.aggrId
#define a_aggrName    a_stat_st.aggrName
#define a_devName     a_stat_st.devName
#define a_device      a_stat_st.device
#define a_type        a_stat_st.type
#define a_nVolumes    a_stat_st.nVolumes
```

The following fields are defined for compatibility.

```
#define a_blocks      a_status.blocks
#define a_blocksize   a_status.blocksize
#define a_totalUsable a_status.totalUsable
#define a_realFree    a_status.realFree
#define a_minFree     a_status.minFree
#define a_volId       a_status.volId
```

### 13.2.8 Aggregate States

The following state bits within *a\_states* in the **aggr** structure in Section 13.2.6 on page 270 are defined. The second, **AGGR\_DELETED**, is obsolete and no longer in use.

```
#define AGGR_EXPORTED 1 /* Exported aggregate */
#define AGGR_DELETED 2 /* Aggregate deleted */
```

### 13.2.9 Fileset Creation Flags

The following are flags for **volCreate()** in the **aggrops** structure in Section 13.2.12 on page 272.

```
#define AGGR_CREATE_ROOT      (0x1)
#define AGGR_CREATE_VALID_FLAGS (0x1)
```

### 13.2.10 Aggregate Attach Flags

The following are flags for **ag\_attach()** in the **aggrops** structure in Section 13.2.12 on page 272.

```
#define AGGR_ATTACH_NOEXPORT (0x1)
#define AGGR_ATTACH_VALID_FLAGS (0x1)
```

### 13.2.11 Aggregate Sync Flags

The following are flags for **ag\_sync()** in the **aggrops** structure in Section 13.2.12 on page 272.

```
#define AGGR_SYNC_FILESYS      1
#define AGGR_SYNC_COMMITMETA  2
#define AGGR_SYNC_COMMITALL    3
```

### 13.2.12 Aggregate Operations Vector

The following is the aggregate operations vector.

```

struct aggrops {
    int (*ag_hold)();
    int (*ag_rele)();
    int (*ag_lock)();
    int (*ag_unlock)();
    int (*ag_stat)();
    int (*ag_volCreate)();
    int (*ag_volInfo)();
    int (*ag_detach)();
    int (*ag_attach)(); /* NOT an AG_ATTACH; Accessed explicitly */
    int (*ag_sync)();
};

```

### 13.2.13 Exported Aggregate Registry Items

The following are the exported globals and functions for dealing with the *Aggregate Registry*.

```

extern struct aggr *ag_root;
extern osi_dlock_t ag_lock;
extern long ag_attached;

#ifdef KERNEL
extern struct icl_set *xops_iclSetp;
#endif /* KERNEL */

extern void ag_Init(void);
extern int ag_PutAggr(struct aggr *), ag_RemoveAggr();
extern struct aggr *ag_GetAggrByDev(dev_t);
extern struct aggr *ag_GetAggr(long);
extern int ag_NewAggr(
    char *aggrNamep,
    long aggrId,
    char *aggrDevnamep,
    char aggrType,
    struct aggrops *aggrOpsp,
    dev_t dev,
    opaque fsdata,
    struct aggr **aggrpp,
    unsigned flags
);

struct aggrops *agOpvec[MAX_AG_TYPE];

```

### 13.2.14 Vnode Ops Classification

The following are broad classifications of the vnode ops types.

```
#define VNOP_TYPE_NOOP          5 /* noop */
#define VNOP_TYPE_READONLY     10 /* read only vnode op */
#define VNOP_TYPE_READWRITE    15 /* read write vnode op */
#define VNOP_TYPE_INVALID      30
```

### 13.2.15 Concurrency for vol\_stat\_st

The following define the concurrency levels for concurrency field in vol\_stat\_st. They are used to control concurrency between volops and vnops.

```
#define VOL_CONCUR_ALLOPS      1
#define VOL_CONCUR_READONLY    10
#define VOL_CONCUR_NOOPS       20
```

The following is a macro to determine if a vnode op can execute through, even if its fileset is busy.

```
#define VOL_VNOP_COMPAT(concurr, type) (((concurr == VOL_CONCUR_ALLOPS) || \
    (type == VNOP_TYPE_NOOP) || \
    ((concurr == VOL_CONCUR_READONLY) && \
    (type == VNOP_TYPE_READONLY)))) ? \
    1 : 0)
```

### 13.2.16 Aggregate Table Entry Format

The following defines the /opt/dcelocal dfstab file and the sizes of the strings it contains.

```
#define ASTAB_INFIX      "/var/dfs/"
#define ASTAB_SFX        "dfstab"

/* Sizes of strings */

#define ASTABSIZE_SPEC   512
#define ASTABSIZE_DIR    512
#define ASTABSIZE_TYPE   64
```

This is the actual **dfstab** table entry structure, defined as a **struct astab**. This structure identifies a **DCE LFS** aggregate that can be exported. For non-LFS partitions only, an additional entry is present - the **Fileset ID**. See Section 13.2.17 on page 274, below for the definition. The **OSF DCE DFS Administration Guide and Reference** provides pertinent information about the **dfstab** file.

```
struct astab {
    char    as_spec[ASTABSIZE_SPEC]; /* block special device name */
    char    as_aggrName[ASTABSIZE_DIR]; /* Aggregate name */
    char    as_type[ASTABSIZE_TYPE]; /* type of physical file system */
    u_long  as_aggrId; /* Aggregate Id */
};
```

### 13.2.17 UFS Entry Extra Data

```

/* File system extra data:  UFS version */

struct ufs_astab {
    hyper    uas_volId; /* 1st & only volume */
    char     uas_mountedon[ASTABSIZE_DIR]; /* name of mounted-on dir */
};

```

### 13.2.18 Values for Filesystem Type

```

#define     ASTABTYPE_UFS      "ufs"
#define     ASTABTYPE_EPI     "lfs"
#define     ASTABTYPE_AIX3    "aix3"
#define     ASTABTYPE_VXFS    "vxfs"

```

### 13.2.19 Adding New Filesystem Types

The following macros are for making the addition of new file system types easier than manual addition.

```

#define ASTABTYPE_SUPPORTS_EFS(atype) ((strcmp(atype, ASTABTYPE_EPI, \
    ASTABSIZE_TYPE) == 0) || (strcmp(atype, ASTABTYPE_VXFS, \
    ASTABSIZE_TYPE) \
    == 0))

#define ASTABTYPE_TO_AGTYPE(atype, agtype) \
    if (!strcmp(atype, ASTABTYPE_UFS, ASTABSIZE_TYPE)) \
        agtype = AG_TYPE_UFS; \
    else if (!strcmp(atype, ASTABTYPE_EPI, ASTABSIZE_TYPE)) \
        agtype = AG_TYPE_EPI; \
    else if (!strcmp(atype, ASTABTYPE_VXFS, ASTABSIZE_TYPE)) \
        agtype = AG_TYPE_VXFS; \
    else \
        agtype = AG_TYPE_UNKNOWN

```

## 13.3 Fileset Data Types

### 13.3.1 Define the Fileset Function Array

A description of the operations available within the fileset function array appears below. The array is an exported data type.

```

struct volumeops {
    /* per-fileset operations */
    int (*vol_hold)();
    int (*vol_rele)();
    int (*vol_lock)();
    int (*vol_unlock)();
    int (*vol_open)();
    int (*vol_seek)();
    int (*vol_tell)();
    int (*vol_scan)();
};

```

```

int (*vol_close)();
int (*vol_destroy)();          /* actually whole-fileset */
int (*vol_attach)();
int (*vol_detach)();
int (*vol_getstatus)();
int (*vol_setstatus)();

/* per-file operations */
int (*vol_create)();
int (*vol_read)();
int (*vol_write)();
int (*vol_truncate)();
int (*vol_delete)();
int (*vol_getattr)();
int (*vol_setattr)();
int (*vol_getacl)();
int (*vol_setacl)();

/* more whole-fileset operations */
int (*vol_clone)();
int (*vol_reclone)();
int (*vol_unclone)();

/* vnode lookup operations */
int (*vol_vget)();
int (*vol_root)();
int (*vol_isroot)();

/* more per-fileset operations */
int (*vol_getvv)();
int (*vol_setdystat)();
int (*vol_freedystat)();
int (*vol_setnewvid)();

/* another per-file operation */
int (*vol_copyacl)();

/* per-fileset operations */
int (*vol_concurr)();
int (*vol_swapids)();
int (*vol_sync)();
int (*vol_pushstatus)();

/* per-file operations */
int (*vol_readdir)();
int (*vol_appenddir)();

/* per-fileset operations */
int (*vol_bulksetstatus)();
int (*vol_getzlc)();

/* another per-file operation */

```

```

int (*vol_getnextholes)();

/* another per-fileset operation */
int (*vol_deplete)();
};

```

### 13.3.2 Volume Operations Definitions

The following are the definitions of the system calls for the volume operations in the **struct volumeops** Function Array. They can be found in the header file, `/src/file/xvolume/vol_init.h`.

```

#define      VOLOP_OPEN          4      /* Slots 0-3 are reserved */
#define      VOLOP_SEEK          5
#define      VOLOP_TELL          6
#define      VOLOP_SCAN          7
#define      VOLOP_CLOSE         8
#define      VOLOP_DESTROY       9
#define      VOLOP_GETSTATUS     10
#define      VOLOP_SETSTATUS     11
#define      VOLOP_CREATE        12
#define      VOLOP_READ           13
#define      VOLOP_WRITE         14
#define      VOLOP_TRUNCATE      15
#define      VOLOP_DELETE        16
#define      VOLOP_GETATTR       17
#define      VOLOP_SETATTR       18
#define      VOLOP_GETACL        19
#define      VOLOP_SETACL        20
#define      VOLOP_CLONE         21
#define      VOLOP_RECLONE       22
#define      VOLOP_VGET          23
#define      VOLOP_ROOT          24
#define      VOLOP_ISROOT        25
#define      VOLOP_UNCLONE       26
#define      VOLOP_FCLOSE        27      /* XXXX */
#define      VOLOP_SETVV         28
#define      VOLOP_SWAPVOLIDS    29
#define      VOLOP_COPYACL       30
#define      VOLOP_AGOPEN        31
#define      VOLOP_SYNC          32
#define      VOLOP_PUSHSTATUS    33

#define      VOLOP_PROBE         34
#define      VOLOP_LOCK          35
#define      VOLOP_UNLOCK        36

#define      VOLOP_READDIR       37
#define      VOLOP_APPENDDIR     38

#define      VOLOP_BULKSETSTATUS 39
#define      VOLOP_GETZLC        40
#define      VOLOP_GETNEXTHOLES 41
#define      VOLOP_DEplete       42

```

### 13.3.3 The Fileset NextHole Structure

A description of the operations available within the fileset NextHole structure appears below. It can be found in the header file, `/src/file/xvolume/vol_init.h`.

```

/* Parameter for VOLOP_GETNEXTHOLES */
#define VOLHOLE_MAX_HOLES 10

struct vol_NextHole {
    struct afsHyper startPoint; /* Return all holes starting here or later */
                                /* Zero this to start */
    unsigned long flags;
    unsigned long outCount; /* count of returned holes */
    unsigned long spare1;
    unsigned long spare2;
    struct vol_holeDesc {
        struct afsHyper holeStart; /* byte address in file where hole starts */
        struct afsHyper holeLen; /* byte count of the hole */
    } holes[VOLHOLE_MAX_HOLES];
};
/* IN flags: */
/* none yet */
/* OUT flags: */
#define VOLHOLE_FLAG_LAST 0x1 /* no more holes in file past these */

```

### 13.3.4 Define the vol\_stat\_st Structure

The following structure defines the static portion of a fileset's status.

```

struct vol_stat_st {
    afsHyper volId; /* on-disk */
    afsHyper parentId; /* on-disk */
    afsTimeval cloneTime; /* on-disk */
                                /* when fileset was made via clone, reclone */
    afsTimeval vvCurrentTime; /* when lazy replica's VV was known current */
    afsTimeval vvPingCurrentTime; /* when last tried for VV of lazy replica */
    long type;
    u_long accStatus; /* volops we'll perform in this trans (bit mask) */
    u_long accError; /* error to return for incompatible vnode ops */
    u_long states; /* on-disk lots of individual VOL_XXX bits */
    u_long reclaimDally; /* for the R/W backing some lazy-rep R/O */
    long tokenTimeout;
    long activeVnops; /* count of active vnode ops */
    long volMoveTimeout;
    long procID; /* pid of process making fileset busy */
    long spare5;
    long spare6;
    long spare7;
    long spare8;
    long spare9;
    char volName[VOLNAME_SIZE]; /* on-disk (VOLNAME_SIZE is 112) */
    unsigned char concurrency; /* level of concurrency for vnode ops */
    char cspares[15];

```

};

The portion of a fileset's status which (for the most part) is changed by **DCE DFS** and not by the LFS. Fields tagged with an `"/* on-disk .. */` must be stored in permanent, on-disk storage. Except as discussed in the chapter *Fileset (Volume) Operations Interface* or in the chapter *Aggregate Operations Interface*, these status fields are ignored by the LFS layer. Several of the fields appear to be obsolete (and unused). The above fields have the following meanings.

<code>.volId</code>	The ID of this fileset.
<code>.parentId</code>	For a <code>.backing</code> (VOL_BACKUP) or <code>.readonly</code> (VOL_READONLY) fileset, the ID of the fileset that it was originally built from (in other words, the "read-write" member of a fileset group). For the RW fileset, this field is set to <b>0,0</b> .
<code>.cloneTime</code>	The time at which a cloned fileset (the backing half) is initially created or re-cloned.
<code>.vvCurrentTime</code>	Time that the lazy replica's volume version was known to be current.
<code>.vvPingCurrentTime</code>	Time when DFS tried to last obtain the lazy replica's volume version.
<code>.type</code>	Set to one of VOL_RW or VOL_READONLY; see Section 13.3.9 on page 282 for the encodings of these.
<code>.accStatus</code>	While a fileset is open, set by DFS to the open-type argument that was passed to the <code>vol_open()</code> call.
<code>.accError</code>	While a fileset is open, set by DFS to the error code to be returned to any conflicting opens that fail.
<code>.states</code>	Fileset state bits. See Section 13.3.9 on page 282. Controls the amount of time to wait (in other words, daily) before actually removing deleted files on read-only replica filesets.
<code>.tokenTimeout</code>	Set to the volume version token expiration time by DFS (at the same time that the <code>.vvCurrentTime</code> and <code>.vvPingCurrentTime</code> fields are set).
<code>.activeVnops</code>	The number of active vnode operations (from either the protocol exporter or the vfs and vnode glue code) currently operating on this fileset.
<code>.volMoveTimeout</code>	For filesets actively involved in a move operation (either the source or destination), the estimated time of completion. (Used to recover from failures during the operation.)
<code>.procID</code>	The identity of the process (thread, actually) that has the fileset open.
<code>.volName</code>	The name of the fileset.
<code>.concurrency</code>	For an open fileset, one of the VOL_CONCURR_xxx values determining the set of vnode operations that are permitted to proceed concurrently.
<code>.states</code>	The meanings of the state bits (VOL_xxx flags) within this field are listed in Section 13.3.4 on page 277.



### 13.3.5 Define the vol\_stat\_dy Structure

The following structure defines the dynamic portion of a fileset's status.

```
struct vol_stat_dy {
    afsTimeval creationDate; /* on-disk */
                                /* when this fileset was created */
    afsTimeval updateDate; /* on-disk */
                                /* when any update happened in this fileset */
    afsTimeval accessDate; /* on-disk */
                                /* when any access was made to this fileset */
    afsTimeval backupDate; /* spare */
    afsTimeval copyDate; /* on-disk */
                                /* when dump was made that was restored here */
    afsHyper volversion; /* on-disk */
    afsHyper backupId; /* on-disk */
    afsHyper cloneId; /* on-disk */
    afsHyper llBackId; /* on-disk */
                                /* low-level backing ID */
    afsHyper llFwdId; /* on-disk */
                                /* low-level forward ID */
    afsHyper allocLimit; /* on-disk */
    afsHyper allocUsage; /* on-disk */
    afsHyper visQuotaLimit; /* on-disk */
    afsHyper visQuotaUsage; /* on-disk */
    long fileCount;
    long minQuota;
    long owner;
    long unique; /* on-disk */
                                /* necessary if ``version'' is afsHyper? */
    long index; /* Redundant */
    long rwIndex;
    long backupIndex;
    long parentIndex;
    long cloneIndex;
    long nodeMax;
    long aggrId;
    long spare2;
    long spare3;
    long spare4;
    long spare5;
    long spare6;
    u_long tag;
    u_long msgLen;
    char statusMsg[128]; /* on-disk */
    char cspares[16];
};
```

The portion of a fileset's status which (for the most part) is changed by **DCE DFS** and not by the LFS. Fields tagged with an `"/>* on-disk .. */` must be stored in permanent, on-disk storage. Except as discussed in the chapter *Fileset (Volume) Operations Interface* or in the chapter *Aggregate Operations Interface*, these status fields are ignored by the LFS layer. Several of the fields appear to be obsolete (and unused).

**Note:** For fields holding time values, Transarc's Episode only stores the .sec portions of these times on-disk. The .usec field is set to zero on get info requests.

The above fields have the following meanings.

<i>.creationDate</i>	The time, set by the LFS, at which this fileset was created.
<i>.updateDate</i>	Set by the LFS to the current time whenever a fileset's status is changed. This includes the creation and deletion of files as well as whenever the volume version is incremented.  <b>Note:</b> For performance reasons, there is freedom in the granularity of updates to this field (e.g.,: every second or couple of seconds). Therefore, this field need not be absolutely accurate.
<i>.accessDate</i>	Although the original intent of this field is unclear, under Transarc's Episode it is updated at exactly the same instants that the above <i>.updateDate</i> field is. (Exception: at fileset creation, via <i>ag_volCreate()</i> , these two fields can initially be set to two different values, although they'll synchronize themselves almost immediately.) Unless or until the meaning of these fields changes in the future, it should be possible to implement them as a single on-disk field.  <b>Note:</b> Transarc considers it unlikely that the <i>.accessDate</i> field will ever be given a distinct meaning.
<i>.backupDate</i>	Appears to be currently unused.
<i>.copyDate</i>	When a fileset is restored, this field is set to the time at which the dump (or backup) from which it is restored was originally taken.
<i>.volversion</i>	This field is incremented by the LFS whenever any change is made to any file on the fileset (to either the status or data portion of the file).  <b>Note:</b> It appears this field is not updated for fileset status updates.
<i>.backupId</i>	The ID of any .backup (VOL_BACKUP, created via "fts clone") fileset within the fileset group of which this one is a member.
<i>.cloneId</i>	The ID of any .readonly (VOL_READONLY, release or scheduled replicas) filesets within the fileset group of which this one is a member.
<i>.llBackId</i>	This field, and the one following are related. If non-zero, this field identifies the fileset backing that is backing this one.
<i>.llFwdId</i>	This field and the preceding one maintain the linked list that maintains the backing relationship between cloned filesets. If non-zero, this field identifies the immediate fileset that this one is backing. Refer to Section 15.3 on page 328 for further details if desired.

**Note:** The next four quota related fields are more thoroughly described in the section, *Fileset quotas*, Section 15.7 on page 333. When the Limit values are initially presented to the file system by a *vol\_setattr()* or *ag\_volCreate()* call, they should be rounded up to cover an integral number of internal "allocation units". The values returned by *vol\_getattr()* should reflect these rounded up values and can, therefore, differ from what was originally specified.

If, to a *vol\_setattr()* or *ag\_volCreate()* operation, a quota limit is supplied that is larger than the maximum value that can be represented internally, the maximum value that can be held internally should be used instead.

Refer to the section *Fileset Quotas*, Section 15.7 on page 333.

The following continue to define the meanings of the fields in the **struct vol\_stat\_dy** structure.

<i>.allocLimit</i>	The maximum size, in bytes, that a fileset's "allocated" usage can grow to.
<i>.allocUsage</i>	The "allocated" measure, in bytes, of the amount of disk space being used by the fileset.
<i>.visQuotaLimit</i>	The maximum size, in bytes, that a fileset's "visible" usage can grow to.
<i>.visQuotaUsage</i>	The "visible" measure, in bytes, of the amount of disk space being used by the fileset.
<i>.fileCount</i>	Currently unused.
<i>.minQuota</i>	Currently unused.
<i>.owner</i>	Currently unused.
<i>.unique</i>	Episode uses this per- fileset counter for the maintenance of inode generation numbers. Other LFS implementations may or may not use this field. Consult the section on <i>Inode Generation Numbers</i> , Section 15.8 on page 334, for further information.
<i>.index</i>	The index of this fileset on the aggregate. This value has the same interpretation and significance as that used by the aggregate <i>ag_vollInfo()</i> operation. See Section 14.4 on page 312 if desired.
<i>.rwIndex</i>	Currently unused.
<i>.backupIndex</i>	Currently unused.
<i>.cloneIndex</i>	Currently unused.
<i>.parentIndex</i>	This field is always returned as 0 by Transarc's Episode (although there appears to be latent support for it).  <b>Note:</b> It appears that this field was or is intended to provide the index (see the <i>.index</i> field above) of the fileset that is backing this one. (Not the index of the <i>vol_stat_st .parentID</i> fileset.)
<i>.nodeMax</i>	The canonical (based on VOL_ROOTINO) index of the highest inode (its index, that is) allocated within the fileset. (This field indicates the size of the dynamically allocated inode table.) There does not appear to be any use of this field.
<i>.aggrId</i>	The Id of the aggregate on which this fileset resides.
<i>.tag</i>	Currently unused.
<i>.msgLen</i>	Currently unused.
<i>.statusMsg</i>	A status message -- <i>not</i> necessarily null terminated.

### 13.3.6 Define the `vol_status` Structure

The following structure defines the combined static and dynamic portions of a fileset's status.

```
struct vol_status {
    struct vol_stat_st vol_st;
    struct vol_stat_dy vol_dy;
};
```

### 13.3.7 Define the `vol_statusDesc` Structure

The following structure defines the status description structure used by such operations as `vol_bulksetstatus()`.

```
#define VOL_MAX_BULKSETSTATUS    3

typedef struct vol_statusDesc {
    union {
        int          volDesc;
        struct volume *volp;
    }               vsd_volId;
    u_long          vsd_mask;
    u_long          vsd_spare;
    struct vol_status vsd_status;
} vol_statusDesc_t;
#define vsd_volDesc    vsd_volId.volDesc
#define vsd_volp      vsd_volId.volp
```

### 13.3.8 Define Transient Error

The following macros determine if an error is transient or persistent.

```
#define VOL_ERROR_IS_TRANSIENT(error) (error < VOLERR_TRANS_HIGHEST && \
    error > VOLERR_TRANS_LOWEST)
```

### 13.3.9 Volume States (for `vol_stat_st`)

The following general status flags are for the state field of the `struct vol_stat_st` defined in Section 13.3.4 on page 277. Following these definitions is an indication of which are stored in on-disk, permanent storage, in Section 13.3.10 on page 283.

```
/* These really belong to v_voltypes */
#define VOL_TYPEBITS    0x00003 /* low-level type of volume this is */
#define VOL_RW         0x00001 /* R/W volume */
#define VOL_READONLY   0x00002 /* ReadOnly Volume */

/* Some of these (that is, VOL_OK) are composites */
#define VOL_BUSY       0x00004 /* Volume is Busy */
/* The following three definitions are obsolete */
#define VOL_OFFLINE    0x00008 /* Volume is offline */
#define VOL_DELONSALVAGE 0x00010 /* Delete On Salvage */
#define VOL_OUTOFSERVICE 0x00020 /* Out Of service */
#define VOL_DEADMEAT    0x00040 /* About to be deleted */
#define VOL_LCLMOUNT    0x00080 /* Volume is mounted locally */

/* replication-specific state bits */
```

```

#define VOL_REPFIELD      0x00f00 /* which replication created this? */
#define VOL_REP_NONE     0x00000 /* none */
#define VOL_REP_RELEASE  0x00100 /* ``vos release'' */
#define VOL_REP_LAZY     0x00200 /* lazy replication */
#define VOL_IS_COMPLETE  0x01000 /* this volume instance is complete */
#define VOL_HAS_TOKEN    0x02000 /* our VVage should always be zero. */
#if 0
/* no longer used -- can be reclaimed for some other purpose */
#define VOL_KNOWDALLY    0x04000 /* know reclaimDally val-this vol */
#endif
#define VOL_NOEXPORT     0x08000 /* Do not export this volume! */

#define VOL_TYPEFIELD    0xf0000 /* what high-lvl type vol this is */
#define VOL_TYPE_RW     0x10000 /* read-write (ordinary) */
#define VOL_TYPE_RO     0x20000 /* ``.readonly'' */
#define VOL_TYPE_BK     0x30000 /* ``.backup'' */
#define VOL_TYPE_TEMP   0x40000 /* temporary use, dumps, moves, etc. */

#define VOL_GRABWAITING  0x100000 /* grab is pending for this volume */
#define VOL_LOOKUPWAITING 0x200000 /* lookup is pending for this volume */
#define VOL_REPSERVER_MGD 0x400000 /* managed by repserver */
#define VOL_MOVE_TARGET  0x800000 /* this vol is target for move op */
#define VOL_MOVE_SOURCE  0x1000000 /* this vol is source for move op */
#define VOL_ZAPME        0x2000000 /* delete vol with extreme prejudice */
#define VOL_CLONEINPROG  0x4000000 /* this vol is partially cloned */
#define VOL_IS_REPLICATED 0x8000000 /* this vol is replicated */
#define VOL_OPENDONE     0x10000000 /* back from VOL_OPEN procedure */
#define VOL_OK          ~(VOL_BUSY|VOL_OFFLINE|VOL_DELONSALVAGE| \
                          VOL_OUTOFSERVICE| VOL_ZAPME|VOL_CLONEINPROG| \
                          VOL_OPENDONE)

```

### 13.3.10 On-disk States

The following set of state definitions is stored in on-disk, permanent storage. All others are not stored on-disk for the fileset (and hence, not seen by the LFS layer).

```

VOL_RW
VOL_READONLY
VOL_DELONSALVAGE
VOL_IS_COMPLETE
VOL_NOEXPORT
VOL_REPSERVER_MGD
VOL_MOVE_TARGET
VOL_MOVE_SOURCE
VOL_ZAPME
VOL_CLONEINPROG
VOL_IS_REPLICATED

```

```

VOL_REPFIELD
VOL_REP_NONE
VOL_REP_RELEASE
VOL_REP_LAZY
VOL_TYPEFIELD
VOL_TYPE_RW
VOL_TYPE_RO
VOL_TYPE_BK
VOL_TYPE_TEMP

```

### 13.3.11 Kernel Maintained State Bits

These are the bits that the kernel maintains alone, not by VOLOP\_SETSTATUS.

```

#define VOL_BITS_NOSETSTATUS (VOL_BUSY|VOL_DEADMEAT|VOL_GRABWAITING| \
    VOL_LOOKUPWAITING|VOL_DELONSALVAGE| \
    VOL_OPENDONE)

/* Don't turn off VOL_DELONSALVAGE when attaching a volume. It is a */
/* kernel-maintained bit that is also persistent, unlike the other bits */
/* in VOL_BITS_NOSETSTATUS */

#define VOL_BITS_INITIALMASK (VOL_BITS_INITIAL | (VOL_BITS_NOSETSTATUS \
    & ~VOL_DELONSALVAGE))

/* During an identity swap, preserve these bits in their original struct
 * volume, because any waiters will be waiting on the address of that
 * structure's v_states field. */

#define VOL_BITS_NOSWAP (VOL_GRABWAITING|VOL_LOOKUPWAITING)

```

### 13.3.12 Useful Macros in Volume Pointers

```

#define VOL_READWRITE(volp) ((volp)->v_stat_st.states & VOL_RW)
#define VOL_EXPORTED(volp) (!((volp)->v_stat_st.states & VOL_NOEXPORT))

```

### 13.3.13 Mask Bits for VOL\_SETSTATUS MASK Argument

The setstatus operations take a mask that specifies the set of status fields to be updated. Bits within this mask are defined as follows.

```

#define VOL_STAT_VOLNAME          0x00000001
#define VOL_STAT_VOLID           0x00000002
#define VOL_STAT_VERSION         0x00000004
#define VOL_STAT_UNIQUE          0x00000008
#define VOL_STAT_OWNER           0x00000010
#define VOL_STAT_TYPE            0x00000020
#define VOL_STAT_STATES          0x00000040
#define VOL_STAT_ALLOCLIMIT      0x00000080
#define VOL_STAT_BACKUPID        0x00000100

```

```

#define VOL_STAT_PARENTID          0x00000200
#define VOL_STAT_CLONEID          0x00000400
#define VOL_STAT_CREATEDATE       0x00000800
#define VOL_STAT_UPDATEDATE       0x00001000
#define VOL_STAT_ACCESSDATE       0x00002000
#define VOL_STAT_COPYDATE         0x00004000
#define VOL_STAT_NODEMAX          0x00008000
#define VOL_STAT_VISLIMIT         0x00010000
#define VOL_STAT_MINQUOTA         0x00020000
/* spare VOL_STAT_SIZE           0x00040000 */
#define VOL_STAT_INDEX            0x00080000
#define VOL_STAT_BACKVOLINDEX     0x00100000
#define VOL_STAT_STATUSMSG        0x00200000
#define VOL_STAT_CLONETIME        0x00400000
#define VOL_STAT_VVCURRTIME       0x00800000
#define VOL_STAT_VVPINGCURRTIME  0x01000000
/* spare VOL_STAT_ACCERROR       0x02000000 */
#define VOL_STAT_BACKUPDATE       0x04000000
#define VOL_STAT_RECLAIMDALLY     0x08000000
#define VOL_STAT_LLBACKID         0x10000000
#define VOL_STAT_LLFWDID          0x20000000
#define VOL_STAT_VOLMOVETIMEOUT  0x40000000

```

### 13.3.14 Define the Fileset (volume) Structure

This structure is a memory resident structure maintained for each fileset. It is allocated by the file system independent layer (in the xvolume manager) at fileset attach time (*vol\_attach()*) and is passed as an argument to each of the fileset operations.

```

/*
 * In-core volume structure
 */

struct volume {
    struct squeue      v_lruq; /* lruq + free list queue */
    struct volume     *v_next;
    struct aggr       *v_paggrp; /* points to aggregate this */
                                /* fileset belongs to */
    struct volumeops  *v_volOps; /* points to volumeops vector */
    struct lock_data  v_lock; /* lock variable
    long              v_count; /* ref count for storage */
    opaque            v_fsDatap; /* points to file sys dependent */
                                /* PRIVATE data */
    struct vol_stat_st v_stat_st;
};

```

### 13.3.15 Define Volume Fields as v\_\*

By defining all volume's fields as v\_\* they can be freely moved around without changing the code all over.

```
/* stuff in vol_stat_st */
#define v_volId          v_stat_st.volId
#define v_parentId      v_stat_st.parentId
#define v_cloneTime     v_stat_st.cloneTime
#define v_vvCurrentTime v_stat_st.vvCurrentTime
#define v_vvPingCurrentTime v_stat_st.vvPingCurrentTime
#define v_voltype       v_stat_st.type
#define v_accStatus     v_stat_st.accStatus
#define v_accError      v_stat_st.accError
#define v_concurrency   v_stat_st.concurrency
#define v_states        v_stat_st.states
#define v_reclaimDally  v_stat_st.reclaimDally
#define v_volName       v_stat_st.volName
#define v_activeVnops   v_stat_st.activeVnops
#define v_procID        v_stat_st.procID
```

### 13.3.16 Define the vol\_Dirent Structure

This is the structure used by *vol\_readdir()* and *vol\_appenddir()*. See Appendix H on page 403 for information about the fields within this structure.

```
typedef struct vol_dirent {
    int32    offset;      /* Logical directory offset */
    int32    vnodeNum;   /* Canonical vnode number */
    u_int32  codesetTag; /* Tag for the name's codeset */
    u_short  reclen;     /* Len of this rec, 4-byte boundary */
    u_short  namelen;    /* Length of 'name' */
    char     name[OSI_MAXNAMLEN+1]; /* File name, variable sized */
} vol_dirent_t;
```

### 13.3.17 Define the Fileset Error Codes

The following define the fileset related error codes. Currently these overlay the *errno* error space; in fact, they are returned from the *afs\_syscall* interface via *errno*.

```
#define VOL_ERR_EOF          1
#define VOL_ERR_DELETED     2
#define VOL_ERR_EOW         8
```

### 13.3.18 Define the Fileset Open Operation Types

When an administrative, user-space utility *opens* a fileset, it supplies a bitmask that indicates the types of operations that will be performed on the open fileset. This is used in order to determine the class of concurrent operations that will be permitted on the fileset while it is open. Bits within this mask are defined as follows.

```
/*
 * VOLOP bit definitions
 */

#define VOL_OP_SEEK          0x00000001
```



```

#define VOL_OP_TELL          0x00000002
#define VOL_OP_SCAN         0x00000004
#define VOL_OP_DESTROY      0x00000008
#define VOL_OP_GETSTATUS    0x00000010
#define VOL_OP_SETSTATUS    0x00000020
#define VOL_OP_CREATE       0x00000040
#define VOL_OP_READ         0x00000080
#define VOL_OP_WRITE        0x00000100
#define VOL_OP_TRUNCATE     0x00000200
#define VOL_OP_DELETE       0x00000400
#define VOL_OP_GETATTR      0x00000800
#define VOL_OP_SETATTR      0x00001000
#define VOL_OP_GETACL       0x00002000
#define VOL_OP_SETACL       0x00004000
#define VOL_OP_CLONE        0x00008000
#define VOL_OP_RECLONE      0x00010000
#define VOL_OP_UNCLONE      0x00020000
#define VOL_OP_SETNEWVID    0x00040000
#define VOL_OP_COPYACL      0x00080000
#define VOL_OP_SWAPIDS      0x00100000
/* now, for the piece of SETSTATUS that changes the fileset ID */
#define VOL_OP_SETSTATUS_ID 0x00200000
#define VOL_OP_SYNC         0x00400000
#define VOL_OP_PUSHSTATUS   0x00800000
#define VOL_OP_READDIR     0x01000000
#define VOL_OP_APPENDDIR    0x02000000
#define VOL_OP_GETZLC       0x04000000
/* Forcibly (architecturally) deny concurrent vnode ops */
#define VOL_OP_NOACCESS     0x08000000
#define VOL_OP_GETNEXTHOLES 0x10000000
#define VOL_OP_DETACH       0x20000000
#define VOL_OP_DEPLETE      0x40000000

```

### 13.3.19 Define the Fileset Operations VOL\_SYS\_XXX

```

/*
 * VOL syscall to VOLOP translation
 * Translates syscalls to volop collections;
 * Depends on xvolume/vol_init.c implementation.
 */

#define VOL_SYS_SEEK          (VOL_OP_SEEK)
#define VOL_SYS_TELL          (VOL_OP_TELL)
#define VOL_SYS_SCAN         (VOL_OP_SCAN)
#define VOL_SYS_DESTROY      (VOL_OP_DESTROY)
#define VOL_SYS_GETSTATUS    (VOL_OP_GETSTATUS)
#define VOL_SYS_SETSTATUS    (VOL_OP_SETSTATUS)
#define VOL_SYS_CREATE       (VOL_OP_CREATE)
#define VOL_SYS_READ         (VOL_OP_READ)
#define VOL_SYS_WRITE        (VOL_OP_WRITE)
#define VOL_SYS_TRUNCATE     (VOL_OP_TRUNCATE)
#define VOL_SYS_DELETE       (VOL_OP_DELETE)
#define VOL_SYS_GETATTR      (VOL_OP_GETATTR)

```

```

#define VOL_SYS_SETATTR          (VOL_OP_SETATTR)
#define VOL_SYS_GETACL          (VOL_OP_GETACL)
#define VOL_SYS_SETACL          (VOL_OP_SETACL)
#define VOL_SYS_CLONE           (VOL_OP_CLONE)
#define VOL_SYS_RECLONE         (VOL_OP_RECLONE)
#define VOL_SYS_UNCLONE         (VOL_OP_UNCLONE)
#define VOL_SYS_SETVV           (VOL_OP_GETSTATUS | VOL_OP_SETSTATUS)
#define VOL_SYS_SWAPVOLIDS      (VOL_OP_SWAPIDS)
#define VOL_SYS_COPYACL         (VOL_OP_COPYACL)
#define VOL_SYS_SYNC            (VOL_OP_SYNC)
#define VOL_SYS_PUSHSTATUS      (VOL_OP_PUSHSTATUS)
#define VOL_SYS_READDIR         (VOL_OP_READDIR)
#define VOL_SYS_APPENDDIR       (VOL_OP_APPENDDIR)
#define VOL_SYS_GETZLC          (VOL_OP_GETZLC)
#define VOL_SYS_GETNEXTHOLES    (VOL_OP_GETNEXTHOLES)
/* Forcibly (architecturally) deny concurrent vnode ops */
#define VOL_SYS_NOACCESS        (VOL_OP_NOACCESS)
#define VOL_SYS_DETACH          (VOL_OP_DETACH)
#define VOL_SYS_DEPLETE         (VOL_OP_DEPLETE)

```

### 13.3.20 Sync Type Values for vol\_sync()

```

#define VOL_SYNC_COMMITSTATUS    1
#define VOL_SYNC_COMMITMETA     2
#define VOL_SYNC_COMMITALL      3

```

### 13.3.21 Define the Fileset Handle Structure

The following is the internal volume handle for all VOP\_\* calls. It is an in-kernel structure only; a volume descriptor (long) is passed to the user space to associate it with this structure.

There, as an administrative utility performs actions on an open fileset, certain information regarding its progress (namely which file on the fileset is being processed) is recorded in a **vol\_handle** structure. This is passed as an argument to several of the fileset operations.

The *.fid*, *.index* and *.type* fields below are not interpreted outside of the LFS layer.

```

struct vol_handle {
    struct volume *volp; /* pointer to incore volume */
    long opentype; /* Open volume type */
    /* Mask of VOL_OP_XXX bits specified at vol_open(0 time and */
    /* gives the type ops to perform */
    struct afsFid fid; /* Fid for "active" volume anode */
    long index; /* active Anode's index */
    /* set by vol_open(), vol_seek(), */
    /* vol_scan() and vol_create(). */
    long type; /* active Anode's mode/type */
};

```

### 13.3.22 Define the Root Anode Index

The following represents the special index for the volume's root anode. The special dummy index is required since various physical file systems use different numbers for their root anode (ufs = 2, episode = 15, etcetera). Do *not* change its definition since code depends on it (for instance, scanning -1, 0, 1 ... to maxino is assumed by various programs). Ordinary files are indicated by indices of 0, 1, 2, and so forth.

```
#define VOL_ROOTINO    -1
```

### 13.3.23 Define the Maximum Quota Size

Define the largest 64-bit constant that is evenly divisible by 64K and whose value will fit in 32 bits when divided by 1024. This is the maximum value allowed by LFS when the aggregate fragment size is 1K (the smallest frag. size allowed); although larger quota values are possible when the fragment size is increased. The value must also be evenly divisible by 64K to prevent rounding errors when converting to units of 64K (or less), the largest fragment size allowed by LFS.

The following 64 bit number is approximately 4.4 terabytes.

```
#define VOL_MAX_QUOTA_HIGH    0x3ff
#define VOL_MAX_QUOTA_LOW    0xffff0000
```

## 13.4 Extended Credential Data Types

This section will define the data types exported by the credential package described in *The xcred Package*, Section 17.1 on page 469.

### 13.4.1 Define the xcred\_PListEntry\_t Structure

This structure is used to build a list of all attribute-value pairs for an **xcred**, including derivation relationships. It is used by *xcred\_DeleteEntry()*. Section 17.3 on page 484.

```
/* Property list entry. */

typedef struct xcred_PListEntry {
    struct xcred_PListEntry *nextP; /* Next entry in linear plist */
    struct xcred_PListEntry *prevP; /* Prev entry in linear plist */
    long attrBytes; /* Size of attr (incl. null) */
    char *attributeP; /* Ptr to attr. component */
    long valueBytes; /* Size of value component */
    char *valueP; /* Ptr to value component */
    long flags; /* Internal flags */
} xcred_PListEntry_t;
```

### 13.4.2 Define the `xcred_t` Structure

This refcounted and lockable structure contains all the extended credential info associated with a given user, stored in a property list. Also stored is information on the UNIX cred with which this structure is associated. Included is a generation number, incremented each time someone changes the property list.

```
/* Xcred structure. */

typedef struct xcred {
    struct squeue lruq; /* LRU queue */
    struct lock_data lock; /* Structure lock */
    long refcount; /* Reference count */
    long changeCount; /* # times plist changed */
    long flags; /* Flags described above */
    long uflags; /* More user-defined flags */
    xcred_PListEntry_t *propListP; /* Property list */
    struct ucred *assocUCredP; /* Assoc UNIX cred struct */
    long assocPag; /* Unique identifier in ucred */
} xcred_t;
```

### 13.4.3 Reserved Attributes

The following parts of the attribute space are reserved for the following purposes. Other non-conflicting names may be defined by users.

- *NFSId* contains a long integer representing the user's identity according to an unauthenticated NFS call.

## 13.5 The VFS+ Switch

### 13.5.1 DCE LFS VNOPS Vector Organization

#### Example 13-1 Initializing the Extended VNOPS Vector

```
/*
 * The extended VNOPS vector is initialized as follows:
 *
 * OOPS point to xglue_*
 * XOPS point to the passed in 'axfuncs'
 * NOPS point to the passed in 'ofuncs'
 *
 * The CM calls with 'axfuncs' pointing to the CM VNOPS;
 * EPISODE calls with 'axfuncs' pointing to the EFS VNOPS;
 * Everyone calls with 'ofuncs' pointing to n<os>_ops.
 *
 * This results in an extended VNOPS vector such that
 * VOP_xxx calls on vnodes in the AFS VFS generate a
 * call chain like this:
 *
 * VOP_xxx
 * xglue_xxx() * does tkc synchronization *
```

```

*          n<os>_xxx()    * maps VFS to VFS+          *
*          cm_xxx()
*
* VOP_xxx calls on vnodes in the LFS VFS generate a
* call chain like this:
*
*          VOP_xxx
*          xglue_xxx()    * does tkc synchronization *
*          n<os>_xxx()    * maps VFS to VFS+          *
*          efs_xxx()
*/

```

## 13.5.2 LFS Generating Extended Vnode Operations

### Example 13-2 Generating New Style Vnode Operations From Original

```

/*
* Used by GetNewVnodeOpsFromOld (below) to generate new style vnode
* operations from original-style vnode operations.
* In other words, this is used by ufs (not episode or cm) to
* get extended vnode structures for its vnodes.
*
* We initialize the extended VNOPS vector as follows:
*
*          OOPS  point to xglue_*
*          XOPS  point to xufs_*
*          NOPS  point to the passed in 'aofuns'
*
* The PX calls xvfs_convert(vp) to convert a vnode to a vnode
* with the extended vnode ops. This results in a call
* to this routine, where 'aofuns' points to ufs vnops or
* efs vnops, depending on which VFS was referenced by PX.
*
* This results in an extended VNOPS vector such that
* VOP_xxx calls on vnodes in the AFS VFS generate a
* call chain like this:
*
*          VOP_xxx
*          xglue_xxx()    * does tkc synchronization *
*          ufs_xxx()
*
* or
*
*          VOP_xxx
*          xglue_xxx()    * does tkc synchronization *
*          nosf_xxx()    * maps VFS to VFS+          *
*          efs_xxx()
*
* Calls via VOPX (for example, by PX) generate a call chain like this:
*
*          VOPX_xxx

```

```

*      efs_xxx()
*
* or
*
*      VOPX_xxx
*      xufs_xxx()      * maps VFS+ to VFS      *
*      ufs_xxx
*/

```

### 13.5.3 Converting Old Vnode Operations

#### Example 13-3 GetNewVnodeOpsFromOld Routine

```

/*
 * Enter some old vnode ops into a conversion hash table, which is
 * just a linked list for now.
 */
static struct xvfs_vnodeops *GetNewVnodeOpsFromOld(
    struct osi_vnodeops *aoldOpsp)
{
    register struct xvfs_opspair *tp;
    register struct xvfs_vnodeops *np;

    /* Look for the vnode ops and see if we've already done a conversion
     * from this ops.  If so, return it.
     *
     * Also, we look to see if these vnode ops are already converted.
     * If they are, we just return the same ops pointer.  This
     * path is taken when UFS recycles an inode: it clears v_flag (which
     * clears the converted flag), but doesn't reset the vnode ops pointer.
     * We *really* don't want to convert the vnode a second time.
     */
    for (tp = xvfs_AllOpsPairs; tp; tp = tp->next) {
        if (aoldOpsp == (struct osi_vnodeops *) tp->newOps)
            return (struct xvfs_vnodeops *) aoldOpsp;

        if (aoldOpsp == tp->oldOps)
            return tp->newOps;
    }
}

```

## 13.6 Vnode Data Types

### 13.6.1 Preliminary Items

There is a preliminary check that must be made for Vnode Operations, as follows.

```
#ifndef TRANSARC_XVNODE_H
#define TRANSARC_XVNODE_H
#include <dcedfs/volume.h>
#include <dcedfs/common_data.h>
#include <dcedfs/osi.h>
#include <dcedfs/osi_cred.h>

#ifdef AFS_DEFAULT_ENV
#error "Check that xvfs_genvnode.h is right for you"
#include <dcedfs/xvfs_genvnode.h> /* std VFS systems */
#endif /* AFS_DEFAULT_ENV */
```

### 13.6.2 Converted Vnode Indication

The status flag described in Section C.1 on page 245 is defined here, along with it's relatives, as follows.

```
/* Exported structures and functions are preceded by module
 * name xvfs_ */

/*
 * Define constant used in v_flag field to indicate that the vnode
 * has already been converted.
 */

#ifndef AFS_OSF_ENV
#define V_CONVERTED 0x1000 /* looks out of the way for now */
#endif

#ifdef AFS_AIX31_ENV
#define IS_CONVERTED(vp) ((vp)->v_gnode->gn_flags & V_CONVERTED)
#define SET_CONVERTED(vp) (vp)->v_gnode->gn_flags |= V_CONVERTED
#define SET_UNCONVERTED(vp) (vp)->v_gnode->gn_flags &= ~V_CONVERTED
#else
#define IS_CONVERTED(vp) ((vp)->v_flag & V_CONVERTED)
#define SET_CONVERTED(vp) (vp)->v_flag |= V_CONVERTED
#define SET_UNCONVERTED(vp) (vp)->v_flag &= ~V_CONVERTED
#endif
#define XVFS_CONVERT(vp) (IS_CONVERTED(vp) ? 0 : xvfs_convert(vp))
#define xvfs_PutVolume(volp) if (volp) VOL_RELE(volp)
```

### 13.6.3 Define the `xvfs_attr` Structure

The following defines additional attribute fields that have special meaning in DFS; used by both Episode and UFS. The Fields flagged with an `/* on disk */` must be stored in permanent, on-disk storage for the file.

```

struct Txvattr {
    afsHyper dataVersion; /* on disk */
    afsHyper fileID; /* on disk */
    afsHyper volVersion; /* on disk */
    u_long author;
    u_long callerAccess;
    u_long anonAccess;
    u_long parentVnode;
    u_long parentUnique;
    afsTimeval serverModTime;
    u_long fstype; /* type of the containing aggregate, from aggr.h */
    /* Values that are principally derived from what Episode stores
       for DFS */
    afsUUID objid; /* on disk */
    u_long timeUncertainty; /* on disk */
    /* for Episode, the CFLAGS from epia_GetInfo (anode.h),
       including EPIA_CFLAGS_COPYONWRITE */
    u_long representationFlags;
    /* Magic cookies usable to identify common ancestry. */
    u_long backingIndex; /* on disk */
    u_long backingVolIndex; /* on disk */
    /*
     * These Ix values are cookies that are interpretable only by the
     * underlying representation.
     */
    u_long aclIx; /* on disk */
    u_long initDirAclIx; /* on disk */
    u_long initFileAclIx; /* on disk */
    u_long plistIx; /* on disk */
    u_long uPlistIx; /* on disk */

    u_long clientOnlyAttrs; /* rsvrd for machines that never run FXs */
                          /* on disk */ /* clientOnlyAttrs */

    u_long spare1;
    u_long spare2;
    u_long spare3;
    u_long spare4;
    u_long spare5;
    u_long spare6;
};

```

These extended attributes have the following meanings.

**.dataVersion** This 64-bit version number increases monotonically whenever the file's contents change ( whenever the **mtime** is updated). Its advantage over the **mtime** is that it will never go backwards.

**Note:** The value in this field is used by the DFS cache client (CM).



<i>.fileID</i>	The inode index and inode generation numbers for the file, in that order. These are the same values that would reside in an <b>afsFid</b> naming the file.
<i>.volVersion</i>	Whenever the volume version number is incremented due to a change to this file, that (incremented) volume version is stored with the file.
<i>.author</i>	This unused field is set to <code>-1</code> by the <code>vn_getattr()</code> and <code>vol_getattr()</code> calls.
<i>.callerAccess</i>	This <b>permset_t</b> (see the DFS <b>ACL</b> specification, <i>ACL Storage Format</i> , Chapter 9 on page 165 ) is set to the rights that the caller of the <code>vn_getattr()</code> or <code>vol_getattr()</code> call has to the file.
<i>.anonAccess</i>	This <b>permset_t</b> is equal to the set of rights that an object's <b>ACL</b> grants to everyone, even unauthenticated clients. If an object does not have an <b>ACL</b> , this field will be 0. See the DFS <b>ACL</b> specification mentioned above, and the description of <code>vn_getattr()</code> in Section 16.7 on page 432. This field is used by the client cache manager in some situations to avoid having to perform a remote <b>ACL</b> check.
<i>.parentVnode</i>	This unused field is set to <code>-1</code> by the <code>vn_getattr()</code> and <code>vol_getattr()</code> calls.
<i>.parentUnique</i>	This unused field is set to <code>-1</code> by the <code>vn_getattr()</code> and <code>vol_getattr()</code> calls.
<i>.serverModTime</i>	This field should be set to the file's modified time ( <code>mtime</code> ) by the <code>vn_getattr()</code> and <code>vol_getattr()</code> calls.
<i>.fstype</i>	This field is set to the appropriate <code>AP_TYPE_XXX</code> constant for the LFS by the <code>vn_getattr()</code> and <code>vol_getattr()</code> calls.
<i>.objid</i>	The object UUID for this file.
<i>.timeUncertainty</i>	This field is currently unused; the <code>vol_getattr()</code> and <code>vn_getattr()</code> calls should return 0 for it and the <code>vol_setattr()</code> and <code>vn_setattr()</code> calls should ignore it. It was intended to hold DTS time uncertainties to accompany the file's <code>atime</code> , <code>mtime</code> and <code>ctime</code> .
<i>.representationFlags</i>	This field is set to a set of internal, LFS-dependent state flags.  <b>Note:</b> These are a set of internal Episode flags. It is not known yet whether they have any DFS (non-diagnostic) use and hence require architecting.

The following two fields will only be non-zero for a cloned file. They identify the *backing* file for files in a copy-on-write relationship. For example: if F2 is a read-only clone of the read-write fileset F1 (F1 ==> F2), ( meaning F1 is backed by F2), these fields would be zero for F2 and non-zero for F1.

**Note:** It is not known whether these fields serve any actual purpose or are there for diagnostic (or obsolete) reasons. If the former, it's not clear to what degree a particular LFS has to agree with what Episode puts here.

<i>.backingIndex</i>	The index of the fileset which is holding the backing, copy-on-write file for this one. This index has the same interpretation (and significance) as the one used by the <code>ag_volinfo()</code> aggregate call. Note that filesets related by cloning must reside on the same aggregate.
<i>.backingVolIndex</i>	The index of the file which is serving as the backing, copy-on-write file for this one. This field has the same interpretation (and significance) as the <i>.Vnode</i> field in a file's <b>afsFid</b> (or <i>.fileID</i> , see above).

In order for fileset moves and backups (both implemented via cloning) to be transparent, a file and its backing file will ALWAYS have the same index.

*.clientOnlyAttrs* A set of uninterpreted (by the LFS) flags maintained with the file on behalf of DFS.

The following five fields should be set to a non-zero value by the *vn\_getattr()* and *vol\_getattr()* calls if the file has an associated auxiliary object of the particular type. Although the precise interpretation of the fields is not architected, it is suggested that they somehow identify or name their auxiliary object such that sharing within a fileset (for example, shared ACLs) can be detected via a simple equality check. Obviously, the initial directory and file ACL fields are only required for directory objects.

*.aclIx* Object ACL.  
*.initDirAclIx* For a directory: the *Initial Container* (directory) ACL.  
*.initFileAclIx* For a directory: the *Initial File* ACL.  
*.plistIx* The system property list.  
*.uPlistIx* The user property list.

The "Enhanced" attribute structure: The standard "vattr" structure is simply an overlay (always first) in this structure. This structure consists of both the standard vnode and extended DFS attributes.

```
struct xvfs_attr {
    struct vattr vattr;
    struct Txvattr xvattr;
};
```

#### 13.6.4 Define the Vnode Operations Function Array

A description of the operations available within the vnode function array appears below. The array is an exported data type.

```
/* extended VOPX operations; note that all attr-taking dudes
 * take xvattr, not vattr.
 */
struct xvfs_xops {
    int (*vn_open)();
    int (*vn_close)();
    int (*vn_rdwr)();
    int (*vn_ioctl)();
    int (*vn_select)();
    int (*vn_getattr)();
    int (*vn_setattr)();
    int (*vn_access)();
    int (*vn_lookup)();
    int (*vn_create)();
    int (*vn_remove)();
    int (*vn_link)();
    int (*vn_rename)();
    int (*vn_mkdir)();
    int (*vn_rmdir)();
    int (*vn_readdir)();
    int (*vn_symlink)();
```

```

int (*vn_readlink)();
int (*vn_fsync)();
int (*vn_inactive)();
int (*vn_bmap)();
int (*vn_strategy)();
int (*vn_ustrategy)(); /* assuming stuff already mapped in */
int (*vn_bread)();
int (*vn_brelse)();
int (*vn_lockctl)();
int (*vn_fid)(); /* op for old style fid op */
int (*vn_hold)(); /* maybe we don't need these; revisit when done */
int (*vn_rele)();
/*
 * new ones for us to provide, rather than just existing to
 * make writing the O functions easier (that is, porting).
 */
int (*vn_setacl)();
int (*vn_getacl)();
int (*vn_afsfid)();
int (*vn_getvolume)();
int (*vn_getlength)();
/*
 * Some new ops for AIX 3
 */
int (*vn_map)(); /* also used for SunOS 5 */
int (*vn_unmap)();
/*
 * A new op for OSF/1
 */
int (*vn_reclaim)();
/*
 * Some new ops for SunOS 5
 */
int (*vn_read)();
int (*vn_write)();
int (*vn_realvp)();
void (*vn_rwlock)();
void (*vn_rwunlock)();
int (*vn_seek)();
int (*vn_space)();
int (*vn_getpage)();
int (*vn_putpage)();
int (*vn_addmap)();
int (*vn_delmap)();
int (*vn_pageio)();
#define vn_frlock vn_lockctl /* overlay equivalent ops */
/*
 * Ops for HP/UX
 */
int (*vn_pagein)();
int (*vn_pageout)();
};

```

### 13.6.5 Define the Enhanced Vnode Operations

A description of the operations available within the enhanced (converted) vnode function array appears below. The array is an exported data type.

```
/*
 * Operations on the Enhanced vnodes.
 */
struct xvfs_vnodeops {
    struct osi_vnodeops oops; /* Glued vendor vnodeops */
    struct xvfs_xops xops; /* Enhanced (Decorum) vnodeops */
    struct osi_vnodeops nops; /* Original vendor vnodeops */
};
```

The **osi\_vnodeops** structure is somewhat of a misnomer. It is the native Operating System (OS) **vnode** vector, which is anything but osi (Operating System Independent).

### 13.6.6 Define VOPX\_XXX for Extended Vnodeops

```
/*
 * Macros for extended vnodeops
 */
#define VOPX_OPEN          efs_null
#define VOPX_CLOSE        efs_null
#define VOPX_RDWR         efs_rdwr
#define VOPX_IOCTL        efs_invalid
#define VOPX_SELECT        efs_invalid
#define VOPX_GETATTR      efs_getxattr
#define VOPX_SETATTR      efs_setxattr
#define VOPX_ACCESS        efs_access
#define VOPX_LOOKUP        efs_lookup
#define VOPX_CREATE        efs_create
#define VOPX_REMOVE        efs_remove
#define VOPX_LINK          efs_link
#define VOPX_RENAME        efs_rename
#define VOPX_MKDIR         efs_mkdir
#define VOPX_RMDIR         efs_rmdir
#define VOPX_READDIR       efs_readdir
#define VOPX_SYMLINK       efs_symlink
#define VOPX_READLINK      efs_readlink
#define VOPX_FSYNC         efs_fsync
#define VOPX_INACTIVE      efs_inactive
#define VOPX_BMAP          efs_bmap
#if defined(AFS_AIX31_VM)
#define VOPX_STRATEGY(VP,BP) efs_strategy (BP)
#else
#define VOPX_STRATEGY      efs_null
#endif
#define VOPX_ISTRATEGY     efs_panic
#define VOPX_BREAD         efs_panic
#define VOPX_BRELSE        efs_panic
#define VOPX_LOCKCTL       efs_lockctl
#define VOPX_FID           efs_fid
#define VOPX_HOLD          efs_hold
```

```

#define VOPX_RELE      efs_rele
#define VOPX_GETACL   efsx_getacl
#define VOPX_SETACL   efsx_setacl
#define VOPX_AFSFID   efs_afsfid
#define VOPX_GETVOLUME efs_getvolume
#define VOPX_GETLENGTH efs_getlength
#ifdef AFS_AIX31_ENV
#define VOPX_MAP      efs_map
#define VOPX_UNMAP    efs_unmap
#endif
#ifdef AFS_OSF_ENV
#define VOPX_RECLAIM  efs_reclaim
#endif
#ifdef AFS_SUNOS5_ENV
#define VOPX_READ     efs_vmread
#define VOPX_WRITE    efs_vmwrite
#define VOPX_REALVP   efs_realvp
#define VOPX_RWLOCK   efs_rwlock
#define VOPX_RWUNLOCK efs_rwunlock
#define VOPX_SEEK     efs_seek
#define VOPX_SPACE    efs_space
#define VOPX_GETPAGE  efs_getpage
#define VOPX_PUTPAGE  efs_putpage
#define VOPX_ADDMAP   efs_addmap
#define VOPX_DELMAP   efs_delmap
#define VOPX_PAGEIO   efs_pageio
#define VOPX_FRLOCK   efs_frlock
#endif
#ifdef AFS_HPUX_ENV
#define VOPX_PAGEIN   efs_pagein
#define VOPX_PAGEOUT  vfs_pageout
#endif

```

### 13.6.7 Define VOPX\_UPDATE Flags

```

/*
 * Flags argument for VOPX_UPDATE
 */
#define XVN_ACC      1    /* atime */
#define XVN_UPD      2    /* mtime */
#define XVN_CHG      4    /* ctime */

```

### 13.6.8 Define the Enhanced Operations Vector

This is the enhanced (converted) VFS operations vector.

```

/*
 * The extended VFS op vector
 */
struct xvfs_vfsops {
    struct vfsops xvfsops;
    struct vfsops vfsops;
    int (*vfsgetvolume)();
};

```

### 13.6.9 Define the VFS Operations

The following defines a VFS operations vector.

```
struct vfsops (
    int (*vfs_mount )();
    int (*vfs_unmount )();
    int (*vfs_root )();
    int (*vfs_statfs )();
    int (*vfs_sync )();
    int (*vfs_vget )();
    int (*vfs_getmount )();
};
```

### 13.6.10 Define the Vnode Operation Classifications

The following definitions classify the **vnode** operations into noop, read-only and read-write.

```
/* Macros to classify vnode ops into noop, read-only and read-write */
#define VNOP_LOCK      VNOP_TYPE_READWRITE
#define VNOP_LINK      VNOP_TYPE_READWRITE
#define VNOP_UNLINK    VNOP_TYPE_READWRITE
#define VNOP_MKDIR     VNOP_TYPE_READWRITE
#define VNOP_RMDIR    VNOP_TYPE_READWRITE
#define VNOP_RENAME    VNOP_TYPE_READWRITE
#define VNOP_SYNCGP    VNOP_TYPE_READWRITE
#define VNOP_TRUNC     VNOP_TYPE_READWRITE
#define VNOP_GETVAL    VNOP_TYPE_READONLY
#define VNOP_RWGP      VNOP_TYPE_READWRITE
#define VNOP_STAT      VNOP_TYPE_READWRITE
#define VNOP_UPDATE    VNOP_TYPE_READWRITE
#define VNOP_OPEN      VNOP_TYPE_READWRITE
#define VNOP_CLOSE     VNOP_TYPE_READWRITE
#define VNOP_READLINK VNOP_TYPE_READWRITE
#define VNOP_SYMLINK   VNOP_TYPE_READWRITE
#define VNOP_BMAP      VNOP_TYPE_READWRITE
#define VNOP_NAMEI     VNOP_TYPE_READWRITE
#define VNOP_MKNOD     VNOP_TYPE_READWRITE
#define VNOP_REMOVE    VNOP_TYPE_READWRITE
#define VNOP_LOOKUP    VNOP_TYPE_READONLY

/*
 * Change MAP and UNMAP to NOOP type vnodeops instead of READWRITE type
 * to prevent MAP and UNMAP vnodeops from blocking on busy filesets
 * as the AIX kernel global shared memory lock is held across these
 * calls. If these vnode ops block, any process trying to start or exit
 * will deadlock on the held shared memory lock
 */
#define VNOP_MAP      VNOP_TYPE_NOOP
#define VNOP_UNMAP    VNOP_TYPE_NOOP

#define VNOP_ACCESS   VNOP_TYPE_READONLY
#define VNOP_GETATTR  VNOP_TYPE_READONLY
#define VNOP_SETATTR  VNOP_TYPE_READWRITE
```

```

#define VNOP_FSYNC      VNOP_TYPE_READWRITE
#define VNOP_FTRUNC    VNOP_TYPE_READWRITE
#define VNOP_RDWR      VNOP_TYPE_READWRITE
#define VNOP_LOCKCTL   VNOP_TYPE_READWRITE
#define VNOP_READDIR   VNOP_TYPE_READWRITE
#define VNOP_GETACL    VNOP_TYPE_READONLY
#define VNOP_SETACL    VNOP_TYPE_READWRITE
#define VNOP_PGRD      VNOP_TYPE_READWRITE
#define VNOP_PGWR      VNOP_TYPE_READWRITE
#define VNOP_CREATE    VNOP_TYPE_READWRITE

/* SunOS relies on VFS_ROOT never blocking. If VFS_ROOT blocks, it holds
 * the vnode mutex which prevents anybody trying to obtain a reference
 * to the vnode to block that leads to a deadlock.
 * HP, AIX and Solaris hold a reference to a filesystems's root vnode always
 * and hence there should not be any problem if VFS_ROOT is unglued.
 */
#ifndef AFS_OSF_ENV
#define VFSOP_ROOT      VNOP_TYPE_NOOP
#else
#define VFSOP_ROOT      VNOP_TYPE_READWRITE
#endif

#define VFSOP_UNMOUNT  VNOP_TYPE_READWRITE
#define VFSOP_FHTOVP   VNOP_TYPE_READWRITE
#define VFSOP_VGET     VNOP_TYPE_READWRITE

```

### 13.6.11 Directory Entry Formats

The following defines the native directory entries of various kinds as shown by the `ifdef` statements.

```

#ifdef AIX                /* Native dir entry has an offset */
struct dirent {
    u_long      offset;
    u_long      inode;
    u_short     recordlen;
    namelen;
    char        dir_name[];
};
#elif SunOS-5            /* Native dir entry has an offset */
struct dirent {
    u_long      inode;
    u_long      offset;
    u_short     recordlen;
    u_short     namelen;
    char        dir_name[];
};
#else                    /* Native dir entries do NOT have an offset */
struct dirent {
    u_long      inode;
    u_short     recordlen;
    u_short     namelen;
};

```

```
    char          dir_name[];  
    u_long        offset;  
};
```

The fields have the following meaning:

- .offset*            An offset corresponding to whatever (empty space or another entry) follows this entry in the directory.
- .inode*            An LFS-dependent (interpreted only by the LFS) inode number.
- .recordlen*        Length of this entry, rounded up to a 4 byte boundary. Therefore, the offset from the start of this entry to the next one in a supplied buffer.
- .namelen*          Length of the file name, not including the terminating null char.
- .dir\_name*         Variable sized file name: includes a terminating null char.



# Aggregate Operations Interface

This chapter describes the portion of the VFS+ interface which allows operations on aggregates, the **DCE DFS** data containers which house filesets, as introduced in Section 11.3 on page 235. A description of each aggregate known to the local host is kept in the *Aggregate Registry*, a global kernel table. An aggregate descriptor includes a pointer to a **struct aggrops** function array, namely those operations that may be performed on the aggregate. The corresponding aggregate operation is chosen from the function array and executed. The aggregate operation architecture is remarkably similar to that used to provide access to VFS operations. In fact, one implementation technique for these functions is to define an aggregate file system via a **struct vfs**, suitably extended, and use the standard VFS mechanisms to manipulate them.

A description of the operations available within the aggregate function array appears below. The array is presented as an exported data type, and the associated declaration may be found in Section 13.2.12 on page 272. It is reproduced here for convenience.

```
struct aggrops {
    int (*ag_hold)();
    int (*ag_rele)();
    int (*ag_lock)();
    int (*ag_unlock)();
    int (*ag_stat)();
    int (*ag_volCreate)();
    int (*ag_volInfo)();
    int (*ag_detach)();
    int (*ag_attach)();
    int (*ag_sync)();
};
```

## 14.1 Initialization

Each type or implementation of a DFS LFS file system is identified by a small integer constant within the file **aggr.h** within the DFS source directory **src/file/xaggr**. Section 13.2.4 on page 270, *Valid Aggregate Types* lists the currently defined aggregate types.

### 14.1.1 Identifying a New LFS to DCE DFS

When a new file system is implemented, a new **AG\_TYPE\_xxx** constant must be added to this file in order to identify its type. In this constant, the **xxx** is the actual identifier of the type of aggregate (and hence, LFS) that is being defined. Additionally, as previously mentioned in Section 11.8.1 on page 240, the constant **MAX\_AG\_TYPE** (in Section 13.2.4 on page 270) should be adjusted since this controls the size of the aggregate operations array discussed below.

These aggregate types need not be consistent across multiple vendor platforms. They do not cross the network in any protocols - the aggregate operations are local affairs. It would be beneficial if they were; however, there is no provision currently by which this can be accomplished.

**Note:** Re-using the aggregate type chosen by Transarc for Episode is not an option.

1. It would preclude running with both Episode and the vendor's LFS.

2. It would probably be illegal.
3. The **backup** and **restore** algorithms assume that when moving between aggregates of the same type, directory offsets can be preserved. Undoubtedly, there are other technical differences that would cause difficulties.

### 14.1.2 Registering Aggregate Operations

As previously mentioned in the first paragraph of Chapter 14 on page 303, in the discussion of **Aggregate-level facilities** in Section 11.3 on page 235 and Section 11.8.1 on page 240, there is an *Aggregate Registry* table containing aggregate information that has been made available to **DCE DFS**, by the act of an LFS registering its aggregate operations at initialization time. Registration is made by use of the *ag\_setops()* function call. The result is that a pointer to the aggregate operations vector just identified is inserted in the *Aggregate Registry* in the appropriate offset identified by the aggregate type identifier. This results in the user space DFS components being informed about a new aggregate type.

The aggregate operations for a particular file system are identified by a **struct agprops** entrypoint vector identified in Section 13.2.12 on page 272.

## 14.2 Exporting volumes to DFS

The **dfsexport** command is used to export the filesets contained on an aggregate to DFS. Any aggregates being exported are required to have an entry in the *dfstab* file (*/opt/dcelocal/var/dfs/dfstab*). Each entry in this file identifies an aggregate and, for Transarc LFS file systems (as opposed to exported UFS file systems), contains the following fields:

**Device Name** Pathname to the device file for the aggregate.

**Aggregate Name** A 64-character (byte) aggregate name interpreted by neither DFS nor the LFS.

**File System Type** The **aggregate type**, or more succinctly, the **file system type** that is housing the aggregate. For an exported UFS, this is "**ufs**". For Transarc's Episode, this is "**lfs**".

For a new implementation of an LFS, a new type string should be declared in the DFS source file *src/file/xaggr/astab.h*. For example, the following type for an appropriate "xxx" value would be added to those defined in Section 13.2.18 on page 274:

```
#define ASTABTYPE_XXX "xxx"
```

**Aggregate ID** Each aggregate on a machine is identified by an **aggregate ID**. These IDs, which are selected by the system administrator as the *dfstab* file is constructed, are uninterpreted by DFS/LFS and are simply unique (for that local machine) integers starting from 1.

**Note:** The format of entries in this table can be different for different aggregate types. For example, an exported UFS has a **dfstab** entry that is different from what is shown above. If necessary (or beneficial), a new *lfs* could define private fields holding information to be passed to it at export time. Since DFS already has a mechanism for passing LFS-private data to *ag\_attach()*, only the **dfsexport** command itself would require changes in this event. (Specifically: **dfstab.c** and **export.c** in *src/file/xaggr*.)

Appendix G on page 401 lists the steps involved in exporting the filesets in an aggregate.

### 14.3 Aggregate Mounts

The UNIX operating system, in most implementations, only performs file system updates (syncs) to file systems that it finds in the local mount table. This immediately causes a problem for filesets which are exported to DFS but not mounted locally. To overcome this difficulty, DFS performs a special mount of aggregates as they are attached - just to get them into the mount table. Specifically:

- No provision is made for allowing access below these mount points since the only goal is to get an entry into the mount table.
- Aggregates - not individual filesets - are mounted. The mount points used are of the form:

```
/opt/dcelocal/var/dfs/aggrs/<aggregate-name>
```

- The mount is performed with a new file system type, MOUNT\_AGFS.
- DFS supplies code that implements this file system type (**vfs** and **vnode** operations). A few **vnode** operations are marginally implemented to allow **/bin/ls** operations within the mount directory listed above; most return ENOSYS since they should never be used. The **vfs** operations are implemented on top of the aggregate operations (*ag\_xxx()*) provided by the LFS. They are:

**\_root** Returns a "synthesized" vnode with the V\_ROOT flag set.

**\_statfs** Calls *ag\_stat()* to obtain a statfs structure that gives the space usage on the aggregate.

**\_sync** Calls *ag\_sync()*, with a sync type of AG\_SYNC\_FILESYS, to actually perform the *sync*.

**\_unmount** Updates internal tables to reflect the fact that this aggregate is no longer "mounted".

- These "pseudo-mounts" exist independently of any local mounts that might already exist or be created later for filesets on the aggregate.

### 14.4 Aggregate Array Functions

Descriptions of the members of the **struct aggrops** function array follow. File system independent code above the LFS enforces the requirement that with the exception of *ag\_stat()*, *ag\_vollInfo()* and *ag\_sync()*, these operations can only be issued by ROOT (the local super user).

**NAME**

*ag\_hold* — Increment the aggregate reference count

**SYNOPSIS**

```
int ag_hold(  
    /* IN */ struct aggr *aggd  
);
```

**ARGUMENTS**

*aggd*                   Aggregate descriptor pointer.

**DESCRIPTION**

Increment the reference count of the aggregate described by the *aggd* pointer.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[*error\_status\_ok*] This function always returns success.

**ERRORS**

None.

**NAME**

ag\_rele — Decrement an aggregate's reference count

**SYNOPSIS**

```
int ag_rele(  
    /* IN */ struct aggr *aggd  
);
```

**ARGUMENTS**

*aggd*                   Aggregate descriptor pointer.

**DESCRIPTION**

Decrement the reference count of the aggregate associated with the *aggd* pointer. The aggregate reference count is normally incremented by the *ag\_hold()*, *agg\_Lookup()* and *agg\_Attach()* functions.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**NAME**

`ag_lock` — Lock access to an aggregate

**SYNOPSIS**

```
int ag_lock(  
    /* IN */ struct aggr *aggr,  
    /* IN */ int type  
);
```

**ARGUMENTS**

*aggr* Aggregate descriptor pointer.

*type* Type of lock to be obtained. Allowable values are **READ\_LOCK**, **WRITE\_LOCK** and **SHARED\_LOCK**.

**DESCRIPTION**

Lock access to an aggregate.

**DISCUSSION**

This operation is currently unused by DFS. The semantics of this function are unspecified. It is implementation-dependent. If implemented, this operation could utilize the standard DFS lock primitives as follows:

```
{  
    code = error_status_ok;  
    {  
        if (type == READ_LOCK) {  
            lock_ObtainRead(&aggrp->a_lock);  
        }  
        else if (type == WRITE_LOCK) {  
            lock_ObtainWrite(&aggrp->a_lock);  
        }  
        else if (type == SHARED_LOCK) {  
            lock_ObtainShared(&aggrp->a_lock);  
        }  
        else code = EINVAL;  
    }  
    return code;  
}
```

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:  
[error\_status\_ok] This function was successful.

**ERRORS**

[EINVAL] If *type* is neither **READ\_LOCK**, **WRITE\_LOCK** or **SHARED\_LOCK**.

**NAME**

ag\_unlock — Unlock access to an aggregate

**SYNOPSIS**

```
int ag_unlock(
    /* IN */ struct aggr *aggr,
    /* IN */ int type
);
```

**ARGUMENTS**

*aggr* Aggregate descriptor pointer.

*type* Type of lock obtained at call to *ag\_lock()*. Allowable values are **READ\_LOCK**, **WRITE\_LOCK** and **SHARED\_LOCK**.

**DESCRIPTION**

Unlock access to an aggregate.

**DISCUSSION**

This operation is currently unused by DFS. The semantics of this function are unspecified. It is implementation-dependent. If implemented, this operation could utilize the standard DFS lock primitives as follows:

```
{
    code = error_status_ok;
    {
        if (type == READ_LOCK) {
            lock_ReleaseRead(&aggrp->a_lock);
        }
        else if (type == WRITE_LOCK) {
            lock_ReleaseWrite(&aggrp->a_lock);
        }
        else if (type == SHARED_LOCK) {
            lock_ReleaseShared(&aggrp->a_lock);
        }
        else code = EINVAL;
    }
    return code;
}
```

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

**ERRORS**

[EINVAL] If *type* is neither **READ\_LOCK**, **WRITE\_LOCK**, or **SHARED\_LOCK**.

**NAME**

**ag\_stat** — Returns statistics on the given aggregate

**SYNOPSIS**

```
int ag_stat(  
    /* IN */ struct aggr *aggr,  
    /* INOUT */ struct ag_status *astatd  
);
```

**ARGUMENTS**

<i>aggr</i>	Aggregate descriptor pointer.
<i>astatd</i>	Pointer to an aggregate status structure into which the information will be read.

**DESCRIPTION**

Provide the caller with the statistics block associated with the aggregate described by **aggr** pointer, placing the results at location **astatd**.

**DISCUSSION**

Any status fields not explicitly set should be zeroed.

The static portion of the status can be simply copied from the `aggr->a_stat_st` to `astatp->ag_st`. (This static status is maintained in the aggregate structure.)

The dynamic portion of the status in `astatp->ag_dy` is computed as appropriate for the particular LFS implementation. Fields with no significance (namely, **minFree**, **spares**) should be zeroed.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.



**NAME**

ag\_volCreate — Create a new fileset within an aggregate

**SYNOPSIS**

```
int ag_volCreate(
    /* IN */ struct aggr *aggd,
    /* IN */ struct afsHyper *avolidp,
    /* IN */ struct vol_status *avolstatp,
    /* IN */ long aflags
);
```

**ARGUMENTS**

<i>aggd</i>	Aggregate descriptor pointer.
<i>avolidp</i>	Pointer to the ID of the new fileset.
<i>avolstatp</i>	Pointer to a fileset status block with initial values in fields that require them.
<i>aflags</i>	Miscellaneous flag bits. The only bit which has significance is AGGR_CREATE_ROOT.

**DESCRIPTION**

Create and attach (for DFS export) a fileset within the aggregate associated with the descriptor located at **aggd**. The initial fileset header is derived from the image provided through the **avolstatp** pointer, which has various fields set.

**DISCUSSION**

The following bit is defined to have significance and is defined as:

```
/* defined bits within the flags argument */
#define AGGR_CREATE_ROOT 1
```

See Appendix D on page 317 for further information about this function.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

**ERRORS**

[EINVAL]	If not exactly one of VOL_RW and VOL_READONLY are set in the initial status <b>.v_states</b> field. The new fileset has not been created in this case.
[ENOSPC]	There is not enough space to create the fileset. The fileset has not been created in this case.
[EDQUOT]	The disk quota has been exceeded. The fileset has not been created in this case.

## NAME

ag\_volInfo — Enumerate the filesets within the given aggregate

## SYNOPSIS

```
int ag_volInfo(
    /* IN */ struct aggr *aggr,
    /* IN */ long index,
    /* OUT */ struct volume *vold
);
```

## ARGUMENTS

<i>aggr</i>	Aggregate descriptor pointer that is being enumerated.
<i>index</i>	Index of fileset within the aggregate that is to be examined for this call.
<i>vold</i>	Fileset descriptor into which to copy information about the fileset located at slot <b>index</b> within the aggregate.

## DESCRIPTION

Determine if there is a fileset at position **index** within the aggregate described by the **aggr** pointer. If there is, put information about it in the fileset descriptor (**vold**). If there is no fileset at position **index**, return **ENOENT**. The first valid fileset index within an aggregate is zero. To enumerate all filesets within an aggregate, *ag\_volInfo()* is called repeatedly, with **index** starting at zero and incremented at each subsequent call. The iteration completes when *ag\_volInfo()* returns **VOL\_ERR\_EOF**.

*ag\_volInfo()* can be compared with the fileset operation, *vol\_getstatus()* (see Section 15.14 on page 355). Both return information about a fileset. But when the former is called, the aggregate is not necessarily attached, and consequently the aggregate op does not look up the fileset in the Fileset Registry, nor make any attempt to gather dynamic information about how the fileset is being used.

## DISCUSSION

For information relative to the fields within the fileset structures, see Appendix E on page 321.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

## ERRORS

[ENOENT] There is no fileset at the position selected by the *index* argument.

[VOL\_ERR\_EOF] All filesets on the aggregate have index values smaller than the one supplied on this call. The calling software should stop looping through the filesets, as they have all been enumerated.

## SEE ALSO

*vol\_getstatus()*.

**NAME**

ag\_detach — Detach the given aggregate

**SYNOPSIS**

```
int ag_detach(  
    /* IN */ struct aggr *aggd  
);
```

**ARGUMENTS**

*aggd* Aggregate descriptor pointer.

**DESCRIPTION**

Detach the aggregate associated with the **aggd** pointer, taking all its filesets off-line.

**DISCUSSION**

This call should release, via *osi\_Free()*, any private storage that was allocated by the *af\_attach()* function. In this event, the aggregate structure (**struct aggr** field *.a\_fsDatap*) should be zeroed as well. See Section 13.2.6 on page 270 for information pertinent to the fields in this structure.

Other file system dependent clean up might be required as well, keeping in mind that filesets on the aggregate might still be mounted for local use.

**Note:** Although Transarc's Episode file system guarantees that an aggregate will not be *ag\_detach()*'d until it is no longer in use (by either DFS or locally mounted filesets), other implementations of the LFS and **mount** and **umount** commands need not behave in this manner.

Consult Appendix F on page 323 for further information.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

**ERRORS**

[ENXIO] Most failures that are I/O related should be returned with this code.

**SEE ALSO**

Section 11.8.3 on page 241 defines the signature for the *osi\_Free()* fP function.

## NAME

ag\_attach — Attach the given aggregate

## SYNOPSIS

```
int ag_attach(
    /* IN */ dev_t dev,
    /* IN */ struct vnode *bdevvp,
    /* IN */ u_long flags,
    /* IN */ caddr_t data,
    /* OUT */ opaque *fsdatap,
    /* OUT */ long *fsdatalen
);
```

## ARGUMENTS

<i>dev</i>	Device representing the partition on which the aggregate is located.
<i>bdevvp</i>	Vnode for the block device of this aggregate that is being attached. The device number is obtained from this vnode.
<i>flags</i>	Various bits. If the flag <code>AGGR_ATTACH_NOEXPORT</code> is present, this indicates that the aggregate is not actually being exported to DFS (Instead, one of its filesets is being mounted locally).
<i>data</i>	A data area (in user space), holding filesystem-specific data about the aggregate. In <b>DCE LFS</b> this parameter is not used. For UFS aggregates it holds the fileset ID of the aggregate's one fileset, and the name of the directory on which the corresponding filesystem is mounted. If present, it should be copied in via <i>copyin()</i> .
<i>fsdatap</i>	Pointer to a place to put a filesystem-specific control block.
<i>fsdatalen</i>	Length of the filesystem-specific control block, for use in freeing it, in case this is done in the course of aborting <i>ag_Attach()</i> rather than (as is usually the case) by <i>ag_Detach()</i> .

## DESCRIPTION

Attach the aggregate specified, putting all its filesets on-line.

## DISCUSSION

Consult Appendix F on page 323 for further information.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

## ERRORS

[ENXIO] Most failures that are I/O related should be returned with this code. For instance, this failure indicates that the aggregate is in an inconsistent (corrupt) state and needs to be run through a disk salvager (such as FSCK) before it can be accessed.

**NAME**

ag\_sync — Sync the given aggregate

**SYNOPSIS**

```
int ag_sync(
    /* IN */ struct aggr *aggr,
    /* IN */ int syncType
);
```

**ARGUMENTS**

*aggr* Aggregate descriptor pointer.

*syncType* Describes the type of sync to be done, choose from one of: **AGGR\_SYNC\_FILESYS**, **AGGR\_SYNC\_COMMITMETA** or **AGGR\_SYNC\_COMMITALL**.

**DESCRIPTION**

Sync the state of the aggregate to the permanent storage, according to the *syncType* parameter:

**AGGR\_SYNC\_FILESYS**

All dirty data on the aggregate (and possibly all aggregates of the same type) is (scheduled to be asynchronously) written to permanent storage, consistent with the behavior of **sync(2)**. This is useful to periodically flush dirty data for filesets that have no local mount point, and thus, are not affected by the system's periodic calls to **sync(2)**.

**AGGR\_SYNC\_COMMITMETA**

All dirty meta-data (for instance, file status) is written to permanent storage. The call does not return until the I/O has completed (It is synchronous.).

**AGGR\_SYNC\_COMMITALL**

All dirty data (both meta-data and user data) is written to permanent storage. The call does not return until the I/O has completed.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

**ERRORS**

[EINVAL] The **syncType** argument is not one of **AGGR\_SYNC\_FILESYS**, **AGGR\_SYNC\_COMMITMETA** or **AGGR\_SYNC\_COMMITALL**.



## Information Pertinent to `ag_volCreate()`

In the process of creating and attaching a fileset on an aggregate, the following fields from the supplied initial volume status are stored on disk for the newly created fileset. On-disk fileset status fields not mentioned here should be set to zero.

### D.1 Static Status

From the static status portion, `statusp->vol_st`:

`.volName`

`.volId`            The `volIdp` argument is used instead of this.

`.parentId`

`.cloneTime`

`.states`            The following fileset state bits are sampled and stored on-disk.

```
VOL_READONLY        \ exactly 1 of these
VOL_RW              /  should be set
VOL_IS_COMPLETE
VOL_DELONSALVAGE
VOL_ZAPME
VOL_CLONEINPROG
VOL_REPFIELD
VOL_TYPEFIELD
VOL_REPSERVER_MGD
VOL_NOEXPORT
VOL_IS_REPLICATED
VOL_MOVE_SOURCE
VOL_MOVE_TARGET
```

### D.2 Dynamic Status

From the dynamic status portion, `.vol_dy`

`.creationDate`    If the seconds (`.sec`) portion of this is 0, the fileset's creation time is set to the current time.

Observation: although the hi-level Episode code samples this and leaves it in the memory structure, the low-level code seems to ignore it and set the on-disk field to "now".

`.updateDate`

`.accessDate`

`.copyDate`

`.volversion`

`.backupId`

.cloneId  
 .llBackId  
 .llFwdId  
 .allocLimit  
 .visQuotaLimit  
 .unique            Observation: although the hi-level Episode code seems to sample *.unique* and leave it in the memory structure, the low-level code seems to ignore it and set the on-disk field to 0.  
 .statusMsg

### D.3 Other Items

The ID of the new fileset is taken from the *volIdp* argument.

If the *AGGR\_CREATE\_ROOT* flag is specified, a properly initialized root directory (containing "." and "..", owner and group taken from the callers credentials, without ACL or property lists) is created. Otherwise, the fileset is left in an inconsistent state (namely, waiting for a clone or restore) and the *VOL\_DELONSALVAGE* volume state flag should be set in both the supplied volume status structure and the on-disk storage for this fileset.

The *.vol\_dy.index* field in the supplied status structure is set to the value that, if passed as the index argument to the aggregate *ag\_volInfo()* operation, would select this newly created fileset.

Observation: this *.vol\_dy.index* field may in fact be obsolete or un-needed.

If, for some reason, the operation fails with the new fileset left in an intermediate state (for instance, cleanup was not possible), the *VOL\_DELONSALVAGE* flag should be set in both the supplied status *.vol\_st.states* word and in on-disk storage. (Such an intermediate state, while sometimes unavoidable, is obviously undesirable.)

### D.4 Attaching the New Fileset

If a new fileset is created and left on the aggregate (even if in an inconsistent state due to an error), this operation must export it to DFS via the DFS function *vol\_Attach()* as follows:

```
/* ... the fileset has been created */
error = vol_Attach(volIdp, statusp, aggrp, &vol_ops_vector);
/* ... and when all done */
return(error);
```

**Note:** This is NOT the VFS+ volume *vol\_attach()* operation, but a different one. DFS has chosen to implement this up-call function in *ag\_volCreate()* rather than having the code that called *ag\_volCreate()* handle it.

The first three arguments are the same as those passed to *ag\_volCreate()* although in a different order. The fourth is the address of the volume operations entypoint vector. This call will itself make calls to the volume operations *vol\_setdystat()* and *vol\_attach()* fP functions. Note that there is no need to examine the error code it returns (specifically: if it fails, the fileset is not deleted or marked inconsistent).



The *vol\_Attach* call, which is not listed in any of the **DCE DFS** header files, has the following signature:

```
vol_Attach(  
    afsHyper          *volIdp,  
    struct vol_status *statusp,  
    struct aggr       *aggrp,  
    struct volumeops  *volopsp  
);
```

No checking is performed to see if a volume with the specified volume ID or name already exists. The LFS is thus relying upon DFS to avoid a problem in this instance. Presently, the *vol\_Attach()* up-call will fail if an aggregate with the given ID already exists.

**Note:** For the sake of per-aggregate consistency, it seems reasonable to suggest that an LFS implementation check within an aggregate for duplicate fileset names and IDs at creation (returning, perhaps, EEXIST).



## Information Pertinent to `ag_volInfo()`

### E.1 Volume Structure Fields

Fields within the volume structure (**struct volume**) should be filled in as follows. (See Section 13.3.14 on page 285 for information on the contents of this structure):

- `.v_paggrp`        Set to the **aggrp** argument passed to this call. (That is, this field points to the **aggr** structure for the aggregate that the fileset resides on.)
- `.v_volOps`        Set to point at the **struct volumeops** vector of volume operations provided by the LFS.

### E.2 Static Status

Fields within the volume `.v_stat_st` structure should be filled in from the fileset status stored on disk as follows. Any field not mentioned should be zeroed.

- `.volName`
- `.volId`            The *volIdp* argument is used instead of this.
- `.parentId`
- `.cloneTime`
- `.states`            Only the following fileset state bits are set and stored on-disk. (Others should be zeroed.)
  - VOL\_READONLY        \ exactly 1 of these
  - VOL\_RW              /    should be set
  - VOL\_IS\_COMPLETE
  - VOL\_DELONSALVAGE
  - VOL\_ZAPME
  - VOL\_CLONEINPROG
  - VOL\_REPFIELD
  - VOL\_TYPEFIELD
  - VOL\_REPSERVER\_MGD
  - VOL\_NOEXPORT
  - VOL\_IS\_REPLICATED
  - VOL\_MOVE\_SOURCE
  - VOL\_MOVE\_TARGET
- `.type`            The states of the VOL\_READONLY and VOL\_RW flags from the `.v_states` field are replicated here.



## Information Pertinent to *ag\_[de,at]tach()*

### F.1 Making an Aggregate Available

The first operation that must be performed in order for an aggregate to be available for use by DFS is to attach it to DFS, putting all its filesets on-line. No other aggregate (or volume) operations will occur on an aggregate until after this has been done. On a successful return from the *ag\_attach()* call, an aggregate struct (**struct aggr**) will be allocated by higher level software.

### F.2 Private Storage

The *af\_attach()* call can optionally allocate and initialize a block of private storage which will then be made available to all subsequent aggregate operations. This storage must be allocated via a call to *osi\_Alloc()*. A pointer to it along with its size are returned in the **fsdatap** (private data) and **fsdatalen** (private data length) arguments. (The length is required so that higher level software can release the storage if an error occurs later in the attach algorithm.) This pointer is deposited in the **aggr** structure *.a\_fsDatap* field, from where it can be retrieved on subsequent aggregate operations. If no such storage is allocated, return values of NULL and 0 should be supplied for **fsdatap** and **fsdatalen**, respectively.

Section 11.8.3 on page 241 defines the signature for *osi\_Alloc()*.

### F.3 Aggregate Structure

Aggregates currently attached (see Note below, or containing filesets mounted locally) are known to DFS by means of an **struct aggr** structure that is remembered. This *ag\_attach()* call is made by DFS whenever it encounters an aggregate that it doesn't already know about. Once this call has been made for a fileset, it will not be made again until an intervening *ag\_detach()* occurs.

**Note:** The **AGGR\_ATTACH\_NOEXPORT** flag is not normally issued by DFS. It's only user is the local mount command for the Episode LFS, which invokes the DFS attach code to allocate an **aggr** structure. There doesn't appear to be any use to which an LFS can put this flag: it will be present if the first time DFS encounters an aggregate it is for a local mount (as opposed to a DFS export). This is either:

- a. Merely an artifact of the Episode implementation.
- b. DFS requires that it "have" appropriate **aggr** structures for any filesets mounted locally regardless of whether or not they are currently attached by DFS. However, this isn't being done for UFS file systems.

At the current time, the first alternative is assumed. If that turns out not to be the case, an interface into DFS for use by the LFS **mount** and **unmount** commands will be needed by any vendor-supplied LFS.



## Fileset (Volume) Operations Interface

This chapter describes the facilities provided for the VFS+ Fileset Operations. It provides an overview, discusses fileset types, clones, fileset and clone requirements, fileset indices, LFS modification of fileset status, zero link count files, quotas, anode generation numbers, file identifiers, vnode to LFS association, and lastly, fileset and fileset registry operations.

### 15.1 Overview

This chapter describes the portion of the VFS+ interface which allows operations on filesets, as introduced in Section 11.6 on page 239 that must be supported by a **DCE LFS**. The *Fileset Registry* is a global kernel table which contains descriptions of each fileset local to the machine. A fileset descriptor includes a pointer to a **struct volumeops** function array, namely those operations that can be performed on the fileset. The associated declaration may be found in Section 13.3.1 on page 274.

An agent that wishes to perform a fileset operation will typically use the fileset ID to hash into the *Fileset Registry* and thus acquire the appropriate fileset descriptor. It will then perform a sequence of operations, calling functions obtained from the descriptor's function array.

In **inode**-based file systems, a fileset represents an **inode** table, and individual files are accessed via **inode indexes**. In **DCE LFS**, a fileset represents an **anode** volume list (avl), and individual files are accessed via **anode indexes**. In either case, agents such as the fileset server will typically iterate over all files and directories in a fileset by iterating over all possible indexes (non-negative integers), stopping when the highest index in the fileset has been processed. It should be borne in mind that the order of iteration, from low index to high, is completely unrelated to any order that might have been derived from the directory hierarchy. There is no guarantee that directories will be processed before the files that are in them, or vice versa, although the root directory will generally be processed first.

A description of the operations available within the fileset function array appears below. The array is an exported data type. (The associated declaration is found in Section 13.3.1 on page 274). It is reproduced here for convenience.

```
struct volumeops {
    /* per-fileset operations */
    int (*vol_hold)();
    int (*vol_rele)();
    int (*vol_lock)();
    int (*vol_unlock)();
    int (*vol_open)();
    int (*vol_seek)();
    int (*vol_tell)();
    int (*vol_scan)();
    int (*vol_close)();
    int (*vol_destroy)();          /* actually whole-fileset */
    int (*vol_attach)();
    int (*vol_detach)();
    int (*vol_getstatus)();
    int (*vol_setstatus)();

    /* per-file operations */
    int (*vol_create)();
};
```

```
int (*vol_read)();
int (*vol_write)();
int (*vol_truncate)();
int (*vol_delete)();
int (*vol_getattr)();
int (*vol_setattr)();
int (*vol_getacl)();
int (*vol_setacl)();

/* more whole-fileset operations */
int (*vol_clone)();
int (*vol_reclone)();
int (*vol_unclone)();

/* vnode lookup operations */
int (*vol_vget)();
int (*vol_root)();
int (*vol_isroot)();

/* more per-fileset operations */
int (*vol_getvv)();
int (*vol_setdystat)();
int (*vol_freedystat)();
int (*vol_setnewvid)();

/* another per-file operation */
int (*vol_copyacl)();

/* per-fileset operations */
int (*vol_concurr)();
int (*vol_swapids)();
int (*vol_sync)();
int (*vol_pushstatus)();

/* per-file operations */
int (*vol_readdir)();
int (*vol_appenddir)();

/* per-fileset operations */
int (*vol_bulksetstatus)();
int (*vol_getzlc)();

/* another per-file operation */
int (*vol_getnextholes)();

/* yet another per-fileset operation */
int (*vol_deplete)();
};
```



### 15.1.1 Classes of Fileset Operations

There are four classes of fileset operations, as delineated by the comments in the above structure declaration:

1. **Per-fileset operations.** These operations are used to get and set fileset status and iterate over all files in the fileset.
2. **Whole-fileset operations.** These operations are similar to the per-fileset operations in that they operate on an entire fileset, but they differ in that they embed iterations over every low-level object in a fileset, even those corresponding to **ACLs** that are not “files” in the usual sense. These operations are used either to destroy a fileset or to manipulate the clone status of filesets.
3. **Per-file operations.** These functions are used to create and delete files, inspect or change the contents of a file, and get and set the file attributes.
4. **Vnode lookup from fileset and file information.** Return a vnode pointer associated with a given fileset and file ID.

The signatures for the *fileset* operation functions appear in the section of this chapter labeled *Fileset Array Functions*. It is Section 15.14 on page 340.

### 15.1.2 Fileset Registry Functions

The following functions manipulate the *fileset registry*:

- volreg\_Enter()* Given fileset and aggregate IDs, along with a pointer to the associated fileset information, create an entry in the Fileset Registry and insert the given information into it.
- volreg\_Delete()* Delete the Fileset Registry entry (if any) that corresponds to the given fileset and aggregate ID pair.
- volreg\_Lookup()* Given a file ID, this function returns a pointer for the fileset descriptor and, optionally, the vnode corresponding to the file.

The signatures for the *fileset registry* functions appear in the section of this chapter labeled *Fileset Registry Array Functions*. It is Section 15.15 on page 396.

## 15.2 Fileset Types Overview

By means of cloning and read-only replication, filesets will in general be grouped into related sets. Within these groupings, filesets can have different types (or roles) as indicated by the `VOL_TYPEFIELD` field of the fileset status *.states* field. In the following types, a grouping is referred to as a *group*. The types (or roles) are:

- VOL\_RW** The read-write fileset, of which there can be only one within a *group*.
- VOL\_READONLY** A read-only replica, of which there can be several within a *group* (at different sites). It's name is the same as the `VOL_RW` one, with a ".readonly" suffix. These filesets are managed by the replication servers. It is possible for a read-only replica to reside within the same aggregate as its read-write master.
- VOL\_BACKUP** A read-only clone of the read(write) fileset, created via the `fts clone` command. It's name is the same as the `VOL_RW` one, with a ".backup" suffix. Because it is a clone, it **MUST** reside within the same aggregate as its read-write master. There can only be one of these within a *group*.

**VOL\_TEMP** A temporary clone of the read-write master used in the process of moving or replicating a fileset.

Volume IDs corresponding to the first three of these types (RW, READONLY and BACKUP) are all allocated when a fileset is first created and are all stored in the FLDB. The fourth (TEMP) type uses a temporary ID since it is only visible to the algorithm performing the move.

### 15.3 Fileset Clone Algorithms

In the following discussion, the types from the previous section, Section 15.2 on page 327, are referred to. Also, throughout this document, the symbol "====>" will be used to denote an "is backed by" relationship.

The fileset cloning mechanism uses copy-on-write techniques to create a "new" fileset that is a snapshot of the current state of an existing one. Usually, the following steps are followed by DFS in the creation of a clone.

- a. An existing fileset, either a VOL\_RW (usually) or VOL\_READONLY one (when updating RO replicas), is to be cloned. Let this fileset be F1.
- b. A new fileset, F2, is created via *ag\_volCreate()*. (See Section 14.4 on page 311.) Since the AGGR\_CREATE\_ROOT flag is not specified to this call, this fileset starts out totally empty.
- c. A series of *vol\_clone(F2, F1, ..)* (description is found in Section 15.14 on page 369) calls are made to actually create clones on F2 for the files on F1.
- d. At this point, the DFS terminology speaks in terms of fileset F2 being backed by F1. That is: "F2 is backed by F1" (represented by F2====>F1). F1, the original fileset, is known as a backing fileset. Files on F1 are known as backing files.

The LFS should NEVER allow modifications to backing files. In this scenario, *only* F2 can be modified.

Throughout this document, the symbol "====>" is used to denote an "is backed by" relationship.

- e. In many (most, perhaps) scenarios, the desire was to create a read-only clone and have the original (F1 as in this discussion) fileset continue to be writeable.

To accomplish this, the fileset identities are swapped via a call to *vol\_bulksetstatus()* (described in Section 15.14 on page 392) that swaps just about every status field in the two filesets. Following this, the relationship has been changed to: "F1====>F2" where F1 is writeable and F2 isn't. Note that the newly created physical fileset (step (b)) is now F1 (it has F1's ID, and so forth.).

- f. This sequence can be repeated again, giving rise to a chain of related (cloned) filesets.
- g. Assume an initial state of "F1====>F2". Following steps (b) through (d), the following relationship holds: "F3====>F1====>F2".
- h. Following another identity swap between F3 and F1, a final result is arrived at, namely: "F1====>F3====>F2".

**Note:** Unless explicitly indicated in the detailed descriptions, the fileset status fields related to fileset clones (and replicas, and so on) should not be examined or used by the LFS since their use by DFS could change in the future.

### 15.3.1 Requirements on Cloning

The following requirements on cloning hold.

1. Under ordinary circumstances, a 3-deep chain of cloned filesets is the "worst" that DFS will ever create. The LFS, however, should be capable of supporting chains of *arbitrary*, or at least, somewhat larger than 3 depth.
2. The LFS must prevent any modifications to backing filesets, which are read-only. For example, given the following state: "F1 ==> F2 ==> F3" *only* F1 is writeable; F2 and F3 are read-only.

Although Transarc's Episode appears to return EIO in such cases, EROFS would seem to be a more appropriate error return.

**Note:** The above restriction only applies to modifications to the file system visible contents of a fileset (files, directories, ACLS, and so forth.) via either the **vnode** or **fileset** interface.

The fileset header of these backing filesets can be modified via the following, for example:

- `vol_setstatus()`
- `vol_swapids()`
- `vol_bulksetstatus()`, and so on.

The following operations are obviously permitted as well:

- `vol_delete()`
- `vol_unclone()`
- `vol_reclone()`.

**Note:** In actuality, there appears to be a number of cases in which Transarc's Episode will allow modifications to certain extended file status on these backing filesets. If so, then perhaps it is the case that DFS is not relying on this behavior.

3. A new clone (as opposed to a reclone operation) will never be made of a backing fileset. For example, given the following state:

```
"F1 ==> F2 ==> F3"
```

*only* F1 can be the starting point in step (a) above.

4. The act of creating a clone is assumed to be a relatively short operation ("small" number of seconds for a "reasonable" sized fileset) and consume a "small" amount of additional disk space. See the discussion *Looping Operations* in Section 15.10 on page 335.
5. Consult the *Copy-on-Write Impacts* section (Section 16.3 on page 421) in the **vnode** operations chapter for additional assumptions regarding cloning.
6. All files on a fileset, even those with a Zero-Link-Count, are cloned.

### 15.3.2 Uses for Clones

Among the various uses to which DFS puts clones, the following are typical.

1. The **fts backup** command creates a backup fileset (with a suffix of ".backup" and a type of VOL\_BACKUP) which can be backed up in parallel to on-going access to the original fileset or left in place as an on-line backup. The steps followed are precisely those detailed in the cloning steps (a) through (e) in Section 15.3 on page 328.
2. The **fts move** command moves a fileset to a new location. It accomplishes this via the following steps:
  - i. Create a temporary clone Ft of the original fileset F1. After the fileset identity swap, the result is "F1 ==> Ft".  
If F1 had already had a backing fileset F2 (perhaps a .backup fileset), the result would be: "F1 ==> Ft ==> F2".
  - ii. The clone Ft is copied to the new location, while on-going access to F1 is allowed.
  - iii. The temporary Ft is deleted.
  - iv. Changes to F1 that occurred during step (ii) are incrementally copied to the new location. During this hopefully short step, fileset F1 is un-accessible for normal access.
  - v. The FLDB is updated to reflect the new location for the fileset. Recall that this new location only applies to VOL\_RW and VOL\_BACKUP types. (Those types were discussed in Section 15.2 on page 327.)
  - vi. The original fileset is deleted. If there is a \.backup fileset to go with it, it is deleted as well.

**Note:** **It was never copied to the new location! This might be changed in the future.**

3. A replicated read-only fileset ( call it Fro) is updated from its master fileset (call it Frw) as follows. Initially, Frw and Fro exist at different (typically) sites.
  - i. The contents of Frw are pushed to a staging fileset co-resident with it. This may be accomplished via a clone created against Frw. Following the usual identity swap, the result is:
 

```
"Frw ==> F1"           Master site
```

 F1 will either have a temporary ID or the read-only ID from the FLDB.
  - ii. A clone is then created at the replica site against Fro. This time, however, the identities are not swapped and the new fileset (F2) remains writeable. The result is:
 

```
"F2 ==> Fro"           Replica site
```

 F2 will have a temporary fileset ID.
  - iii. The contents of F1 (backing Frw) are copied to the new fileset F2. Following this copy, the identities of Fro and F2 are swapped so that references to the read-only replica are now directed to the newly updated fileset. This swap is accomplished via a *vol\_swapids()* fileset call.

**Note:** Unless explicitly indicated in the detailed descriptions, the fileset status fields related to fileset clones (and replicas, and so on) should not be examined or used by the LFS since their use by DFS could change in the future.

### 15.3.3 Some Fileset and Clone Requirements

In addition to the requirements detailed elsewhere in these specifications, DFS makes a number of general assumptions regarding filesets and cloning. They are:

1. Several modes of fileset usage are expected. In some scenarios, small numbers of very large filesets will prevail while in others large numbers of small filesets will prevail.
2. An aggregate is expected to be capable of supporting a "large" number of filesets. For suitably large aggregates, a lower limit in the thousands or tens-of-thousands seems reasonable. (Therefore: if there is a limit, it should scale in some regard with the aggregate size.)
3. Although it has been mentioned elsewhere, this bears repeating. An LFS must dis-allow any modifications (except for those via *vol\_xxx()* fileset operations) to files on a read-only fileset: either a *.backup* clone or a *.readonly* replica.

## 15.4 Fileset Indices

Filesets on an aggregate are identified by a zero-based index about which DFS makes the following assumptions.

- A fileset index is visible outside the LFS in the following ways.
  - The index of the fileset to be queried is passed as an argument to the *ag\_volInfo()* call.
  - The *ag\_volCreate()* call returns fileset index for a newly created fileset in the volume *.vol\_dy.index* field.
  - A fileset index is returned in the **Txvattr** *.backingIndex* field by the *vn\_getattr()* and *vol\_getattr()* calls.

Since it is stored internally within DFS, these indices must remain consistent over time. (Strictly speaking, it ought to be "safe" to allow a fileset's index to change each time its aggregate is attached ... although there could well be problems with such an approach.)

- The general expectation is that when a fileset is created (by *ag\_volCreate()*), the lowest *unused* index is assigned to it. Although using monotonically increasing indexes might seem desirable (not re-using the indices of deleted filesets), this is not really in keeping with the manner in which *ag\_volInfo()* incrementally steps through the filesets on an aggregate.
- A fileset's index stays with it *forever*. It is not changed by the *vol\_clone()*, *vol\_reclone()*, *vol\_swapids()*, *vol\_bulksetstatus()* or *vol\_setstatus()* operations. Although filesets may change or swap their DFS identity (fileset ID), their fileset index remains the same.

## 15.5 LFS Modification of Fileset Status

For the most part, the LFS neither interprets nor modifies volume status fields. The exceptions, which are also indicated under the appropriate aggregate and fileset operations, are listed here.. Refer to Section 13.3 on page 274, *Fileset Data Types*, for details of the fields and flags discussed in this section.

- The VOL\_DELONSALVAGE flag in *vol\_stat\_st.states* is modified by a number of operations.
- The VD\_RDONLY flag is set in a newly created vnode if either the VOL\_READONLY flag is set in its fileset *vol\_stat\_st.states*, the fileset was mounted read-only or the file is backing (clone) another file in a copy-on-write relationship.

**Note:** While the above point is correct, it does not really apply to discussing fileset status modifications.

- The *vol\_stat\_dy .updateDate* field is updated whenever a fileset's status is changed. This includes file creation and deletion, fileset status operations, the clone operations and whenever the volume version is incremented.
- The *vol\_stat\_dy .accessDate* field is updated whenever the *.updateDate* field is.
- The *vol\_stat\_dy .volversion* is advanced whenever any modification is made to the data or status of any file contained within that fileset.

**Note:** Transarc's Episode advances the volume version by simply incrementing it. Other approaches, such as setting it to the current time (given enough precision), would seem to be acceptable as well ... as long as the version changes and never decreases (64-bits) in value.

- The *vol\_stat\_dy .allocUsage* and *.visQuotaUsage* fields are both adjusted as the size of a fileset changes.

## 15.6 Zero Link Count Files

Ordinarily, files whose link count is zero are actually deleted by the LFS when their vnode reference (usage) count goes to zero. (See the discussion under the *vn\_inactive()* vnode operation, Section 16.7 on page 453.)

In certain circumstances, DFS will *protect* a deleted file by artificially holding its vnode until it is certain that any remote clients are through with it. Following a crash or shutdown, any remaining zero-link-count files (their vnodes were still held at crash or shutdown time) must be preserved by the LFS in order to give DFS a chance to re-adopt them once the fileset is attached again. Any FSCK-like programs that are run must also preserve these files.

**Note:** These zero-link-count files are **copied** by the *vol\_clone()* operation into the cloned fileset.

**Note:** It is possible for a zero-link-file to have its link count incremented before it is actually deleted, returning it to a "normal" status. (Although Transarc's Episode allows for this, it is not known how this might actually occur ... if it ever does.)

When a read-write fileset is attached, DFS makes a series of *vol\_getzlc()* volume calls to identify these files with zero link counts. This call returns held vnodes. When and if DFS determines that the file can safely be deleted, it releases the vnode which, as discussed above, actually triggers the deletion.

**Note:** For a readonly replica or clone, DFS never calls `vol_getzlc()`; these zero-link-count files remain intact until they are explicitly deleted (via `vol_delete()`) or until their fileset is destroyed.

Additional discussion can be found under the `vn_inactive()`, `vol_delete()`, `vol_setattr()` and `vol_getzlc()` operations.

An obvious difficulty here is that there is no guarantee that, following a re-boot, DFS will actually trigger the eventual reclamation of these files and the space they occupy. (Specifically, DFS may not be started again.) Although this problem is outside the scope of DFS, one can argue that DFS should provide the framework of a solution. One possible way of dealing with this issue is as follows.

1. Only give this zero-link-count treatment (including retention during FSCK) to files that have actually been accessed by DFS. (Perhaps their fileset has been exported; perhaps the file has been operated on by `vol_xxx()` operations such as `vol_vget()`.)
2. Employ some sort of (perhaps LFS specific) timeout mechanism to ensure the recovery of such files if DFS has not come along within a "reasonable" amount of time.

**Note:** Although the intent here seems desirable (preserving deleted files until they are no longer in use), practical concerns seem to get in the way. The biggest problem is that the DFS layer, in order to avoid tying up all the vnodes in a system, will only artificially hold up to 220 files at most. Although this should be sufficient in some situations, it will not be enough for large servers.

## 15.7 Fileset Quotas

There are two measures of the amount of disk space being used by a fileset.

**allocated**            The amount of disk space ACTUALLY in use by a fileset.

**visible**             The amount of disk space LOGICALLY used by a fileset. This is what is returned by the `stats()` vfs operation.

Normally, these these measures will be the same. For cloned filesets that are sharing disk space via a copy-on-write mechanism, these measures can differ as follows.

- A given disk block can only show up in the *allocated* measure for a single fileset, regardless of how many filesets are sharing it via copy-on-write.

Given a disk block that is shared (copy-on-write) by multiple filesets, DFS does not specify which fileset's *allocated* usage measure that block show up under.

**Note:** This is one possible mechanism by which a particular LFS (and its cloning algorithms) could be distinguished from the Transarc LFS. Under the Transarc LFS, shared (copy-on-write) blocks appear *allocated* under the backing fileset.

This *allocated* measure is the one that is important from a *how full is this physical aggregate* aspect.

- A given disk block will show up in the *visible* measure of EVERY fileset sharing it via copy-on-write. In one sense, this *visible* measure indicates how much physical disk space a fileset could require if all its copy-on-write linkages were broken.

This *visible* measure is the one that is normally of most use to users controlling the usage of their filesets.

- A fileset's *visible* usage measure will, therefore, always be equal to or larger than its *visible* measure (larger than for a cloned fileset; NEVER less than). Considering all the filesets on an aggregate, the sum of all the *allocated* usage measures will be no greater than the aggregate size while the sum of all the *visible* usage measures may be larger than the aggregate size.

There are several fileset status fields that are related to a fileset's disk usage. Above the LFS, they are represented as 64-bit byte counts. The first two are input to the LFS ("quotas"). The latter two are computed and maintained by the LFS and returned to the higher level's ("usage").

- .allocLimit**        The maximum allowable size (quota) that a fileset's *allocated* usage can grow to.
- .visQuotaLimit**    The maximum allowable size (quota) that a fileset's *visible* usage can grow to.
- .allocUsage**        The *allocated* measure of the disk space in use by a fileset. The *visible* measure of the disk space in use by a fileset.

Operations that would cause the **.allocLimit** or **.visQuotaLimit** values for a fileset to be exceeded fail with an EDQUOT error. The Super User (ROOT) enjoys no special privilege in this regard.

## 15.8 Anode Generation Numbers

**Note:** In the following, "inode" means "anode".

Each inode has an associated generation number which, in combination with an inode index, uniquely identifies a particular instantiation or use of an inode.

An LFS has freedom in how it generates these generation numbers. For example, it can use either a per-inode counter or a per-fileset counter (Transarc's Episode uses the fileset **vol\_stat\_dy.unique** field for this). The key behavior is that each time an inode is re-used (taking into account any inode table shrinkage and re-growth), a different generation number must be used for it.

The generation number for the root directory on a fileset should always be set to 1.

Within DFS, these inode generation numbers are referred to as "uniquifiers" and reside within the *.Unique* field of an **afsFid** structure.

## 15.9 File Identifiers (afsFids)

DFS identifies files via an **afsFid** structure. The fields within an **afsFid** have the following significance to an LFS.

- .Vnode*                Represents a file.
  - .Unique*              An anode (inode) index and generation number. An LFS is free to use whatever encoding is appropriate for it. There appears to be an assumption made that the root directory on a fileset has the lowest (from the set of legal or possible values) inode index. Consult the description under *vol\_open()* (Section 15.14 on page 345) on the setting of the *vol\_handle.index* field.
- A *.Unique* field of -1 is interpreted as a don't care value; it is not required to match the inode generation number currently found on the file.

**Note:** DFS does not appear to rely on this above behavior (*.Unique* of -1), although this is what Transarc's Episode does.



- .Cell*                   An identifier for the DFS cell.
- Although an LFS should not have any reason to examine this field, higher levels will generally set it to 0..1 (.high = 0, .low = 1). It seems safest to follow the same convention for any **afsFids** that are returned by the LFS.
- .Volume*                An identifier for the DFS fileset.
- In general, an LFS should not have any reason to examine this field.

Several of the fileset operations take an **afsFid** argument which identifies a file. In most (perhaps all) cases, these afsFids were originally returned to user space via an earlier *efs\_scan()* or *efs\_create()* call based on a canonical inode position argument (VOL\_ROOTINO specifies the root) passed to them. Normally, therefore, the LFS should always be able to locate the file corresponding to the specified **afsFid**. If for some reason, however, the object in question cannot be located, the operations should fail with either an ESTALE (if the *.Unique* field was the problem) or EIO (any other problem; namely, the given *.Vnode* index not in use) error.

## 15.10 Looping Operation Considerations

The *vol\_deplete()*, *vol\_clone()*, *vol\_reclone()* and *vol\_unclone()* operations all operate over an entire fileset. To avoid executing within the kernel for an arbitrarily long time and causing user-space RPC timeouts, they should periodically return with a code of ELOOP. This return value tells the user space code that this operation hasn't completed and that it should repeat the call. Since these calls do not accept an iteration argument, it may be necessary to record a "where to continue" marker in the LFS private fileset storage that is established in *vol\_setdystat()*. If this mechanism is used, it should be reset to "start at the beginning" at *vol\_open()* time and at the completion (that is, when something other than ELOOP is returned) of any of these operations.

**Note:** As currently exists, these timeouts will occur within a small number of minutes. This document recommends that these looping operations plan on spending no more than 30 seconds or so within the kernel. A certain amount of care may be required since some operations can clearly take longer than this (namely, deleting or cloning a VERY large file.) Some possibilities: processing portions of a file each iteration (as much as will "fit" within 30 seconds) or operating asynchronously with helper daemons to ensure that the process that issued the *vol\_xxx()* operation is able to return to user space every 30 seconds.

## 15.11 Vnode to LFS Association

During normal operation, there is an association between a vnode and the LFS-specific data structures which represent that object.

The fileset operations described in this specification enable an administrative utility to come in and make fundamental changes to the underlying LFS state. Depending on the precise implementation, some of these operations might require that the above mentioned association between a vnode and its LFS state (namely, location information) be broken. Examples of such operations are those involved in cloning, uncloning, recloning, deletion, restoration and movement. In the simplest case, the vnode usage count can be *driven* to zero so that it can be simply released or discarded. If this isn't possible (say a file is open locally), the vnode must be de-coupled from its LFS state -- which can change out from under the vnode. Such dis-association between a vnode and the LFS requires at least the following:

- Flushing any modified cache state back to the LFS.
- Leaving enough location information (inode index, generation number, fileset ID) in the vnode to enable a "hook up" with the LFS at a later time.

For similar reasons, any cached state (in the VM or naming cache, for example) must be considered in the face of volume operations that might modify or examine that state.

Once such a vnode has been dis-associated with the underlying LFS state, a corresponding mechanism is required to re-bind the vnode to the LFS (location information, status) before it is used again.

If the fileset no longer resides at the machine in question, an [ENODEV] error should be returned from the operation. If the fileset is still local to the machine but the file itself no longer exists, an [EIO] error should be returned from the operation.

**Note:** The precise error values returned in these cases appear not to be that crucial. DFS has suggested that returning [ESTALE] in both cases might be better.

### 15.11.1 Determining whether a Fileset is Local

Although there are numerous possible LFS-specific mechanisms for determining if a fileset is still local, Transarc's Episode does the following.

- Initializes an afsFid from the volume ID recorded in the vnode at dis-association time.
- Calls `volreg_Lookup(&fid, &volume_ptr)` in the DFS to let it make the determination. This call returns [ENODEV] if the fileset is not in the *fileset registry*.
- If this call succeeds, the hold count on the returned volume structure is incremented. The LFS must decrement it by calling through the fileset ops vector as follows:

```
(*volume_ptr->v_volOps->vol_rele)(volume_ptr);
```

**Note:** The fileset might reside within a different LFS entirely!):

If the volume structure is not needed by the LFS, this step can be avoided by calling `volreg_Lookup()` with NULL as the second argument -- indicating that the caller does not desire a returned volume structure.

- The anode index and generation number recorded in the vnode are used to look up the object on the fileset.

### 15.11.2 Handling Dis-sociations and Re-associations

There are several alternatives for scheduling these dis-associations and re-associations.

The dis-associations can be handled in several ways. Two possibilities:

- [d1] At *vol\_open()* time, all the vnodes on a fileset can be dis-associated.
- [d2] At individual *vol\_xxx()* operations, the vnodes affected by that call can be dis-associated.

Likewise, the re-associations can be handled in several ways. Two possibilities:

- [a1] At *vol\_close()* time, re-associate any vnodes for the fileset in question. Vnodes that could not be re-associated are marked such that any use of them (checked for at individual *vn\_xxx()* operations) fails.
- [a2] Before using any vnode (say, as input to individual *vn\_xxx()* operations), check to see if that vnode requires re-association. Again, if the re-association does not succeed, the operation fails.

The difference between [d1] and [d2] appear to be minor. One difference between [a1] and [a2] is whether it is possible for [a1] to encounter a vnode which cannot be re-associated until some later time (a subsequent *vn\_xxx()*). This should never occur (either it can be re-associated at *vol\_close()* time or it will never be able to be).

Usually, vnodes will be unused at the time of such destructive operations. Hence, they can be cleanly discarded at *vol\_open()* time or at re-association time if a problem is encountered. However, there are two exceptions.

1. Vnodes for files being accessed locally can be in use at *vol\_open()* time. These should not cause a problem since DFS will not delete or move such a locally mounted fileset. (However, it would be wise for an LFS to take the proper precautions and not rely on DFS to protect it.)
2. Normally, the DFS file exporter locks out conflicting fileset operations while it is holding a vnode (the vnod--to-volume synchronization mechanism). In one situation, however, (in the read-write path waiting for an RPC data pipe) it retains a held vnode and temporarily allows fileset operations (of any type) to occur. If a destructive fileset operation does manage to interrupt it, there will be a vnode which
  - cannot be discarded at *vol\_open()* time (since it is held)
  - cannot be re-associated at *vol\_close()* time (due to its deletion)
  - needs to be detected when the exporter comes back and tries to use it again

Therefore, even if the re-associations are performed as in [a1], the checks in [a2] are still required.

Either way, attempts to use a vnode which could not be re-associated should fail with [ESTALE].

Additional care might be required if the implementation is such that *vol\_xxx()* operations themselves instantiate or re-associate vnodes. Although DFS normally scans through the vnodes on a fileset from beginning to end, this need not be so. Another round of dis-association might be required under approach [d1] above if a *vol\_scan()*, *vol\_seek()* or *vol\_create()* are observed to back up to a previously processed object or vnode.

### 15.11.3 Complications in Dis-sociations and Re-associations

A related complication needs to be kept in mind as well. A number of *vol\_xxx()* operations delete files. For example, *vol\_delete()*, *vol\_deplete()* and some of the cloning ones. Any dis-associated vnodes which result (in other words, vnodes which, for whatever reasons, could not be outright released or freed) must be treated carefully if a new file with the same inode index and a different generation number (uniquifier) is then created. These dis-associated, *stale* vnodes must not be inadvertently used by a lookup operation working from an anode index (only) in a directory entry. These lookups are only interested in the one, *true* object that might exist on the fileset. The permissible operations on these *stale* vnodes are:

*vn\_inactive()* when their usage count finally goes to 0.

*vol\_create()* if it attempts to re-create the object with the same fileset ID, anode index and generation number. (Depending upon the implementation, *vol\_reclone()* might encounter this situation as well.)

Because dis-associated vnodes hold both an anode index and generation number (see above), determining their validity against the actual contents of the fileset is not difficult. The key point is that there can be several vnodes in existence for a given fileset ID and anode index (but with different generation numbers). At most one of these is valid; the remainder are *stale* ... waiting for *vn\_inactive()* to finally retire them.

The cause of this is situation that the fileset operations aren't bound by many rules. They can delete objects that are still being used (presumably, locally mounted). They might (based on the implementation) give rise to multiple vnodes for the same anode index and fileset. While a fileset is open (blocking ordinary accesses), they can violate the invariant that a given **fid** resolves to one-and-only-one object (the same one every time). What is important is that when the fileset is finally closed, no inconsistencies are observed to remain (other than the fact that a file locally open has been deleted and perhaps replaced).

Another complexity to be aware of is that calls to *vn\_inactive()* are not synchronized with fileset operations (this is discussed in the vnode operations chapter, *VFS (Vnode) Interface and Operations*, Chapter 16 on page 415). Therefore, this operation can come in while a fileset is open and being operated upon (by *vol\_xxx()* operations).

**Note:** This will change in the future.

### 15.11.4 Fileset Moves

A brief summary of how the DFS client handles moved filesets follows:

- Fileset movement between machines is only cleanly handled for requests that come in via the protocol exporter (from a remote cache manager). As alluded to previously in Section 15.11 on page 336, the DFS **fts** command refuses to delete, zap or move (non-locally) a fileset that is mounted locally.

When a remote file request arrives, the exporter first determines if the desired fileset (identified by a volume ID) is in fact local and available. This determination is made by a routine (*volreg\_Lookup()*, see Section 15.15 on page 396) that consults a table, maintained by DFS, of locally attached and exported filesets. If the desired fileset is available and local, it is *held* while the operation proceeds. While the fileset is held this way, any conflicting fileset operations are delayed.

If it turns out that the fileset is currently in use by a conflicting fileset operation, this request is stalled until it is able to proceed.

If, alternatively, the fileset is not found locally, a [VOLERR\_PERS\_DELETED] (derived from a preliminary [ENODEV]) error is returned to the remote client. That client, on receipt of this error, attempts to re-bind to a correct location (from an updated FLDB) before retrying its request. (This occurs in the cache manager routine *cm\_Analyze()*.)

- Fileset motion within a given machine must be handled entirely within the LFS (The above DFS mechanism is not triggered since the fileset is still local.). This mechanism also allows local, non-DFS access to be unaffected by local moves.

The key point is that any LFS implementation needs to carefully consider the VOL\_OP\_XXX bits (defined in Section 13.3.18 on page 286) in the type argument passed to *vol\_open()* and decide, on a case-by-case basis, which warrant vnode dis-association and which might require cache purging, and so forth. Each VOL\_OP\_XXX bit warns of an impending *vol\_XXX()* operation whose potential effect on file system state requires analysis.

DFS guarantees, with a few exceptions, that while a fileset is open for one of these operations, conflicting vnode operations are not allowed to come in. See *Synchronization Between Vnode and Fileset Operations* Section 16.5 on page 422, in the **vnode** operations chapter titled *VFS (Vnode) Interface and Operation*, Chapter 16 on page 415 of this document. As discussed there, DFS presently does allow some vnode operations to proceed. The LFS has the responsibility for these of dealing with in-progress volume operations and dis-associated vnodes.

## 15.12 Private LFS Fileset Data

As each fileset is exported to DFS, file system independent code makes calls to the *vol\_setdystat()* and *vol\_attach()* fileset operations. If an LFS requires that some private data or a handle be made available to it on subsequent volume operations, it can arrange for this as follows.

- Allocate space during the *vol\_setdystat()* operation and place a pointer to this in the *.v\_fsDatap* fileset **struct volume** structure field.
- In subsequent calls, this pointer can be retrieved from the supplied fileset (volume) structure.

Filesets are *detached* via calls to *vol\_freedystat()* and *vol\_detach()*. Any storage allocated by *vol\_setdystat()* must be released in *vol\_freedystat()*.

## 15.13 Fileset Operations

File system independent code above the LFS enforces the requirement that these operations can only be issued by ROOT (the local super user). For this reason, the credentials argument to many of these calls can be ignored except as indicated in the individual descriptions in Section 15.14 on page 340, *Fileset Array Functions*.

The descriptions do not list all possible error cases and return values. The primary intent has been to identify the DFS specific behavior and situations that would not arise in an ordinary UNIX (UFS) file system implementation.

With the exception of *vol\_open()*, *vol\_setdystat()*, *vol\_hold()* and *vol\_rele()*, all of these fileset operations have been preceded by a corresponding call to *vol\_open()*.

Many of these operations accept a position or anode (inode) index argument. These are intended to ONLY select ordinary files and directories. ACLs and property lists are not visible as distinct objects via any of the standard LFS interfaces (fileset, vnode, aggregate operations). If they are in fact implemented as objects on the disk, they must not appear in the "inode space" of visible filesets.

Except where explicitly noted in the following descriptions, these fileset operations do *not* have the side affect of updating the accessed, modified or changed times (atime, mtime, ctime) of a file.

### 15.14 Fileset Array Functions

Descriptions of the members of the **struct volumeops** function array follow. There are four classes of fileset operations in this array; per-fileset, whole-fileset, per-file, and vnode lookup. The function array definition in Section 15.1 on page 325 delineate them as comments in the structure itself. Refer to it for information about which class the operations belong to.

**NAME**

vol\_hold — Increment the reference count of a fileset descriptor

**SYNOPSIS**

```
int vol_hold(  
    /* IN */ struct volume *vold  
);
```

**ARGUMENTS**

*vold*                   fileset descriptor pointer.

**DESCRIPTION**

Hold a volume structure by incrementing the reference count of the **vold** fileset descriptor.

**DISCUSSION**

This operation should be performed under the volume lock as follows:

```
lock_ObtainWrite(&volp->v_lock);  
volp->v_refCount++;  
lock_ReleaseWrite(&volp->v_lock);
```

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**NAME**

vol\_rele — Decrement the reference count of a fileset descriptor

**SYNOPSIS**

```
int vol_rele(  
    /* IN */ struct volume *vold  
);
```

**ARGUMENTS**

*vold*                   fileset descriptor pointer.

**DESCRIPTION**

Release a volume structure by decrementing the reference count of the **vold** fileset descriptor. Wake up anyone that might be waiting on the structure.

**DISCUSSION**

This operation should be performed under the volume lock as follows.

```
{  
lock_ObtainWrite(&volp->v_lock);  
if (volp->v_count == 1)  
    vol_VolInactive(volp);  
else  
    volp->v_refCount--;  
    lock_ReleaseWrite(&volp->v_lock);  
}
```

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.



**NAME**

vol\_lock — Lock access to a fileset

**SYNOPSIS**

```
int vol_lock(
    /* IN */ struct volume *vold,
    /* IN */ int           type
);
```

**ARGUMENTS**

*vold*                   Fileset descriptor pointer.

*type*                   type of lock to be obtained. Allowable values are **READ\_LOCK**, **WRITE\_LOCK**, and **SHARED\_LOCK**.

**DESCRIPTION**

Lock access to a fileset.

**DISCUSSION**

Although this function is exported to user-space **DCE DFS**, there are no users of it.

This operation should utilize the standard DFS lock primitives as follows:

```
if (type == READ_LOCK)
    lock_ObtainRead(&volp->v_lock);
else if (type == WRITE_LOCK)
    lock_ObtainWrite(&volp->v_lock);
else if (type == SHARED_LOCK)
    lock_ObtainShared(&volp->v_lock);
```

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL]               The value of *type* is not **READ\_LOCK**, **WRITE\_LOCK** or **SHARED\_LOCK**.

**NAME**

vol\_unlock — Unlock access to a fileset

**SYNOPSIS**

```
int vol_unlock(  
    /* IN */ struct volume *vold,  
    /* IN */ int          type  
);
```

**ARGUMENTS**

*vold*                   Fileset descriptor pointer.

*type*                   type of lock to be released. Allowable values are **READ\_LOCK**, **WRITE\_LOCK**, and **SHARED\_LOCK**.

**DESCRIPTION**

Unlock access to a fileset.

**DISCUSSION**

Although this function is exported to user-space **DCE DFS**, there are no users of it.

This operation should utilize the standard DFS lock primitives as follows:

```
if (type == READ_LOCK)  
    lock_ReleaseRead(&volp->v_lock);  
else if (type == WRITE_LOCK)  
    lock_ReleaseWrite(&volp->v_lock);  
else if (type == SHARED_LOCK)  
    lock_ReleaseShared(&volp->v_lock);
```

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL]               The value of *type* is not **READ\_LOCK**, **WRITE\_LOCK** or **SHARED\_LOCK**.

**NAME**

vol\_open — Begin a sequence of user space administrative operations in the given fileset

**SYNOPSIS**

```
int vol_open(
    /* IN */ struct volume      *vold,
    /* IN */ long               opentype,
    /* IN */ long               errtoreturn,
    /* OUT */ struct vol_handle *handle
);
```

**ARGUMENTS**

<i>vold</i>	Fileset descriptor pointer.
<i>opentype</i>	mask of operations to be carried out in this sequence containing a set of VOL_OP_XXX bits that indicate the types of operations that will be performed on the open fileset.
<i>errtoreturn</i>	This argument is unused. The error code to return for operations that attempt to access this fileset between this <i>vol_open()</i> and the associated <i>vol_close()</i> .
<i>handle</i>	pointer to the fileset iterator structure created for the given fileset. It should be filled in as per the DISCUSSION section.

**DESCRIPTION**

This call is used to begin a sequence of operations on the fileset described by *vold*. The *opentype* parameter declares the kinds of fileset operations to be carried out. The *errtoreturn* parameter is used to set the error code to return should other operations attempt to access this fileset between this *vol\_open()* call and the bracketing *vol\_close()*.

The value of the input *opentype* parameter is constructed as the bitwise inclusive-OR of the VOL\_OP\_XXXX definitions from the include file **volume.h**, such as VOL\_OP\_SCAN. This value is useful as a hint that allows the file system to prepare a fileset to perform the operations. If such preparation is expensive for some operations, that preparation need not be carried out unless the sequence in question will include those operations.

This is the same operation bit-mask that is passed to the *vol\_concurr()* operation, which is responsible for selecting the kinds of vnode operations that may be carried out concurrently with the given sequence of operations.

At most one sequence of fileset operations may be in progress at any time. Calls to *vol\_open()* that are made while another sequence is in progress will fail, being given as their error code the value passed as *errtoreturn* when beginning the sequence in progress.

A sequence of fileset operations is terminated by a call to *vol\_close()*.

**DISCUSSION**

More details pertaining to this operation can be found in Appendix I on page 405.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

Appendix I on page 405.

**NAME**

vol\_seek — Seek to a particular index within a fileset

**SYNOPSIS**

```
int vol_seek(
    /* IN */   struct volume    *vold,
    /* IN */   long             position,
    /* INOUT */ struct vol_handle *handle
);
```

**ARGUMENTS**

*vold* Fileset descriptor pointer.

*position* canonical anode index within a fileset.

*handle* pointer to the fileset iterator structure created for the given fileset to be filled in by this call.

**DESCRIPTION**

This call is used to seek to a desired *position* within the file table of the fileset associated with *vold*. The *handle* iterator structure is updated to reflect the new position.

**DISCUSSION**

The canonical inode number supplied in the *position* argument is converted to an LFS relative one and stored in the **vol\_handle** *.index* field. As a general rule, this indicates the next file that will be operated upon.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL] The supplied anode index is an illegal value (It is less than VOL\_ROOTINO.).

[VOL\_ERR\_EOF] There are no valid objects on the fileset with an inode index greater than or equal to the supplied position argument.

[VOL\_ERR\_EOW] Returned by the UFS volume operations when the position argument is greater than the size of the statically configured inode table.

**SEE ALSO**

*vol\_open()*.

**NAME**

**vol\_tell** — Return the current index of the selected file within a fileset

**SYNOPSIS**

```
int vol_tell(  
    /* IN */ struct volume    *vold,  
    /* IN */ struct vol_handle *handle,  
    /* OUT */ long            *position  
);
```

**ARGUMENTS**

*vold* Fileset descriptor pointer.

*handle* pointer to the fileset iterator structure created for the given fileset holding current position.

*position* on return, set to the current anode index contained within *handle*.

**DESCRIPTION**

Given an iteration structure *handle* for the fileset described by *vold*, set *position* to the current index pointed to by *handle*.

**DISCUSSION**

This LFS relative anode index stored in the **struct vol\_handle** field (by *vol\_seek()*, *vol\_scan()* or *vol\_create()*) is returned in the *position* argument. As discussed in *vol\_seek()*, this is an LFS-specific (non canonical) inode index.

**Note:** Note the discrepancy with the behavior of *vol\_seek()*.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

*vol\_seek()* for discrepancy in behavior with this (*vol\_tell()*) operation.

**NAME**

vol\_scan — Read the file table for the given fileset at the current position, verifying that there is a file at the position

**SYNOPSIS**

```
int vol_scan(
    /* IN */ struct volume *vold,
    /* IN */ long position,
    /* IN */ struct vol_handle *handle
);
```

**ARGUMENTS**

*vold* Fileset descriptor pointer.

*position* canonical anode index within a fileset.

*handle* pointer to the fileset iterator structure associated with the given fileset described by *vold* to be filled in by this call.

**DESCRIPTION**

Set file descriptor fields in *handle* to describe the current file for the fileset described by *vold*, as specified by the associated iterator *handle*. If the position within *handle* is a legal one (it corresponds to an actual file within the fileset), then information about the file is written in *handle*. If the position within *handle* is *not* legal, then *vol\_scan()* returns [VOL\_ERR\_EOF] (if the position is past the end of the fileset's file table) or [VOL\_ERR\_DELETED] (if the position happens to be empty).

**DISCUSSION**

This operation performs the same functions carried out by *vol\_seek()* and, additionally, returns in the supplied **vol\_handle** structure some information on the file selected by position.

The **vol\_handle** .*type* field is set to the file's mode bits (including the type -- S\_IFMT). The *.fid.Vnode* and *.fid.Unique* fields in the **vol\_handle** are filled in (for *.Vnode*, from the canonical anode index in the position argument).

**Note:** The significance of the *vol\_handle.fid.Vnode* field isn't actually mandated by DFS. This *vol\_handle.fid* field will subsequently be passed back to *vol\_xxx()* operations requiring a *fid* argument. The *vol\_scan()* and *vol\_create()* calls can use any encoding they wish for the *.Vnode* field (either a canonical/DFS or private/LFS index) as long as the other *vol\_xxx()* operations accept the same encoding in *fids* passed to them.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[VOL\_ERR\_EOF] The supplied position argument is past the end of the fileset's file table.

[VOL\_ERR\_DELETED] The supplied position argument is empty. The *vol\_SEEK()* error values take precedence.

**SEE ALSO**

*vol\_seek()* for additional error values returned.

**NAME**

vol\_close — Finish an iteration for the given fileset

**SYNOPSIS**

```
int vol_close(  
    /* IN */ struct volume *vold,  
    /* IN */ struct vol_handle *handle,  
    /* IN */ long isabort  
);
```

**ARGUMENTS**

*vold* Fileset descriptor pointer.

*handle* pointer to the fileset iterator structure associated with the given fileset described by *vold*.

*isabort* If non-zero, this call is being made in order to abort an open fileset request (The *vol\_open()* call has already occurred).

**DESCRIPTION**

Indicate to the system that file-level operations on the fileset described by *vold* have completed. *handle* is updated so that it no longer refers to the fileset to be closed. Upon a successful *vol\_close()*, the given fileset is put back on line.

**DISCUSSION**

If this is an ordinary close (*isabort* equals 0) and the **vol\_handle** *.voltype* field indicates that the VOL\_DELONSALVAGE flag was set at *vol\_open()* time, this flag is cleared in both *volp->v\_stat\_st.states* and in the on-disk fileset storage.

If *isabort* is non-zero, the fileset may be left marked in an inconsistent state (VOL\_DELONSALVAGE) after it is closed.

See the discussion under *vol\_open()*.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

See the DISCUSSION in *vol\_open()*.



**NAME**

vol\_destroy — Delete an empty fileset

**SYNOPSIS**

```
int vol_destroy(  
    /* IN */ struct volume *vold  
);
```

**ARGUMENTS**

*vold* Fileset descriptor pointer.

**DESCRIPTION**

This operation deletes a fileset that is empty. It contains no files or objects.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EIO] An error was encountered (the fileset wasn't empty) and the operation did not complete successfully.

**NAME**

vol\_deplete — Destroy a fileset

**SYNOPSIS**

```
int vol_deplete(
    /* IN */ struct volume *vold
);
```

**ARGUMENTS**

*vold* Fileset descriptor pointer.

**DESCRIPTION**

Destroy the fileset described by *vold*, including all of its files. All in-memory vnodes associated with this fileset are also deleted, and the *Fileset Registry*

In the interest of preventing system calls from taking an uncomfortable length of time, this operation may return before the fileset is completely destroyed. If this has occurred, the operation should return code [ELOOP], and the caller should simply call it again, until it stops returning [ELOOP].

**DISCUSSION**

This operation deletes all files on the specified fileset. Note that files are unconditionally deleted and do not undergo the normal Zero-Link-Count processing discussed in Section 15.6 on page 332. The fileset itself is not deleted: that task is performed by *vol\_delete()*. See the discussion under *Looping Operations* in Section 15.10 on page 335.

When this operation is first called to begin deleting a fileset, the VOL\_DELONSALVAGE flag should be set in both *volp->v\_stat\_st.states* and in the on-disk fileset storage.

While DFS guarantees that this call will not be made on a fileset which is backing for another fileset, this fileset might currently be backed itself. Therefore, using the notation described in Section 15.3 on page 328 ,if F1 is the fileset being deleted:

```
F1 ==> Fxxx    can happen
while
Fxxx ==> F1    will never occur
```

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[ELOOP] The operation has not finished. The *vol\_deplete()* call should be made again.

[EIO] An error was encountered and the operation did not complete successfully.

**SEE ALSO**

Section 15.3 on page 328, *Fileset Clone Algorithms*, discusses backing filesets. Also, see the discussion under *Looping Operations* in Section 15.10 on page 335.

**NAME**

vol\_attach — Initialize a fileset for use by DFS as the aggregate is attached

**SYNOPSIS**

```
int vol_attach(  
    /* IN */ struct volume    *vold  
);
```

**ARGUMENTS**

*vold*                      Fileset descriptor pointer.

**DESCRIPTION**

Perform filesystem-specific initialization associated with bringing a fileset on-line while attaching an aggregate.

**DISCUSSION**

As each fileset is exported to DFS, a call is made to *vol\_setdystat()* and *vol\_attach()* - in that order. This occurs at aggregate attach time (following calls to *ag\_attach()* and *ag\_vollInfo()*) and at fileset creation time (during the call to *ag\_volCreate()* -- consult its description).

Any LFS actions necessary before a fileset can be exported should be performed here. Such operations should be performed while a write lock is held on **volp->v\_lock**.

This call succeeds if the VOL\_DELONSALVAGE flag is set, even though the fileset may be in a corrupt state.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

*ag\_volCreate()* and *ag\_vollInfo()*.

**NAME**

`vol_detach` — Make a fileset unavailable for use by DFS as the aggregate is detached

**SYNOPSIS**

```
int vol_detach(  
    /* IN */ struct volume *vold,  
    /* IN */ int localrefs  
);
```

**ARGUMENTS**

*vold* Fileset descriptor pointer.  
*localrefs* Count of references held by the caller.

**DESCRIPTION**

Perform filesystem-specific takedown associated with taking a fileset off-line while detaching an aggregate or destroying a fileset. Fails if the fileset descriptor is not idle, where “idle” is defined as having a reference count equal to the count of long-term references for the structure plus *localrefs*, the count of fileset references held directly in the caller.

**DISCUSSION**

As a fileset is made unavailable for use by DFS, calls are made to *vol\_detach()* and *vol\_freedystat()* - in that order.

Any EBUSY checks and other LFS-specific activity should be performed while a write lock is held on *volp->v\_lock*.

A fileset can be exported (attached) to DFS and (or) locally mounted. Presumably, an LFS will defer any low-level “disconnect” until a fileset is neither locally mounted nor exported. If there are any vnodes that are in-use by DFS, this call should fail with EBUSY. (The *Vnode to LFS Association* section, Section 15.11 on page 336, mentions one case in which the file exporter allows volume operations such as this to come in while it has a held vnode.) Of course, it may not be possible to differentiate vnodes being used by DFS from those being used locally.

If a fileset is not mounted locally, any in-use vnodes are due to DFS. If the fileset is mounted locally, it isn’t crucial that DFS-held vnodes be detected here.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[*error\_status\_ok*] This function returns success.

**ERRORS**

[EBUSY] An error was encountered (the fileset descriptor wasn’t idle) and the operation did not complete successfully.

**SEE ALSO**

Section 15.11 on page 336.

**NAME**

vol\_getstatus — Read specified fileset status block and return status

**SYNOPSIS**

```
int vol_getstatus(  
    /* IN */ struct volume      *volp,  
    /* OUT */ struct vol_status *statusp  
);
```

**ARGUMENTS**

*volp* Fileset descriptor pointer.  
*statusp* pointer to the fileset status structure to be filled in by this call.

**DESCRIPTION**

Place the current status of the fileset described by *volp* into the *statusp* fileset status block.

**DISCUSSION**

See Appendix J on page 409 for details concerning the processing this operation does.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EIO] An error was encountered and the operation did not complete successfully.

**SEE ALSO**

Appendix J on page 409.

**NAME**

`vol_setstatus` — Write (set) specified fileset status block

**SYNOPSIS**

```
int vol_setstatus(  
    /* IN */    struct volume    *volp,  
    /* IN */    long                mask,  
    /* IN */    struct vol_status  *statusp  
);
```

**ARGUMENTS**

*volp*                   Fileset descriptor pointer.  
*mask*                   A mask consisting of a set of VOL\_STAT\_XXX bits.  
*statusp*                pointer to the supplied fileset status structure.

**DESCRIPTION**

Set the fileset status block for the fileset described by *volp* to the contents of *statusp*, modifying only the status information indicated by bits set in *mask*.

**DISCUSSION**

See Appendix K on page 411 for details concerning the processing this operation does.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL]               A supplied status parameter is illegal or cannot be handled. For instance, if an attempt is being made to set both the VOL\_RW and VOL\_READONLY status flags.

[EIO]                   An error was encountered and the operation did not complete successfully.

**SEE ALSO**

Appendix K on page 411.

**NAME**

vol\_create — Create a file (or directory, symlink, and so forth) at a given position within a fileset

**SYNOPSIS**

```
int vol_create(
    /* IN */ struct volume      *volp,
    /* IN */ long              position,
    /* IN */ struct xvfs_attr  *xvattrp,
    /* INOUT*/ struct vol_handle *handlep,
    /* IN */ struct ucred      *credp
);
```

**ARGUMENTS**

<i>volp</i>	Fileset descriptor pointer.
<i>position</i>	A canonical anode index representing the position at which to create a file.
<i>xvattrp</i>	The vnode attribute structure giving the attributes to be applied to the new file.
<i>handlep</i>	A pointer to the fileset iterator structure associated with the fileset described by <i>volp</i> .
<i>credp</i>	A pointer to the credentials structure.

**DESCRIPTION**

Create and initialize a file at the given *position* within the fileset described by *volp*. The initial type, mode, and other attributes of the new file are given by *vattrp*. Anode descriptor information in *handle* is set to describe the new file. If the file at the given *position* exists, *vol\_create()* returns [EEXIST] and does not reinitialize it.

**DISCUSSION**

See Appendix L on page 413 for a more detailed description of the processing done by this operation.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL]	A supplied position parameter is illegal or cannot be handled.
[EEXIST]	An error was encountered and the operation did not complete successfully. Specifically, the file already exists at the specified position.

**SEE ALSO**

Appendix L on page 413.

## NAME

vol\_read — Read data from the given file

## SYNOPSIS

```
int vol_read(
    /* IN */   struct volume      *volp,
    /* IN */   struct afsFid      *Fidp,
    /* IN */   struct afsHyper    offset,
    /* IN */   long                length,
    /* IN */   char                *readbuff,
    /* IN */   struct ucred        *credp,
    /* OUT */  long                *amt_read
);
```

## ARGUMENTS

<i>volp</i>	Fileset descriptor pointer.
<i>Fidp</i>	An <b>afsFid</b> identifying the file ID of the file to be read.
<i>offset</i>	Byte position in the file at which the read is to commence.
<i>length</i>	The number of bytes to be read.
<i>readbuf</i>	A pointer to the buffer in which the file's data are to be deposited.
<i>credp</i>	A pointer to the credentials structure.
<i>amt_read</i>	Upon return, set to the number of bytes actually read from the specified file.

## DESCRIPTION

Given a file described by *Fidp* within the fileset associated with *volp*, attempt to read *length* bytes, starting at byte position *offset*. The file data are placed into the buffer pointed to by *readbuff*, which must be at least *length* bytes long. Upon completion, *amt\_read* is set to the number of bytes actually placed into *readbuff*. Other than the standard return code indicating success, and [EINVAL], *vol\_read()* may return [EIO] if a lower-level read error occurred.

## DISCUSSION

This call is used to read the contents of a symbolic link. It is the ONLY *vol\_xxx()* operation that can do this.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

## ERRORS

[EINVAL]	An attempt was made to read a FIFO, Block special device, CHAR special device or a directory.
[EIO]	An actual read error occurred.



**NAME**

vol\_write — Write data to the given file

**SYNOPSIS**

```
int vol_write(
    /* IN */ struct volume      *volp,
    /* IN */ struct afsFid     *Fidp,
    /* IN */ struct afsHyper   offset,
    /* IN */ long              length,
    /* IN */ char               *writebuf,
    /* IN */ struct ucred       *credp
);
```

**ARGUMENTS**

<i>volp</i>	Fileset descriptor pointer.
<i>Fidp</i>	An <b>afsFid</b> identifying the file ID of the file to be written.
<i>offset</i>	Byte position in the file at which the write is to commence.
<i>length</i>	The number of bytes to be written.
<i>writebuf</i>	A pointer to the buffer holding the data to be written to the file.
<i>credp</i>	A pointer to the credentials structure.

**DESCRIPTION**

Given a file described by *Fidp* within the fileset associated with *volp*, attempt to write *length* bytes, starting at byte position *offset*. The file data are taken from the buffer pointed to by *writebuf*, which must be at least *length* bytes long. Other than the standard return code indicating success, and [EINVAL], *vol\_srite()* may return [EIO] if a lower-level write error occurred.

**DISCUSSION**

This call is used to write the contents of a symbolic link. It is the ONLY *vol\_xxx()* operation that can do this.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL]	An attempt was made to write a FIFO, Block special device, CHAR special device or a directory.
[EINVAL]	An attempt was made to write to a file offset that is larger than the maximum supported file size. This can occur if the <i>length</i> argument is larger than 4 Gigabytes and the filesystem is limited to 32 bits.
[EIO]	An actual read error occurred.

**NAME**

vol\_truncate — Truncate a file to a specified size

**SYNOPSIS**

```
int vol_truncate(  
    /* IN */    struct volume    *volp,  
    /* IN */    struct afsFid    *Fidp,  
    /* IN */    struct afsHyper    newsize,  
    /* IN */    struct ucred    *credp  
);
```

**ARGUMENTS**

<i>volp</i>	Fileset descriptor pointer.
<i>Fidp</i>	An <b>afsFid</b> identifying the file ID of the file to be truncated.
<i>newsize</i>	New size of <i>Fidp</i> 's file, in bytes.
<i>credp</i>	A pointer to the credentials structure.

**DESCRIPTION**

In the fileset described by *vold*, truncate the file associated with *Fidp* to *newsize* bytes.

**DISCUSSION**

The specified file is truncated to the given size. A file can have its size increased by this call. If so, the newly acquired file space is "sparse" and does not show up as allocated space (see the quota discussion in Section 15.7 on page 333) in the fileset.

If a directory is being truncated, the *newsize* argument **MUST** be zero. In this case, the directory is left in an "empty" state without entries for "." and ".

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL]	A <i>newsize</i> of something other than 0 is specified for a directory.
[EINVAL]	A <i>newsize</i> argument is supplied that is larger than the maximum supported file size. This can occur if the <i>length</i> argument is larger than 4 Gigabytes and the filesystem is limited to 32 bits.

**SEE ALSO**

Section 15.7 on page 333.

**NAME**

vol\_delete — Delete a file at a given position within a fileset

**SYNOPSIS**

```
int vol_delete(
    /* IN */ struct volume      *volp,
    /* IN */ struct afsFid     *Fidp,
    /* IN */ struct ucred      *credp
);
```

**ARGUMENTS**

*volp* Fileset descriptor pointer.

*Fidp* An **afsFid** identifying the file ID of the file to be truncated.

*credp* A pointer to the credentials structure.

**DESCRIPTION**

Delete the file described by *Fidp* in the fileset associated with *volp*.

**DISCUSSION**

The specified file is unconditionally deleted. It does not undergo the normal Zero-Link-Count processing. The deletion is immediate. Any vnodes for the file will be in a dis-associated state (consult *Vnode to LFS Association*, Section 15.11 on page 336, and the discussion in *Zero Link Count Files*, Section 15.6 on page 332).

If a directory is being truncated, the newsize argument **MUST** be zero. In this case, the directory is left in an "empty" state without entries for "." and ".

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

Section 15.11 on page 336, *Vnode to LFS Association*, and also *Zero Link Count Files*, Section 15.6 on page 332.

## NAME

vol\_getattr — Get a file's attributes

## SYNOPSIS

```
int vol_getattr(
    /* IN */    struct volume      *volp,
    /* IN */    struct afsFid      *Fidp,
    /* OUT */   struct xvfs_attr   *xvattrp,
    /* IN */    struct ucred       *credp
);
```

## ARGUMENTS

*volp*               Fileset descriptor pointer.

*Fidp*                An **afsFid** identifying the file ID of the file.

*xvattrp*            Pointer to a returned vnode structure in which to record the file's attributes.

*credp*               A pointer to the credentials structure.

## DESCRIPTION

Read the attributes of the file indicated by *Fidp* in fileset *volp*, into the attribute structure *xvattrp*.

## DISCUSSION

Except as noted below, the standard vnode attributes at *xvattrp->vattp* are returned in an obvious manner.

*.va\_fsid*            The value returned here should be chosen with care since it needs to be unique (per fileset). Note that a given aggregate (device) can hold several filesets. Although a individual LFS is free to use any appropriate algorithm to construct this field, the following is done by Episode:

```
.va_fsid = (volumeID.low    << 16) |
            (device major # << 8)  |
            (device minor #)
```

*.va\_blocks*        The blocks used value is based on the logical amount of disk space used by the file, as if there were no copy-on-write (clone) sharing taking place. This corresponds to the visible (vs actual allocated) quota on a fileset.

The extended attributes at *xvattrp->xvattp* are returned as described in the *VFS (Vnode) Interface and Operations*, Chapter 16 on page 415. The *credp* argument is used in the computation of the *.callerAccess* field.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

## ERRORS

None.

## SEE ALSO

*VFS (Vnode) Interface and Operations*, Chapter 16 on page 415.

**NAME**

vol\_setattr — Set a file's attributes

**SYNOPSIS**

```
int vol_setattr(
    /* IN */ struct volume      *volp,
    /* IN */ struct afsFid     *Fidp,
    /* IN */ struct xvfs_attr  *xvattrp,
    /* IN */ struct ucred      *credp
);
```

**ARGUMENTS**

<i>volp</i>	Fileset descriptor pointer.
<i>Fidp</i>	An <b>afsFid</b> identifying the file ID of the file.
<i>xvattrp</i>	Pointer to a vnode structure from which to set the file's attributes.
<i>credp</i>	A pointer to the credentials structure.

**DESCRIPTION**

Set the attributes of the file indicated by *Fidp* in the fileset *volp*, according to the vnode attribute structure *xvattrp*. In this structure, a -1 in a field indicates that no new value is to be set for the corresponding attribute; this is the same convention as is used in *VOP\_SETATTR*.

**DISCUSSION**

The following file attributes can be set by this call.

- From the standard vnode attributes at *xvattrp->vattr*:

<i>.va_ctime</i>	If the <i>.tv_sec</i> sub-field is not equal to -1.
<i>.va_nlink</i>	If not equal to -1. If the link count is being set to zero, the file will be deleted by <i>vn_inactive()</i> when its vnode reference count goes to zero. As described in the topic, <i>Zero Link Count Files</i> , Section 15.6 on page 332, such files are preserved across reboots and FSCks until <i>vn_inactive()</i> explicitly deletes them.
<i>.va_mode</i>	If not equal to -1.
<i>.va_uid</i>	If not equal to -1.
<i>.va_gid</i>	If not equal to -1.
<i>.va_size</i>	If not equal to -1. If a directory is being operated on, an EISDIR is returned.
<i>.va_mtime</i>	If the <i>.tv_sec</i> sub-field equals 0 and the <i>.tv_usec</i> sub-field equals -1, the file's atime and mtime are both set to the current time. (System V style <i>utime()</i> call.) Alternatively, if the <i>.tv_sec</i> field is not equal to -1, the file's mtime is set.
<i>.va_atime</i>	If the <i>.tv_sec</i> sub-field is not equal to -1 and the above <i>.va_mtime</i> case does not supercede this one, the file's atime is set.

- From the extended vnode attributes at *xvattrp->xvattr*:

<i>.volVersion</i>	If the <i>.high</i> or <i>.low</i> sub-field is not equal to -1.
<i>.dataVersion</i>	If the <i>.high</i> or <i>.low</i> sub-field is not equal to -1.
<i>.fileI</i>	The anode generation value, from the <i>.low</i> sub-field, can be changed if it is not equal to -1. The <i>.high</i> field, the anode index, is ignored.

*.clientOnlyAttrs* If not equal to -1.

The file's ctime is NOT advanced by this call (except for an explicit setting).

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EISDIR] An attempt is made to change the length of a directory via this call.

[EROFS] The fileset is read-only.

**SEE ALSO**

Section 15.6 on page 332.

**NAME**vol\_getacl — Get a file's **ACL****SYNOPSIS**

```
int vol_getacl(
    /* IN */ struct volume      *volp,
    /* IN */ struct afsFid     *Fidp,
    /* OUT */ struct dfs_acl   *aclp,
    /* IN */ long              which,
    /* IN */ struct ucred      *credp
);
```

**ARGUMENTS**

<i>volp</i>	Fileset descriptor pointer.
<i>Fidp</i>	An <b>afsFid</b> identifying the file ID of the file.
<i>aclp</i>	Pointer to a structure into which the <b>ACL</b> and its length are to be stored.
<i>which</i>	A value indicating the type of <b>ACL</b> that is being read.
<i>credp</i>	A pointer to the credentials structure.

**DESCRIPTION**

Read the **ACL** of the file indicated by *Fidp* in fileset *volp*, into the **ACL** structure *aclp*. For regular files this is unambiguous, but directories may have up to three different **ACLs**, and *which* indicates which one is to be read: the directory's own **ACL**, the default **ACL** for subdirectories, or the default **ACL** for subfiles.

**DISCUSSION**

The LFS **ACL** specification chapter should be consulted for information about the values discussed here. It is Chapter 12 on page 251. There is further information on DFS **ACLs** in Chapter 8 on page 155.

The **ACL** designated by the *which* argument is encoded as described in the LFS **ACL** specification chapter, Chapter 12 on page 251, and returned. The *which* argument can take the following values.

**VNX\_ACL\_REGULAR\_ACL**

The **ACL** for the object itself is returned.

**VNX\_ACL\_DEFAULT\_ACL**

The initial directory **ACL** for a directory object is returned. If the object in question is not a directory, a zero-length **ACL** is returned instead.

**VNX\_ACL\_INITIAL\_ACL**

The initial file **ACL** for a the object in question is not a directory. A zero-length **ACL** is returned instead.

If the requested **ACL** does not exist, a zero length **dfs\_acl** is returned (*aclp->dfs\_acl\_len*).

**Note:** If *which* is **VNX\_ACL\_REGULAR\_ACL**, the internal Episode implementation of **ACLs** makes itself apparent here.

Normally, **ACL** operations under Episode behave as follows. When one sets an **ACL**, the file's mode bits are changed to agree with the required **ACL** entries. When a file's mode bits are changed (**chmod()**), its **ACL** is not physically changed on disk. Instead, the mode bits are "merged" into the **ACL** dynamically during **ACL** checks and when **ACLs** are returned to user space.

The *vol\_getacl()* operation fetches the **ACL** as it resides on disk, without coercing the **ACL**'s required entries to agree with the file's mode bits (which might be different due to a **chmod()** operation that was performed after the **ACL** was initially applied).

Consult the **DFS ACL** information in Chapter 8 on page 155 under the topic, *Interaction of Filesystem ACLs with UNIX Permission Bits*, and also Appendix A on page 221, *Mapping DFS ACLs to UNIX mode bits*. See Chapter 9 on page 165 for information about the structures used by **DFS** for **ACLs** in memory.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL] The *which* argument does not equal one of the values mentioned in **DISCUSSION**.

**SEE ALSO**

Chapter 12 on page 251, Chapter 8 on page 155 and Appendix A on page 221.



**NAME**vol\_setacl — Set a file's **ACL****SYNOPSIS**

```
int vol_setacl(
    /* IN */ struct volume      *volp,
    /* IN */ struct afsFid     *Fidp,
    /* IN */ struct dfs_acl    *aclp,
    /* IN */ long              index,
    /* IN */ long              which,
    /* IN */ struct ucred      *credp
);
```

**ARGUMENTS**

<i>volp</i>	Fileset descriptor pointer.
<i>Fidp</i>	An <b>afsFid</b> identifying the file ID of the file.
<i>aclp</i>	Pointer to a structure from which the ACL is to be applied ( or used).
<i>index</i>	Index of another file from which to copy the file's new ACL. This argument is presently UNUSED by this operation. The argument <i>aclp</i> is always used.
<i>which</i>	A value indicating the type of ACL that is to be set.
<i>credp</i>	A pointer to the credentials structure.

**DESCRIPTION**

Set the **ACL** of the file indicated by *Fidp* in the fileset *volp*. If *aclp* is non-null, the new **ACL** is read from it. If *aclp* is null, the new **ACL** is copied from the **ACL** associated with the file at *index*. Parameter *which* indicates which of the file's **ACLs** is to be set; for details see the description of *vol\_getacl()*.

**DISCUSSION**

As described under the *vol\_getacl()* operation, the *which* argument indicates the type of **ACL** on the specified object that is to be set. The mode bits on the target object are NOT changed to agree with the **ACL** being applied: see the note below.

**Note:** If *which* equals `VNX_ACL_REGULAR_ACL`, the internal Episode implementation of **ACLs** makes itself apparent here.

Normally, **ACL** operations under Episode behave as follows. When one sets an **ACL**, the file's mode bits are changed to agree with the required **ACL** entries. When a file's mode bits are changed (**chmod()**), its **ACL** is not physically changed on disk. Instead, the mode bits are "merged" into the **ACL** dynamically during **ACL** checks and when **ACLs** are returned to user space.

The *vol\_setacl()* operation applies the **ACL** to the file without changing the file's mode bits (which might be different due to a **chmod()** operation that was performed). Later, when the **ACL** is used, the "current" mode bits *will* actually be used.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EIO] The *which* argument equals `VNX_ACL_DEFAULT_ACL` or `VNX_ACL_INITIAL_ACL` and the specified object is not a directory. This is

returned by Episode only.

[EINVAL] The *which* argument equals VNX\_ACL\_DEFAULT\_ACL or VNX\_ACL\_INITIAL\_ACL and the specified object is not a directory.

[EINVAL] The *which* argument does not equal one of the values mentioned in DISCUSSION.

[EINVAL] The supplied **ACL** is missing a required field or is in an inconsistent or illegal state. Specifically, if:

- The **ACL**'s length is inconsistent with its contents.
- The *.mgr\_type\_field* is incorrect.
- The required entries are not all present.
- The *entry\_type\_foreign\_user* and *entry\_type\_foreign\_group* do not contain a realm UUID that differs from the default one.
- The *user\_obj* entry does not grant *perm\_control* rights.

**SEE ALSO**

*vol\_getacl()*. There is also a set of references there that are pertinent to this operation as well.

**NAME**

vol\_clone — Create a cloned image of the files on a fileset in a different fileset

**SYNOPSIS**

```
int vol_clone(
    /* IN */ struct volume *destvolp,
    /* IN */ struct volume *srcvolp,
    /* IN */ struct ucred *credp
);
```

**ARGUMENTS**

*destvolp* Fileset descriptor pointer for the destination fileset.

*srcvolp* Fileset descriptor pointer for the source fileset.

*credp* A pointer to the credentials structure.

**DESCRIPTION**

Make progress in making the entire fileset described by *destvolp* be a set of copy-on-write clones of all objects in the fileset described by *srcvolp*. The destination fileset is presumed to be empty before this call is first made.

In the interest of preventing system calls from taking an uncomfortable length of time, this operation may return before the fileset is completely processed. If this has occurred, the operation should return code [ELOOP], and the caller should simply call it again, until it stops returning [ELOOP].

**DISCUSSION**

Refer to the earlier topics on *Looping Operation Considerations*, Section 15.10 on page 335, and *Fileset Clone Algorithms*, Section 15.3 on page 328.

When the clone operation is first started, the target fileset specified by *destvolp* will be completely empty. When the clone operation has successfully completed, this target fileset will contain a "snapshot" of the files that reside on the source fileset identified by *srcvolp*.

The *.volversion* (VV) on the target (*destvolp*) fileset is not updated by this operation. The quota usage fields *.visQuotaUsage* and *.allocUsage* on the target (*destvolp*) fileset are updated as appropriate.

Both fileset header and individual file status on the source (*srcvolp*) fileset are unaffected by this call. When a clone operation is first started, the VOL\_DELONSALVAGE flag should be set in both the target's on-disk fileset storage and in the *vol1p->v\_stat\_st.states*. When the clone has been successfully completed, the VOL\_DELONSALVAGE flag should be cleared again. If the clone fails part way through, the fileset is left in a potentially inconsistent state with the VOL\_DELONSALVAGE flag set.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[ELOOP] The clone operation has not yet completed. This call should be made again.

[EDQUOT] The quota on the target *destvolp* has been exceeded. (This error should not occur if the quotas on the destination fileset have been properly established.)

[ENOSPC] The space on the aggregate has been exhausted.

**SEE ALSO**

Refer to the topics on *Looping Operation Considerations*, Section 15.10 on page 335, and *Fileset Clone Algorithms*, Section 15.3 on page 328 for more information.

**NAME**

vol\_reclone — Update a cloned image of the files on a fileset

**SYNOPSIS**

```
int vol_reclone(
    /* IN */ struct volume *srcvolp,
    /* IN */ struct volume *destvolp,
    /* IN */ struct ucred *credp
);
```

**ARGUMENTS**

*destvolp* Fileset descriptor pointer for the destination fileset.

*srcvolp* Fileset descriptor pointer for the source fileset.

*credp* A pointer to the credentials structure.

**DESCRIPTION**

This call performs actions on each of a pair of filesets that are presumed to be in the copy-on-write relationship established by *vol\_clone()*. The *destvolp* parameter should refer to the fileset in which copy-on-write pointers were established by *vol\_clone()* (as its *destvolp* parameter), and the *srcvolp* parameter should refer to the fileset to which the copy-on-write pointers were made by *vol\_clone()* (as its *srcvolp* parameter).

All changes made to objects in the copy-on-write fileset (*destvolp*) are propagated to corresponding objects in the backing fileset (*srcvolp*). Any conflicting data that had been in the backing fileset is truncated and (or) deleted. The copy-on-write fileset is made to contain no data itself, but to be filled with copy-on-write pointers to the objects in the backing fileset.

In the interest of preventing system calls from taking an uncomfortable length of time, this operation may return before the fileset is completely processed. If this has occurred, the operation should return code [ELOOP], and the caller should simply call it again, until it stops returning [ELOOP].

**DISCUSSION**

**Note:** Note that the first two arguments have their order swapped from that in the *vol\_clone()* operation.

The fileset identified by the *srcvolp* argument must be directly backed by the fileset identified by argument *destvolp* ( F1 ==> F2 ).

This operation renews the clone relationship so that the backing fileset (*destvolp*) is a snapshot of the source (*srcvolp*) as it exists currently. Note that since the clone was originally created, the source fileset can have been modified with files being either modified or deleted.

If the filesets in question are named F1 and F2 (from *srcvolp* and *destvolp*, respectively), there can in general be other filesets in the backing hierarchy which are logically unaffected by this operation. For example, filesets Fxxx and Fyyy in the following relationship are not affected:

```
Fxxx ==> F1 ==> F2 ==> Fyyy
```

**Note:** Although DFS will never normally generate the above scenario, it could arise in error situations.

The *.volversion* (VV) on the backing (*destvolp*) fileset is not updated by this operation. The quota usage fields *.visQuotaUsage* and *.allocUsage* on the backing (*destvolp*) fileset are updated as appropriate.

Both fileset and individual file status on the source (*srcvolp*) fileset are unaffected by this call.

When a reclone is first started, the VOL\_DELONSALVAGE flag is set in both the on-disk fileset storage and in the *destvolp->v\_stat\_st.states*. When the reclone has been successfully completed, the VOL\_DELONSALVAGE flag is cleared again. If the reclone fails part way through, the fileset is left in a potentially inconsistent state with the VOL\_DELONSALVAGE flag set.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[ELOOP] The clone operation has not yet completed. This call should be made again.

[EDQUOT] The quota on the target *destvolp* has been exceeded. (This error should not occur if the quotas on the destination fileset have been properly established.)

[ENOSPC] The space on the aggregate has been exhausted.

[EINVAL] A file on a fileset indicated by *srcvolp* (F1) is backed by a different fileset than the one indicated by the *destvolp* (F2) argument.

[EFBIG] A file on the backing fileset required deleting and cannot be deleted.

**SEE ALSO**

Refer to the topics on *Looping Operation Considerations*, Section 15.10 on page 335, and *Fileset Clone Algorithms*, Section 15.3 on page 328 for more information.

**NAME**

vol\_unclone — Prepare to destroy a cloned image of the files on a fileset

**SYNOPSIS**

```
int vol_unclone(
    /* IN */ struct volume *srcvolp,
    /* IN */ struct volume *destvolp,
    /* IN */ struct ucred *credp
);
```

**ARGUMENTS**

*srcvolp* Fileset descriptor pointer for the source fileset.

*destvolp* Fileset descriptor pointer for the destination fileset.

*credp* A pointer to the credentials structure.

**DESCRIPTION**

This call performs actions on each of a pair of filesets that are presumed to be in the copy-on-write relationship established by *vol\_clone()*. The *destvolp* parameter should refer to the fileset in which copy-on-write pointers were established by *vol\_clone()* (as its *destvolp* parameter). The *srcvolp* parameter should refer to the fileset to which the copy-on-write pointers were made by *vol\_clone()* (as its *srcvolp* parameter).

This call breaks the copy-on-write relationship between the two filesets. All copy-on-write data pointers in the copy-on-write fileset (*destvolp*) are replaced by ordinary data pointers, which are moved there from the backing fileset (*srcvolp*). The pointers are deleted from the backing fileset, thus truncating the objects it contains. After this operation completes successfully, the contents of the former backing fileset (*srcvolp*) are not necessarily meaningful; the former backing fileset is usually destroyed.

In the interest of preventing system calls from taking an uncomfortable length of time, this operation may return before the fileset is completely processed. If this has occurred, the operation should return code [ELOOP], and the caller should simply call it again, until it stops returning [ELOOP].

**DISCUSSION**

**Note:** Note that the first two arguments have their order swapped from that in the *vol\_clone()* operation.

The fileset identified by *srcvolp* must be directly backed by the fileset identified by *destvolp* ( F1 ==> F2 ).

This operation is made prior to deleting the backing fileset (*destvolp*) in order to "break" the copy-on-write relationship between files on it, files on the source fileset (*srcvolp*) and on any fileset that it is backing. Since the fileset is left in an inconsistent state (VOL\_DELONSALVAGE) and is about to be deleted, the internal state of it and its files is undefined at completion. The important point is that the copy-on-write relationship with upstream and downstream filesets be broken and their quotas be adjusted. It is unspecified as to whether the no longer needed raw disk space is given back to the aggregate at *vol\_unclone()* or at *avol\_delete()* time.

This operation should never fail due to a disk full condition. Specifically, it should NOT require additional disk space beyond what is already in use.

In theory, the quota on an upstream or downstream fileset could be exceeded as the copy-on-write relationship is broken. In this case, the unclone should proceed, exceeding (although updating) whatever quotas it encounters.

**Note:** Another possibility would be for the unclone to fail with [EDQUOTE].

With the exception of copy-on-write related fields and quotas, this call does not affect the fileset status of upstream and downstream filesets (specifically: *.llFwdID*, and *.llBackId*).

After the unclone operation completes successfully, the backing fileset will actually be deleted via *vol\_destroy()*.

Consider a fileset backing relationship as follows:

```
F1 ==> F2 ==> F3
```

Prior to deleting F2, DFS will call *vol\_unclone(F1, F2, ..)* to break the copy-on-write relationships involving F2, resulting in the following:

```
F1 ==> F3          F2
```

If F3 were being deleted instead, DFS would call *vol\_unclone(F2, F3, ..)* to break the copy-on-write relationship between F2 and F3, resulting in:

```
F1 ==> F2          F3 .
```

When an unclone is first started, the VOL\_DELONSALVAGE flag is set in both the on-disk fileset storage and in the *destvolp->v\_stat\_st.states*. It remains in this state until the fileset is deleted.

#### RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

#### ERRORS

[ELOOP] The clone operation has not yet completed. This call should be made again.

[EINVAL] A file on a fileset indicated by *srcvolp* (F1) is backed by a different fileset than the one indicated by the *destvolp* (F2) argument.

[EFBIG] A file on the backing fileset required deleting and cannot be deleted.

#### SEE ALSO

Refer to the topics on *Looping Operation Considerations*, Section 15.10 on page 335, and *Fileset Clone Algorithms*, Section 15.3 on page 328 for more information.



**NAME**

vol\_vget — Return the vnode for a designated afsFid

**SYNOPSIS**

```
int vol_vget(
    /* IN */ struct volume *volp,
    /* IN */ struct afsFid *fidp,
    /* OUT */ struct vnode *vnoddep
);
```

**ARGUMENTS**

*volp* Fileset descriptor pointer for a fileset.

*fidp* Fileset descriptor (**afsFid**) for the desired file.

*vnoddep* A pointer to the vnode corresponding to the given **afsFid**, upon completion.

**DESCRIPTION**

Given a fileset descriptor pointer and a file ID, generate the associated vnode structure.

**DISCUSSION**

The returned vnode should be held (its *.v\_count* incremented).

This call should return a vnode for a file even if its link count is zero. This is not an error. Consult the topic on *Zero Link Count Files*, Section 15.6 on page 332.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[ESTALE] The *.Unique* generation number in the **afsFid** does not agree with what is found in the selected vnode. Also, this code is returned on any errors that might be attributable to a bad *.vnode* or *.unique afsFid* field.

**SEE ALSO**

Refer to the topic on *Zero Link Count Files*, Section 15.6 on page 332.

**NAME**

vol\_root — Return the vnode for the root directory on a particular fileset

**SYNOPSIS**

```
int vol_root(  
    /* IN */    struct volume    *volp,  
    /* OUT */   struct vnode     *vnodep  
);
```

**ARGUMENTS**

*volp*                 Fileset descriptor pointer for a fileset.  
*vnodep*               A pointer to the vnode corresponding to the given *volp*, upon completion.

**DESCRIPTION**

Given a fileset descriptor pointer, generate the vnode structure for the root of the fileset.

**Note:** This fileset operation is not used and may be obsolete.

**DISCUSSION**

The returned vnode should be held (its *.v\_count* incremented).

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

None.

**NAME**

vol\_isroot — Determine if a supplied **afsFid** represents the root directory on a fileset

**SYNOPSIS**

```
int vol_isroot(  
    /* IN */    struct volume    *volp,  
    /* IN */    struct afsFid     *fidp,  
    /* OUT */   long             *flagp  
);
```

**ARGUMENTS**

<i>volp</i>	Fileset descriptor pointer for a fileset.
<i>fidp</i>	Fileset descriptor ( <b>afsFid</b> ) for some file.
<i>flagp</i>	A pointer to the returned (boolean) flag, upon completion.

**DESCRIPTION**

Set the flag pointed to by *flagp* to 1 if the *Fid* presents the root of the fileset, or to 0 if it represents some non-root file.

**DISCUSSION**

Only the *.Vnode* and *.Unique* fields of the **afsFid** are examined.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

None.

**NAME**

vol\_getvv — Get Fileset Version number

**SYNOPSIS**

```
int vol_getvv(  
    /* IN */ struct volume *volp,  
    /* OUT */ struct afsHyper *vnp  
);
```

**ARGUMENTS**

*volp* Fileset descriptor pointer for a fileset.  
*vnp* The returned fileset version number.

**DESCRIPTION**

Put the fileset version number of the volume represented by *volp* in the location pointed to by *vnp*. (This is a special case of *vol\_getstatus()*.)

**DISCUSSION**

None.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**NAME**

vol\_setdystat — Allocate private per-fileset filesystem-specific data area (if needed)

**SYNOPSIS**

```
int vol_setdystat(
    /* INOUT */ struct volume      *volp,
    /* IN */   struct vol_stat_dy  *dystatp
);
```

**ARGUMENTS**

*volp*                   Fileset descriptor pointer for a fileset.

*dystatp*                Pointer to dynamic area of fileset status block that has already been established by a call to *ag\_volInfo()* or within the *ag\_volCreate()* call.

**DESCRIPTION**

Allocate storage for file-system-specific data about a fileset that is to be attached, and put a pointer to that storage in the *Fileset Registry* entry that is being prepared. For **DCE LFS** filesets, the file-specific data includes the index of the fileset within the aggregate, the highest valid file index within the fileset, and an **DCE LFS handle** token representing the open fileset). The fileset index field is initialized at this time from the corresponding field in *dystatp*, but the other fields are initialized during *vol\_attach()*.

**DISCUSSION**

If private LFS per-fileset data is required, it should be allocated and a pointer to it stored in the fileset structure (**struct volume**) *.v\_fsDatap* field.

**Note:** Higher level DFS functions do not examine this field, so any desired storage allocation mechanism can be used.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

Refer to the topic, *Private LFS Fileset Data* in Section 15.12 on page 339.

## NAME

vol\_freedystat — Free private per-fileset filesystem-specific data area

## SYNOPSIS

```
int vol_freedystat(  
    /* INOUT */ struct volume    *volp  
);
```

## ARGUMENTS

*volp*                    Fileset descriptor pointer for a fileset.

## DESCRIPTION

Free the file-specific data area that was allocated by *vol\_setdystat*.

## DISCUSSION

Any storage allocated by *vol\_setdystat()* should be freed. The **struct volume** *.v\_fsDatap* field should be cleared as well.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

## ERRORS

None.

## SEE ALSO

*vol\_setdystat()*.

**NAME**

vol\_setnewvid — Set new Fileset's fileset ID

**SYNOPSIS**

```
int vol_setnewvid(
    /* INOUT */ struct volume    *volp,
    /* IN */    struct afsHyper  *idp
);
```

**ARGUMENTS**

*volp*                   Fileset descriptor pointer for a fileset.  
*idp*                    The returned Fileset ID.

**DESCRIPTION**

Set the ID of the fileset represented by *volp* to the ID in *idp*. This is called twice to carry out an exchange of fileset identities, generally after a clone fileset is created. For **DCE LFS**, it modifies entries in the mount table as well as modifying the fileset data structures.

**DISCUSSION**

The fileset ID of the specified fileset is changed to the specified value. Only LFS internal and on-disk structures should be updated by this call; specifically, the fileset ID in the **struct volume** structure should not be changed.

**Note:** This call is obsolete; *vol\_swapids()* appears to be in use these days. It should be implemented, however, since DFS could choose to use it again in the future (to swap the IDs of filesets on different aggregates).

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

## NAME

vol\_copyacl — Copy an ACL from one file to another

## SYNOPSIS

```
int vol_copyacl(
    /* IN */ struct volume *volp,
    /* IN */ struct afsFid *Fidp,
    /* IN */ long destw,
    /* IN */ long index,
    /* IN */ long srcw,
    /* IN */ struct ucred *credp
);
```

## ARGUMENTS

<i>volp</i>	Fileset descriptor pointer for a fileset.
<i>fidp</i>	Fileset descriptor ( <b>afsFid</b> ) identifying the destination file to which the <b>ACL</b> should be applied.
<i>destw</i>	A value indicating the type of <b>ACL</b> on the destination that is to be set.
<i>index</i>	A canonical anode index, as used by <i>vol_seek()</i> or <i>vol_scan()</i> , that identifies the source file from which the <b>ACL</b> is to be copied.
<i>srcw</i>	A value indicating the type of <b>ACL</b> on the source file that is to be copied (from).
<i>credp</i>	A pointer to a credentials structure.

## DESCRIPTION

Set an **ACL** of the file indicated by *Fidp* in the fileset *volp* by copying it from an **ACL** of the file indicated by *index*. The **ACL** copying may be performed by sharing a reference to an **ACL**, if the underlying file system supports such shared references, but there is no semantic requirement for such sharing. Parameter *destw* indicates which of the file's **ACLs** is to be set. Parameter *srcw* indicates which of the source file's **ACLs** is to be copied. For details, see the description of *vol\_getacl()*.

## DISCUSSION

The *srcw* argument selects an **ACL** on the source object, identified by the *index* argument:

VNX_ACL_REGULAR_ACL	<b>ACL</b> on the object itself.
VNX_ACL_DEFAULT_ACL	Initial directory <b>ACL</b> on a directory.
VNX_ACL_INITIAL_ACL	Initial file <b>ACL</b> on a directory.

This **ACL** is copied to the destination object, identified by the *Fidp* argument. The type of **ACL** to be set is indicated by the *destw* argument which has the same VNX\_ACL\_XXX values listed above.

It is legal for the source **ACL** to not exist; in this case, the destination **ACL**, if any, is deleted and the destination is left without an **ACL** of the specified type.

Except as noted further below, the *srcw* and *destw* arguments can each take any of the legal values. For example, it is legal to copy an *initial directory ACL* of a directory to the *regular ACL* of a file.

**Note:** If *srcw* or *destw* equal VNX\_ACL\_REGULAR\_ACL, the internal DFS implementation of **ACLs** makes itself apparent here. The normal, vnode **ACL** operations under Episode behave as follows:



When one sets an **ACL**, the file's mode bits are changed to agree with the required entries. When a file's mode bits are changed (*chmod()*), its **ACL** is not physically changed on disk. Instead, the mode bits are merged with the **ACL** dynamically during **ACL** checks and when **ACLs** are returned to user space.

As is the case for *vol\_getacl()* and *vol\_setacl()*, this operation entirely ignores the mode bits on the file (either merging in if *srcw* is `VNX_ACL_REGULAR_ACL` or setting if *destw* is `VNX_ACL_REGULAR_ACL`).

#### RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

#### ERRORS

[VOL\_ERR\_EOF] There are no valid objects on the fileset with an anode index greater than or equal to the supplied source index.

[EINVAL] The supplied source *index* is an illegal value. It is a too-small or a negative value other than `VOL_ROOTINO`.

[VOL\_ERR\_DELETED]

No valid file has been selected by the position argument (*index*). The preceding errors, [VOL\_ERR\_DELETED] and [EINVAL] take precedence.

[EINVAL] If either the *srcw* or *destw* arguments are not one of the legal `VNX_ACL_XXX` values; or the *srcw* or *destw* argument is either `VNX_ACL_INITIAL_ACL` or `VNX_ACL_DEFAULT_ACL` and the destination object is not a directory; or a supplied **ACL** is inconsistent (see *vol\_setacl()*).

#### SEE ALSO

*vol\_getacl()*, *vol\_setacl()*. Consult the **DFS ACL** information in Chapter 8 on page 155 under the topic, *Interaction of Filesystem ACLs with UNIX Permission Bits*, and also Appendix A on page 221, *Mapping DFS ACLs to UNIX Mode Bits*. See **Chapter 9 on page 165 for information about the structures used by DFS for ACLs in memory.**

**NAME**

`vol_concurr` — Determine allowable concurrency on a fileset (being opened)

**SYNOPSIS**

```
int vol_concurr(  
    /* IN */    struct volume    *volp,  
    /* IN */    long            type,  
    /* IN */    long            errorType,  
    /* OUT */   char            *concurr  
);
```

**ARGUMENTS**

<i>volp</i>	Fileset descriptor pointer.
<i>type</i>	A mask consisting of VOL_OP_XXX bits that indicate the activities that will be performed on the open fileset.
<i>errorType</i>	This argument is unused.
<i>concurr</i>	Location to store descriptor for the allowable kinds of calls that can take place concurrently with the set of planned operations. They are VOL_CONCUR_ALLOPS, VOL_CONCUR_READONLY, and VOL_CONCUR_NOOPS.

**DESCRIPTION**

This function is given a bit mask giving the set of fileset operations to be carried out on the given fileset. It returns an abbreviated description of the class of vnode operations that may be executed concurrently with that set of fileset operations.

The value of the input *type* parameter is constructed as the bitwise inclusive-OR of the VOL\_OP\_XXX definitions from the include file `<volume.h>`, such as VOL\_OP\_SCAN.

**DISCUSSION**

This operation examines the *type* argument and returns a VOL\_CONCURR\_XXX value according to the same rules that are listed under *vol\_open()*.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

*vol\_open()*.

**NAME**

vol\_swapids — Swap fileset identifiers of two filesets

**SYNOPSIS**

```
int vol_swapids(
    /* IN */ struct volume *vol1p,
    /* IN */ struct volume *vol2p,
    /* IN */ struct ucred *credp
);
```

**ARGUMENTS**

*vol1p* Fileset descriptor pointer for a fileset.  
*vol2p* Fileset descriptor pointer for a another fileset.  
*credp* A pointer to the credentials structure.

**DESCRIPTION**

This operation exchanges the fileset IDs of the two filesets pointed to. Both filesets are of the same type.

**DISCUSSION**

This operation is usually called following a clone operation in order to swap the identities of two filesets. In addition to updating any relevant LFS-specific data structures, the volume IDs (as found in the fileset (**struct volume**) structure *.vol\_stat\_st.volld* field) should be swapped in on-disk storage for the two filesets. The fileset structure itself should not be changed; that will be taken care of by the caller of this operation.

The two filesets in question will usually reside on the same aggregate (by virtue of being clones).

**Note:** Currently, Episode DOES allow the filesets to reside on different aggregates. It is the discretion of the LFS to either restrict the two filesets to the same aggregate or not.

Assume that *vol1p* and *vol2p*, prior to this call, referred to filesets as follows:

*vol1p* fileset index F1 on some aggregate, where say, the fileset ID = X

*vol2p* fileset index F2 on some aggregate, where say, the fileset ID = Y

After this call completes, fileset F1 will have an ID of Y and fileset F2 will have an ID of X. An active vnode which, prior to this call, referred to an object on fileset F1 will, after this call, refer to an object on fileset F2.

The contents of the *vol1p* and *vol2p* volume structures, with the exception of the fileset ID fields, will be swapped by the DFS layer once this call returns.

**Note:** If, in the future, filesets are allowed to reside on different aggregates, atomicity concerns arise unless the aggregates can be updated atomically. In this case, the aggregate holding *vol1p* should be updated before the one holding *vol2p*. This rule permits higher level recovery code to operate and take advantage of this ordering.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**NAME**

vol\_sync — Sync the given fileset

**SYNOPSIS**

```
int vol_sync(  
    /* IN */ struct volume *vold,  
    /* IN */ int guarantee  
);
```

**ARGUMENTS**

*vold* Fileset descriptor pointer.

*guarantee* Describes the type of sync to be done, choose from one of: **VOL\_SYNC\_COMMITSTATUS**, **VOL\_SYNC\_COMMITMETA** or **VOL\_SYNC\_COMMITALL**.

**DESCRIPTION**

Sync the state of the fileset to the permanent storage, according to the *guarantee* parameter:

**VOL\_SYNC\_COMMITSTATUS**

The fileset's status information is written to permanent storage. The call does not return until the I/O has completed.

**VOL\_SYNC\_COMMITMETA**

All dirty meta-data is written to permanent storage. The call does not return until the I/O has completed.

**VOL\_SYNC\_COMMITALL**

All dirty data (both meta-data and user data) is written to permanent storage. The call does not return until the I/O has completed.

In all cases, this call does not return until the I/O operations have completed. The ability to write *just* the requested types of data benefits certain of the DFS fileset operations.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

**ERRORS**

[EINVAL] The **guarantee** argument is not one of **VOL\_SYNC\_COMMITSTATUS**, **VOL\_SYNC\_COMMITMETA** or **VOL\_SYNC\_COMMITALL**.

**SEE ALSO**

The definitions of the **VOL\_SYNC\_XXX** types can be found in Section 13.3.20 on page 288.

**NAME**

vol\_pushstatus — Serialize fileset status updates

**SYNOPSIS**

```
int vol_pushstatus(  
    /* IN */ struct volume *vold,  
    );
```

**ARGUMENTS**

*vold*                      Fileset descriptor pointer.

**DESCRIPTION**

Serialize subsequent changes to a fileset's status with respect to prior changes. This operation guarantees that any preceding modifications of the fileset's status will be serialized with respect to later meta-data updates.

**DISCUSSION**

When this call returns, the LFS guarantees that any operations started subsequently on the fileset will not commit to disk unless any that ended prior to this call commit as well.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

## NAME

vol\_readdir — Read entries from a directory

## SYNOPSIS

```
int vol_readdir(
    /* IN */    struct volume *vold,
    /* IN */    struct afsFid *Fidp,
    /* IN */    u_long bufSize,
    /* IN */    char *bufferp,
    /* IN */    struct ucred *credp,
    /* INOUT */ struct afsHyper *positionp,
    /* OUT */   u_long *numEntriesp
);
```

## ARGUMENTS

<i>vold</i>	Fileset descriptor pointer.
<i>Fidp</i>	File ID of the directory.
<i>bufSize</i>	Size, in bytes, of <i>bufferp</i> .
<i>bufferp</i>	Pointer to buffer of <i>bufSize</i> bytes into which the directory entries are to be deposited.
<i>credp</i>	Credential structure.
<i>positionp</i>	Position at which to start reading. This parameter should be set to zero (0) on the first call for a given directory. The value of <i>positionp</i> will be modified by each call, so that, if passed in to a subsequent call, the directory entries will be fetched starting where the last call ended.
<i>numEntriesp</i>	Set to indicate the number of directory entries that are being returned.

## DESCRIPTION

This operation returns directory entries in a system-independent representation. The buffer is filled with **vol\_dirent** structures, as defined in Section 13.3 on page 286. Each directory entry is returned on a 4 byte boundary.

This operation is meant to be called in a loop until *\*numEntriesp* is zero (0).

## DISCUSSION

Appendix H on page 403 contains the fields as they should be filled in by this function.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

## ERRORS

[EINVAL]	(This error is optional. The <b>DCE LFS</b> does not return this value.) If the supplied buffer is not aligned on a 0-modulo-4 byte boundary.
[EINVAL]	The specified object is not a directory.
[EINVAL]	A position argument is supplied that is larger than the maximum supported file size.

**Note:** The intent is to catch position arguments that are larger than 4-Gigabytes on file systems limited to 32-bits.

## NAME

vol\_appenddir — Append entries to a directory

## SYNOPSIS

```
int vol_appenddir(
    /* IN */    struct volume *vold,
    /* IN */    struct afsFid *Fidp,
    /* IN */    u_long *numEntriesp
    /* IN */    u_long bufSize,
    /* IN */    char *bufferp,
    /* IN */    int preserveOffsets,
    /* IN */    struct ucred *credp,
);
```

## ARGUMENTS

<i>vold</i>	Fileset descriptor pointer.
<i>Fidp</i>	File ID of the directory.
<i>numEntriesp</i>	Set to indicate the number of directory entries that are to be appended.
<i>bufSize</i>	Size, in bytes, of <i>bufferp</i> .
<i>bufferp</i>	Pointer to buffer of <i>bufSize</i> bytes into which the directory entries are to be appended.
<i>preserveOffsets</i>	If true, try to preserve the entries' offsets; otherwise, strictly append.
<i>credp</i>	Credential structure.

## DESCRIPTION

This operation appends entries to a directory. The buffer is expected to contain directory entries in the system-independent representation defined by **struct vol\_dirent** (see Section 13.3 on page 286).

## DISCUSSION

Consult the description under the *vol\_readdir()* operation. This function is the inverse of that operation. Directory entries in the format of the entries in **struct vol\_dirent** in the supplied buffer, *bufferp*, are added to the directory.

The **inode** link counts for files being appended are not updated by the operation of this function. If the *preserveOffsets* argument is zero, the entries can be placed anywhere within the directory. If the *preserveOffsets* argument is non-zero, the entries should be placed at the directory location indicated by the *.offset* field in the **struct vol\_dirent**. This will only be the case if the source (the *vol\_readdir()*) and the destination (the *vol\_appenddir()*) aggregates are of the same type.

Entries for "." and ".." can be added with this call since directories created by the *vol\_create()* function operation initially do not contain them.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

## ERRORS

[EINVAL] The specified object is not a directory.

**Note:** Episode currently does not perform this check.



- [EEXIST] A name being added already exists in the directory.  
**Note:** Episode currently does not perform this check.
- [E2BIG] The end of the supplied buffer is reached before the specified number of directory entries is processed.
- [EINVAL] An inconsistency is encountered in a directory entry. For example, if the *.offset* field in the **struct vol\_dirent** is not a legal encoding.
- [EINVAL] *preserveOffsets* is non-zero and the entry cannot be placed at the desired location (for instance, something else is there).

**SEE ALSO**

*vol\_readdir()*.

## NAME

vol\_bulksetstatus — Set statuses of multiple filesets atomically

## SYNOPSIS

```
int vol_bulksetstatus(
    /* IN */    unsigned int        arrayLen,
    /* IN */    struct vol_statusDesc *statusArray,
);
```

## ARGUMENTS

*arrayLen*            The number of elements in *statusArray*. This number can range from 1 through VOL\_MAX\_BULKSETSTATUS

*statusArray*        An array containing one status descriptor for each volume being updated. It describes which filesets' statuses should be updated and what their new values should be.

## DESCRIPTION

This operation atomically updates the status of multiple filesets and possibly performs a fileset ID swap. If two of the **vol\_statusDesc** structures indicate that their respective filesets should have a new ID, the IDs of the two filesets are swapped. See Section 13.3.7 on page 282 for the definition of the **vol\_statusDesc** structure and also the definition of VOL\_MAX\_BULKSETSTATUS.

## DISCUSSION

This operation atomically updates the status for a number of filesets. Either all the changes are made or none of them are. The specified filesets must all reside on the same aggregate.

Each entry in the supplied *statusArray* argument contains the information necessary to update a single fileset. The fields have the following meanings:

*volp*                pointer to a fileset structure

*vsd\_mask*           a mask (VOL\_STAT\_XXX flags) indicating the particular status fields to be set

*vsd\_status*        The volume status being set

If this call is used to *change* the fileset IDs of two filesets (VOL\_STAT\_VOLID mask bits), an identity swap is being performed and the comments found in the *vol\_swapids()* function apply as well.

**Note:**    The emphasis here is on *change*. The LFS can safely ignore attempts to set a fileset ID to its current value.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

## ERRORS

[EXDEV]            The indicated filesets do not all reside on the same aggregate.

[EINVAL]          The two filesets being updated have not had their IDs swapped. Otherwise, the aggregate would be left with two filesets having possibly the same ID. This case can be returned if any fileset ID (VOL\_STAT\_VOLID) is being *changed*, meaning that an identity swap is being performed and the fileset IDs have not been successfully swapped.

**SEE ALSO**

See the description under the *vol\_setstatus()* operation for the details of how each fileset is updated. Also, see the *vol\_swapids()* function for details on identity swap.

**NAME**

vol\_getzlc — Get a vnode from fileset zero-link-count (ZLC) list

**SYNOPSIS**

```
int vol_getzlc(  
    /* IN */      struct volume  *vold,  
    /* INOUT */  unsigned long  *iterator,  
    /* OUT */    struct vnode    **new_vnode,  
);
```

**ARGUMENTS**

<i>volp</i>	Fileset descriptor pointer.
<i>iterator</i>	Iterator used to enumerate the files on the fileset ZLC list.
<i>new_vnode</i>	Held vnode for the specified ZLC file, or NULL if there are no more files to be enumerated.

**DESCRIPTION**

Return a held vnode representing a file on the fileset's zero-linkcount (ZLC) file list. It is used at **attach** time following the *vol\_attach()* volume operation so that fileset ZLC files can be added to the DFS ZLC list for use by remote clients. The iterator, *\*iterator*, should be initialized to zero prior to the first call to this function. It will be updated by every call to this operation so this function can locate the "next" ZLC file for the next call. There are no more files to be enumerated when *new\_vnode* comes back NULL. At this point, the return value will also be zero.

**DISCUSSION**

Consult the Zero-Link-Count file discussion earlier in this document. (See Section 15.6 on page 332.)

Since this operation is performed as each fileset is attached, performance is a consideration and the time required to identify these zero-link-count files must be kept in mind. As mentioned earlier in Section 15.6 on page 332, DFS will not make this call on a readonly (*.readonly* replica or *.backup* clone) fileset.

**RETURN VALUE**

This function is always successful. If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function was successful.

**ERRORS**

None.

**SEE ALSO**

See Section 15.6 on page 332 for information on zero-link-count files.

**NAME**

vol\_getnextholes — Identify allocation holes in files.

**SYNOPSIS**

```
int vol_getzlc(
    /* IN */      struct volume      *vold,
    /* IN */      struct afsFid      *Fidp,
    /* INOUT */   struct vol_NextHole *iterp,
    /* IN */      struct ucred       *credp
);
```

**ARGUMENTS**

*vold*                 Fileset descriptor pointer.

*Fidp*                 An afsFID (file ID) of the directory.

*iterp*                Descriptor for a set of allocation holes. See DISCUSSION below.

*credp*                Credential structure pointer.

**DESCRIPTION**

Refer to Section 13.3.3 on page 277 for a description of the **struct vol\_NextHole** structure used as the INOUT parameter *iterp*.

On each call, the filesystem determines whether there are any allocation holes with byte addresses beginning at or after the value of **startPoint**. If there are any, the filesystem describes between 1 and **VOLHOLE\_MAX\_HOLES** of them in the **holes** array by giving their starting byte addresses and their length in bytes, and by placing the count of filled-in allocation hole descriptions in the **outCount** field. Further, if there are no further allocation holes beyond those that are being described in the **holes** array, the filesystem should set the **VOLHOLE\_FLAG\_LAST** bit in the **flags** field; otherwise, if hole descriptors are being returned and **VOLHOLE\_FLAG\_LAST** is *not* set, the filesystem updates the **startPoint** to point past the last allocation hole being returned. If there are no allocation holes meeting the criteria (starting at or past the given value of **startPoint**), a zero is placed in the **outCount** field and the **VOLHOLE\_FLAG\_LAST** bit is set in the **flags** field.

All allocation holes for a file may be discovered by making a sequence of calls, where for the first call, the **startPoint** field of the **struct vol\_NextHole** structure is zeroed.

**DISCUSSION**

For each hole returned, **.holes[].holeStart** gives the file offset of the hole returned, and **.holes[].holeLen** gives the length in bytes of the hole.

This function is useful for avoiding backing up, copying (during a move) and restoring regions within a file that are sparse.

This function is currently *unimplemented* and returns [ENOSYS].

**RETURN VALUE**

This function presently never succeeds as it is unimplemented. If this function were to succeed (in the future), it would return a value of zero.

**ERRORS**

[ENOSYS]             This function always returns this error as it is currently unimplemented.

## **15.15 Fileset Registry Array Functions**

Descriptions of the *Fileset Registry* Functions follow. These are exported functions that permit a client to address the registry. They consist of functions to enter, delete and look up entries (filesets) in the *Fileset Registry*.

**NAME**

volreg\_Enter — Enter a fileset and its information into the Fileset Registry

**SYNOPSIS**

```
int volreg_Enter(  
    /* IN */ afsHyper      *avolid,  
    /* IN */ struct volume *avolP,  
    /* IN */ char          *avolname  
);
```

**ARGUMENTS**

<i>avolid</i>	The Fileset ID to enter.
<i>avolP</i>	Pointer to the above fileset's descriptor.
<i>avolname</i>	Fileset's name.

**DESCRIPTION**

Given fileset and aggregate IDs, along with a pointer to the associated fileset information, create an entry in the *Fileset Registry* and insert the given information into it.

**DISCUSSION**

None.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EEXIST] The fileset was already in the *Fileset Registry* table.

**NAME**

volreg\_Delete — Delete a fileset entry from the Fileset Registry

**SYNOPSIS**

```
int volreg_Delete(  
    /* IN */ afsHyper    *avolid,  
    /* IN */ char        *avolname  
);
```

**ARGUMENTS**

*avolid*            The Fileset ID to delete.  
*avolname*         Above Fileset's name.

**DESCRIPTION**

Delete the *Fileset Registry* entry (if any) that corresponds to the given fileset and aggregate ID pair.

**DISCUSSION**

None.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.



**NAME**

volreg\_Lookup — Find the entry in the Volume Registry

**SYNOPSIS**

```
int volreg_Lookup(  
    /* IN */ struct afsFid *afidP,  
    /* OUT */ struct volume **avolPP,  
    );
```

**ARGUMENTS**

*afidP*                   The File ID to look up.  
*avolPP*                  Set to the address of the fileset descriptor associated with the given file.

**DESCRIPTION**

Given a file ID, find the entry in the Volume Registry corresponding to the fileset containing the file and return pointers to the associated **struct volume** fileset structure and the address of the file's vnode pointer.

**DISCUSSION**

None.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL]                The file ID has a bad fileset ID.  
[ENODEV]                The file ID could not be found.



## Exporting the Filesets in an Aggregate

### G.1 Fileset Export Steps

The steps involved in exporting an aggregate's filesets are the following.

- The *dfsexport* command reads the **dfstab** file and collects the appropriate aggregate parameters (device pathname, aggregate name, aggregate type, aggregate Id) into an argument structure. See Section 13.2.16 on page 273 for a definition of the **dfstab** entry structure.
- For an exported UFS volume, additional arguments which will be interpreted by the UFS-specific file system code, are collected into an auxiliary structure. See Section 13.2.17 on page 274 for a definition of the **ufs\_astab** entry structure used for UFS volumes.
 

**Note:** Other LFS implementations could, if desired, pass in additional information this way: via an argument structure that is interpreted by only the *dfsexport* command and the file system dependent code (with appropriate modifications to suit this arrangement).
- This aggregate argument structure and optionally, the file system dependent auxiliary one, is passed into the kernel to execute an aggregate attach operation.
- An *ag\_attach()* call to the file system-dependent aggregate code is made.
- This is followed by a series of *ag\_vollInfo()* aggregate operations which are used to obtain the identities of the filesets that are present. Volume calls to *vol\_setdystat()* and *vol\_attach()* are made for each fileset discovered.
- At this point, the filesets are considered exported to DFS and can be accessed at any time.
- In order to have file system syncs work properly for filesets that are DFS exported but not locally mounted, the aggregate (note: aggregate, not fileset!) is specially mounted at **/opt/dcelocal/var/dfs/aggrs/<aggregate-name>**. This is discussed in greater detail in *Aggregate Mounts*, Section 14.3 on page 305.

An LFS may permit individual filesets to be mounted for local access outside of DFS. Episode allows this either before or after they have been exported to DFS.

**Note:** This local mount mechanism is primarily an issue between the local operating system (OS) and the LFS with, hopefully, marginal impact on DFS.

Similar steps occur when an aggregate is detached from (made unavailable to) DFS. What the LFS observes is the following (omitting miscellaneous calls to *ag\_hold()* (and) *ag\_rele()* as well as to *vol\_hold()* (and) *ag\_rele()*):

- The *dfsexport -detach* command first verifies that filesets on the aggregate are idle by revoking all tokens for files on them. An optional force argument to the command specifies that tokens are to be forcibly revoked. This optional option should be used with caution, as it forces an aggregate (or partition) to be detached even if all tokens cannot be revoked. This can result in users accessing the data from the aggregate that are unable to save that data.
- A series of *ag\_vollinfo()* calls are made in order to ascertain the identities of the filesets on the aggregate.

- For each of these filesets, calls to *vol\_detach()* and *vol\_freedystat()* are made -- in that order.
- Finally, an *ag\_detach()* call occurs.

**Note:** DFS is currently re-working the mechanics of fileset export and attachment (more cleanly separating the two notions). It is not known when these will appear in the general DFS product.

## Filled Values for `vol_dirent` Fields

### H.1 Returned Values

Refer to Section 13.3 on page 286 for the definition of the `struct vol_dirent` structure.

Fields within the `vol_dirent` structure should be filled in as follows.

<code>.offset</code>	An indication of where within the directory this entry was found. The intention is that if this offset is passed back via the <code>vol_appenddir()</code> operation during a fileset move, the directory entry will go back to the "same" location. This allows in-progress (interrupted by the move) directory scans to be oblivious to the fact that the directory has moved. This functionality is only required when the <code>vol_readdir()</code> and <code>vol_appenddir()</code> operations were performed on aggregates of the same type. (That is, the encoding of this <code>.offset</code> is LFS-private.)  Additional discussion of directory offsets and their properties can be found under the topic, <i>Directory Offsets</i> , Section 16.5.1 on page 423, and also under <code>vn_readdir()</code> , Section 16.7 on page 447.
<code>.vnodeNum</code>	A canonical inode (anode) index in which <code>-1</code> refers to the root directory, <code>0</code> refers to the first inode after that, and so forth.
<code>.codesetTag</code>	For the current implementation, this should be set to <code>0</code> .
<code>.reclen</code>	The total length of this entry, rounded up to a 4 byte boundary. (Therefore it is the distance from the start of this entry to the next one in a supplied buffer.)
<code>.namelen</code>	The length of the file name field, not including the terminating null character.
<code>.name</code>	A variable sized file name, including a terminating null character.

As many entries as will fit in the supplied buffer, or as remain in the directory, are actually read.

On return, the count pointed at by the `numentriesp` argument is set to the number of actual directory entries returned. If there are no remaining directory entries, this count should be set to `0`.

The `positionp` argument, an `afsHyper`, will be set to zero to indicate that the directory should be read from its beginning. On return, it should be updated by the LFS so that a subsequent call with it will continue with the "next" entry in the directory. DFS makes no assumptions regarding the ordering of entries within a directory. There is no reason to require the first two entries returned for a directory to be `."` and `".."`.

**Note:** Episode rounds the supplied output buffer size down to a multiple of 512 bytes. The *last* directory entry returned has its `.recordlen` field incremented to cover the remaining space within the final 512 byte "block". (It is set to the next 512 byte boundary)

It is believed that this behavior exists only to allow the file exporter to produce meaningful directory offsets on platforms whose native entries do not contain such a field. This scheme is not required. There does not appear to be any reason for an LFS to mimic this 512-byte block behavior.



# Values for `vol_open`

## I.1 The Type Argument

The individual `VOL_OP_XXX` bits in the type argument indicate the types of `vol_xxx()` operations that might be performed on the fileset while it is open. (A user space utility wishing to open a fileset actually specifies a set of `VOL_SYS_XXX` values, each of which has been `#define`'d to a corresponding `VOL_OP_XXX` value.) The `VOL_SYS` definitions can be found in Section 13.3.19 on page 287. Before a subsequent `vol_xxx()` operation is executed, DFS will ensure that the fileset was opened with the appropriate `VOL_OP_XXX` flag. The following table indicates the `VOL_OP_XXX` flags that are required for each `vol_xxx()` operation. Operations that aren't mentioned perform no check (either because they can be performed on any fileset that is open or because they are invoked as part of a larger DFS operation such as attaching or opening a fileset).

<code>vol_seek</code>	<code>VOL_OP_SEEK</code>
<code>vol_tell</code>	<code>VOL_OP_TELL</code>
<code>vol_scan</code>	<code>VOL_OP_SCAN</code>
<code>vol_destroy</code>	<code>VOL_OP_DESTROY</code>
<code>vol_deplete</code>	<code>VOL_OP_DEplete</code>
<code>vol_getstatus</code>	<code>VOL_OP_GETSTATUS</code>
<code>vol_setstatus</code>	<code>VOL_OP_SETSTATUS</code>
<code>vol_create</code>	<code>VOL_OP_CREATE</code>
<code>vol_read</code>	<code>VOL_OP_READ</code>
<code>vol_write</code>	<code>VOL_OP_WRITE</code>
<code>vol_truncate</code>	<code>VOL_OP_TRUNCATE</code>
<code>vol_delete</code>	<code>VOL_OP_DELETE</code>
<code>vol_getattr</code>	<code>VOL_OP_GETATTR</code>
<code>vol_setattr</code>	<code>VOL_OP_SETATTR</code>
<code>vol_getacl</code>	<code>VOL_OP_GETACL</code>
<code>vol_setacl</code>	<code>VOL_OP_SETACL</code>
<code>vol_clone</code>	<code>VOL_OP_CLONE</code>
<code>vol_reclone</code>	<code>VOL_OP_RECLONE</code>
<code>vol_unclone</code>	<code>VOL_OP_UNCLONE</code>
<code>vol_getvv</code>	<code>VOL_OP_GETSTATUS &amp;&amp; VOL_OP_SETSTATUS</code>
<code>vol_copyacl</code>	<code>VOL_OP_COPYACL</code>
<code>vol_swapids</code>	<code>VOL_OP_SWAPIDS</code>
<code>vol_sync</code>	<code>VOL_OP_SYNC</code>
<code>vol_pushstatus</code>	<code>VOL_OP_PUSHSTATUS</code>
<code>vol_readdir</code>	<code>VOL_OP_READIR</code>
<code>vol_appenddir</code>	<code>VOL_OP_APPENDDIR</code>
<code>vol_bulksetstatus</code>	<code>VOL_OP_SETSTATUS</code>
<code>vol_getnextholes</code>	<code>VOL_OP_GETNEXTHOLES</code>

User space code uses the `VOL_OP_NOACCESS` flag as a general mechanism to lock out all other fileset (`vol_xxx()`) and vnode accesses to a fileset.

### I.1.1 Concurrency

Based upon the *type* argument, this call establishes the *volp*->*v\_stat\_st.concurrency* field which controls the extent to which vnode operations can proceed on the fileset while it is open. Higher level code in DFS guarantees that a given fileset will only be in use (*vol\_open()*'d) by one utility (such as **fts**) at a time.

If there are no restrictions on vnode operations which might proceed in parallel, the *.concurrency* field should be set to `VOL_CONCUR_ALLOPS`.

**Note:** Episode does this if only `VOL_OP_GETSTATUS` is set.

If read-only type vnode operations are to be allowed to proceed in parallel, the *.concurrency* field should be set to `VOL_CONCUR_READONLY`.

**Note:** Episode does this if only bits from the following list are specified.

<code>VOL_OP_GETSTATUS</code>	<code>VOL_OP_GETATTR</code>
<code>VOL_OP_GETACL</code>	<code>VOL_OP_SEEK</code>
<code>VOL_OP_READ</code>	<code>VOL_OP_TELL</code>
<code>VOL_OP_SCAN</code>	<code>VOL_OP_READDIR</code>
<code>VOL_OP_GETNEXTHOLES</code>	<code>VOL_OP_SETSTATUS</code>
<code>VOL_OP_SYNC</code>	<code>VOL_OP_PUSHSTATUS</code>

If *no* vnode operations are to be allowed to proceed in parallel, the *.concurrency* field should be set to `VOL_CONCUR_NOOPS`.

**Note:** Episode does this if any other `VOL_OP_XXX` bits are specified.

**WARNING: DFS does allow vnode operations initiated by the process that has the volume open to proceed!**

### I.1.2 Handling Inconsistent State

If any of the bits listed below are set in the *type* argument, the `VOL_DELONSALVAGE` flag should be set in both *volp*->*v\_stat\_st.states* and in the on-disk fileset storage. (This indicates that the volume may be left in an inconsistent state if DFS is not completely successful in the operations it is about to perform.)

<code>VOL_OP_CREATE</code>	<code>VOL_OP_WRITE</code>
<code>VOL_OP_TRUNCATE</code>	<code>VOL_OP_DELETE</code>
<code>VOL_OP_SETACL</code>	<code>VOL_OP_SETNEWVID</code>
<code>VOL_OP_APPENDDIR</code>	

**Note:** Certain volume operations, such as the cloning ones, explicitly set and clear the `VOL_DELONSALVAGE` flag.



## I.2 The Fileset Handle

Refer to Section 13.3.21 on page 288 for the definition of the **struct** `vol_handle` structure.

The caller of this operation will record the type argument into the `vol_handle` `.opentype` field.

The `.index` field of the supplied `vol_handle` argument should be initialized as follows. Note that the caller of this operation will record the `type` argument into the `vol_handle` `.opentype` field.

The assumption is made that the root directory on a fileset will have the lowest possible (legal) inode index. This field is set to one past that root index: the index that is (or would be) used by the "first" non-root object on the fileset. Note that an LFS specific index (as in the `.vnode` component of an **afsFid**) and not a canonical one is used.

**Note:** There is some uncertainty as to whether DFS really relies on the setting of the `.index` field by `vol_open()` or whether it always does a `vol_scan()` or `vol_seek()` before using it. The behavior of Episode is described here, although in the future this document may change to leave this field undefined here.

### I.2.1 Preparing for Fileset Operations

Although conflicting operations (through either the volume or vnode operations vectors) will be blocked while this fileset is held open, it might be necessary to purge certain structures or take other actions in preparation for the volume operations that are about to be performed.

Depending upon the type of open requested, either the purging of cached file and fileset state or the breaking of associations between vnodes and LFS state might be required. This topic is discussed in the topic, *Vnode to LFS association* in Section 15.11 on page 336.



## Status Returned for `vol_getstatus`

### J.1 Fileset Status Set

The following status fields are set from on-disk fileset storage.

```

statusp->vol_st.volId
statusp->vol_st.parentId
statusp->vol_st.cloneTime
statusp->vol_st.volName
statusp->vol_st.states

statusp->vol_dy.creationDate
statusp->vol_dy.updateDate
statusp->vol_dy.accessDate
statusp->vol_dy.copyDate
statusp->vol_dy.volversion
statusp->vol_dy.backupId
statusp->vol_dy.cloneId
statusp->vol_dy.llBackId
statusp->vol_dy.llFwdId
statusp->vol_dy.allocLimit
statusp->vol_dy.allocUsage
statusp->vol_dy.visQuotaLimit
statusp->vol_dy.visQuotaUsage
statusp->vol_dy.unique
statusp->vol_dy.index
statusp->vol_dy.parentIndex
statusp->vol_dy.nodeMax
statusp->vol_dy.aggrId
statusp->vol_dy.statusMsg

```

The `statusp->vol_dy.aggrId` field is set to the aggregate-ID of the aggregate holding this fileset; namely, from the aggregate structure:

```
volp->v_paggrp->a_aggrid
```

The `statusp->vol_st.states` field is treated specially. Any flags not stored on disk for the fileset (see Section 13.3.4 on page 277 and Section 13.3.5 on page 279) are instead taken from whatever value they currently hold in the `volp->v_stat.states` field. The `VOL_LCLMOUNT` flag in this states field is set (otherwise, cleared) if the fileset is mounted locally. Once the `statusp->vol_st.states` field has been computed, it is copied into the fileset **struct volume** structure as well (at `volp->v_stat_st.states`).

The `statusp->vol_dy.nodeMax` field is computed in an LFS-specific manner.

**Note:** Episode always sets `statusp->vol_dy.parentIndex` to 0.

Static status fields in `statusp->vol_st` not mentioned above should be copied from the passed in `volp->v_stat_st.vol_st` structure.

Dynamic status fields in *statusp->vol\_dy* not mentioned above should be zeroed.

If the fileset in question has not been attached, only the static status from *volp->v\_stat\_st* is copied into the returned *vol\_status*.

## Status Set for vol\_setstatus

### K.1 Status Set

The mask argument consists of a set of VOL\_STAT\_XXX flags which indicate which fields from the supplied status are to be applied to the fileset (and stored on-disk).

The following VOL\_STAT\_XXX flags are processed by this call.

#### VOL\_STAT\_VOLID

*statusp->vol\_st.vol\_st.volId* is copied to on-disk fileset storage.

#### VOL\_STAT\_PARENTID

*statusp->vol\_st.vol\_st.parentId* is copied to on-disk fileset storage and also to the volume structure (*volp->v\_stat\_st.vol\_st.parentId*).

**Note:** Why doesn't the DFS caller do this?

#### VOL\_STAT\_CLONETIME

*statusp->vol\_st.vol\_st.cloneTime* is copied to on-disk fileset storage and also to the volume structure (*volp->v\_stat\_st.vol\_st.cloneTime*).

If fewer than 64 bits of time are stored on-disk, the value deposited in the volume structure MUST agree exactly with what is placed on disk.

#### VOL\_STAT\_VVCURRTIME

*statusp->vol\_st.vvCurrentTime* is copied to on-disk fileset storage.

#### VOL\_STAT\_VVPINGCURRTIME

*statusp->vol\_st.vvPingCurrentTime* is copied to on-disk fileset storage.

#### VOL\_STAT\_TYPE

*statusp->vol\_st.type* is copied to on-disk fileset storage.

#### VOL\_STAT\_STATES

flags stored by LFS on-disk (see Section 13.3.4 on page 277 and Section 13.3.5 on page 279) are copied from *statusp->vol\_st.states* to on-disk fileset storage.

With the exception of the following, the supplied state flags are copied into the volume structure as well (*volp->v\_stat\_st.vol\_st.states*):

VOL\_BUSY  
VOL\_DEADMEAT  
VOL\_GRABWAITING  
VOL\_LOOKUPWAITING  
VOL\_DELONSALVAGE  
VOL\_OPENDONE

#### VOL\_STAT\_RECLAIMDALLY

*statusp->vol\_st.reclaimDally* is copied to on-disk fileset storage.

#### VOL\_STAT\_VOLMOVETIMEOUT

*statusp->vol\_stat.volMoveTimeout* is copied to on-disk fileset storage.

#### VOL\_STAT\_VOLNAME

*statusp->vol\_stat.volName* is copied to on-disk fileset storage and also the volume structure

(*volp*->*v\_stat\_st.vol\_st.volName*).

VOL\_STAT\_COPYDATE

*statusp*->*vol\_dy.copyDate* is copied to on-disk fileset storage.

VOL\_STAT\_VERSION

*statusp*->*vol\_dy.volversion* is copied to on-disk fileset storage.

VOL\_STAT\_BACKUPID

*statusp*->*vol\_dy.backupId* is copied to on-disk fileset storage.

VOL\_STAT\_CLONEID

*statusp*->*vol\_dy.cloneId* is copied to on-disk fileset storage.

VOL\_STAT\_LLFWDID

*statusp*->*vol\_dy.llFwdId* is copied to on-disk fileset storage.

VOL\_STAT\_LLBACKID

*statusp*->*vol\_dy.llBackId* is copied to on-disk fileset storage.

VOL\_STAT\_ALLOCLIMIT

*statusp*->*vol\_dy.allocLimit* is copied to on-disk fileset storage.

VOL\_STAT\_VISLIMIT

*statusp*->*vol\_dy.visQuotaLimit* is copied to on-disk fileset storage.

VOL\_STAT\_UNIQUE

*statusp*->*vol\_dy.unique* is copied to on-disk fileset storage.

VOL\_STAT\_STATUSMSG

*statusp*->*vol\_dy.statusMsg* is copied to on-disk fileset storage.

## Processing for *vol\_create*

### L.1 Processing Accomplished

The canonical anode number supplied in the position argument specifies the location (anode index) at which the file is to be created. The anode generation number to be used is taken from the *xvattr->xvattr.fileID.low* field.

**Note:** If the root directory of a fileset is being created, its anode generation number should be set to 1.

The initial file state is taken from the following fields in the initial **struct vattr** at *xvattr->vattr*.

```
.va_mode , .va_type      file type
.va_oid  , .va_gid      file owner
.va_rdev                for device files
```

The link count is set to 1. The atime, mtime and ctime times are all set to the current time.

The file's **uid** should be set to zeroes -- ignoring the value supplied via the *xvattr* argument.

File status fields not mentioned above should be set to 0 (or whatever other default value is appropriate for the LFS).

The newly created file has no associated **ACLs** or property lists. In particular, no initial or inheritance **ACL** is applied.

Ordinary files, directories, symlinks, FIFOs and device files can all be created via this call.

If a directory is created (the supplied *.va\_type* specifies **VDIR**), it should be initialized to an empty state WITHOUT "." and ".." entries (these will be supplied in a subsequent *vol\_appenddir()* operation).

This operation fills in several fields in the *handlep vol\_handle* as follows.

*handlep->fid*        The *.Vnode* and *.Unique* fields are set to the appropriate value for the newly created file. See the discussion under *vol\_scan()*.

*handlep->type*     Set to a combination of the initial mode (from *.va\_mode*) and type (**S\_IFMT** encoding based on the specified *.va\_type*) bits for the newly created file.

*handlep->index*    Set to the same value as the *fid.Vnode* field (LFS-specific anode index).





# VFS (Vnode) Interface and Operations

## 16.1 Overview

As DFS was being designed, there were several goals that related to the VFS interface.

1. That DFS itself be as portable as possible.
2. That DFS "plug" into an OS kernel with minimal (none, if possible) changes to that kernel.
3. If files are being accessed both through DFS as well as a local path (including, say, an NFS exporter), that the local accesses perform the proper DFS token synchronization.

To achieve the first two goals, DFS is itself standardized on a particular definition of the VFS interface which is very similar to SUN's initial design. The **vnode** operations vector for this interface is the **struct xvfs\_xops** defined in Section 13.6.4 on page 296 and is often referred to as the extended vnode interface. All components of DFS (the Cache Manager (CM), Protocol Exporter (PX) and LFS) are written to this set of vnode operations.

The mapping between the format of the native (expected by the OS) vnode operations and this extended set of operations is handled by a set of glue code. Appendix C on page 245 describes typical components of a VFS+ package. In particular, some glue code examples are presented in Section C.3 on page 246. This glue code comes in several types which are invoked in different situations.

- One type of glue code is invoked whenever the native OS performs a vnode operation on the filesystem (that is, NOT a DFS access). Its job is to obtain DFS tokens for the operation in question, call into the file system to perform the operation and then release the DFS tokens before returning. This code is referred to as "token glue".
- A second type of glue converts from the native OS vnode interface (what its caller expects) to the extended vnode interface (what the callee expects).
- A third type of glue converts from the extended vnode interface (what its caller expects) to the native OS vnode interface (what the callee expects). This code is referred to as "vfs glue".

How this glue is invoked will become clear in the following sections.

### 16.1.1 Enhanced Vnode Operations Vector

For files exported to DFS, vnodes are assumed to have a **vnode** ops vector that points at a **struct xvfs\_vnodeops** vector which contains three sets of **vnode** operation pointers (see Section 13.6 on page 293 for definition).

In general terms (the specifics should become clearer below), the three sets of vnode operation vectors contained within this **struct xvfs\_xops** (see Section 13.6.4 on page 296) are the following.

#### O-ops Native OS format

Invoked by the VOPO\_XXX or VOP\_XXX macros

The glue routines pointed at by these pointers are invoked by the native OS as part of ordinary file system activities. For the most part, they behave as shown in the following pseudo-code:

```
xglue_xxx( )
if ((the fileset is exported to DFS) &&
    (the fileset is read/write))
    obtain the needed DFS tokens

/* Perform the operation via indirecting
 * through the .nops vector */
VOPN_xxx( )

if (we obtained tokens above)
    return the tokens
```

For a typical system, the source for these routines can be found in a location such as **src/file/xvnode/<system>/xvfs\_osglue.c**.

### **X-ops Extended format**

Invoked by the VOPX\_XXX macros.

On a file server the DFS Protocol Exporter, which doesn't need to obtain DFS tokens, calls through this vector to perform operations on exported files.

If the physical file system is an LFS that directly supports the extended vnode operations (as Episode does), the entries in this vector point directly at LFS operations.

If the physical file system does not support the extended vnode operations (a native UFS), this vector points at a set of glue routines which attempt to implement the extended operations in terms of the native (or "original") ones pointed at by the **N-ops** vector. For a typical system, the source for this glue (the *xufs\_xxx()* routines) can be found in **src/file/xvnode/<system>/xvfs\_vfs<xxx>.c**.

### **N-ops Native OS format**

Invoked by the VOPN\_XXX macros.

Operations pointed at by this vector are only ever directly invoked by "glue" code in either the **O-ops** or **X-ops** routines mentioned above.

If the physical file system is an LFS that directly supports the extended **vnode** operations (as Episode does), this vector points at a set of glue routines that convert from native OS vnode operations to extended ones. For a typical system, the source for this glue (call them, say, *nux\_xxx()* routines) can be found in **src/file/xvnode/HPUX/xvfs\_os2vfs.c**.

If the physical file system does not support the extended vnode operations (as say, in a native UFS), this is the original (vendor supplied) vector that points into the filesystem code.

The manner in which these three vectors, along with the different types of glue code, are used is most easily demonstrated by a few examples.

**Example 16-1** Protocol Exporter Access to a Fully Functional LFS

The Protocol Exporter already has any tokens that are required. Therefore, it simply calls the through the **X-ops** vector:

```
/* PX routine responding to a remote DFS request */
VOPX_xxx()
/* which lands in the LFS directly */
lfs_xxx()

=>Perform the operation<=
```

**Example 16-2** Protocol Exporter Access to a non-LFS Filesystem

This would be a file system such as UFS. The **X-ops** vector, which the PS calls through as in the case above, must point at glue which "reformats" the operation into one that the native filesystem can perform:

```
/* PX routine responding to a remote DFS request */
VOPX_xxx()
/* which lands in the xufs_xxx glue */
xufs_xxx()
/* Convert to a native vnode operation and invoke */
VOPN_xxx()
/* which lands in the native f.s. */
ufs_xxx()
=>Perform the operation<=
```

**Example 16-3** Local OS to a Fully Functional LFS

(Not through the DFS protocol exporter.) The native OS, knowing nothing about DFS, calls through the **O-ops** vector which lands in the token-obtaining glue:

```
/* native OS */
VOP_xxx() /* equivalently: VOPO_xxx() */
/* which lands in the xglue_xxx() glue */
xglue_xxx()
    Obtain tokens, if necessary
    /* Invoke operation */
    VOPN_xxx()
        /* which lands in the nux_ glue */
        nux_xxx()
            /* Invoke extended vnode operation */
            VOPX_xxx()
                /* which lands in the LFS */
                lfs_xxx()
                    =>Perform the operation<=
```

**Example 16-4** Local OS Access to a Non-LFS Filesystem Such as UFS.

(Not through the DFS protocol exporter.) The native OS, knowing nothing about DFS, calls through the **O-ops** vector which lands in the token-obtaining glue:

```
/* native OS */
VOP_XXX()          /* equivalently: VOPO_XXX() */

/* which lands in the xglue_XXX() glue */
xglue_XXX()
/* Obtain tokens, if necessary */
/* Invoke operation */
VOPN_XXX()
/* which lands in the native f.s. */
ufs_XXX()
=>Perform the operation<=
```

**16.1.2** Converted Vnodes

Vnodes which contain an enhanced (3-part, see above) operations vector are said to be converted. They are identified by the fact that the `V_CONVERTED` bit is set in their flags word. The easiest way to obtain such an enhanced vnode vector is to make a call to the DFS procedure `xvfs_InitFromXOps()` at initialization time. Once a vector has been constructed, individual vnodes can be pointed at it and have their `V_CONVERTED` flag set via one of the mechanisms described below:

```
xvfs_InitFromXOps(&lfs_xops, &enh_ops_vector);
```

The first argument is a **struct xvfs\_xops** vector of extended vnode operations provided by the LFS.

The second argument is a **struct xvnodeops** that is filled in by the call as listed below:

**O-ops**    pointed at `xglue_XXX()` glue routines  
**X-ops**    pointed at supplied LFS operations (from the first argument)  
**N-ops**    pointed at `nux_XXX()` glue routines

As mentioned above, files being accessed via DFS must have their **vnodes** converted to guarantee that the proper token synchronization occurs for all LFS accesses.

**Vnodes** for the LFS can always be converted as described above. Although the `xglue_XXX()` and `nux_XXX()` glue are ALWAYS invoked for local accesses, even if DFS is not in the picture, they "pass through" relatively quickly and don't have a significant impact on performance.

**Note:** This is the approach taken by Episode.

### 16.1.3 Enhanced Vfs Vector

In a similar manner, DFS assumes that the vfs operations vector has been converted into an enhanced **struct xvfs\_vfsops** as follows.

```
struct xvfs_vfsops {
    struct vfsops      xvfsops;
    struct vfsops      vfsops;
    int (*vfsgetvolume)();
};
```

The components of this structure have the following significance:

**.xvfsops** With the few exceptions listed below, the pointers in this vector are the same as those found in the *.vfsops* vector for the LFS.

The *\_vget*, *\_root* and *\_unmount* pointers are replaced by pointers to glue code: *xglue\_vget()*, *xglue\_root()* and *xglue\_unmount()*.

The *xglue\_root()* and *xglue\_vget()* interludes both invoke the LFS operation via the *VFSX\_xxx* macro and, if necessary (the fileset is exported and writeable), converts the returned vnode as discussed above. Note that this behavior is only relevant for exported native UFS filesets; for a fully functional LFS, the assumption is that the *vnodes* it returns are already converted.

The *xglue\_unmount()* interlude flushes any held tokens for the fileset in question if necessary (the fileset is exported and writeable). Unlike the *xglue\_root()* and *xglue\_vget()* functions, this does serve a purpose for fully functional LFS filesystems.

This vector is used by the native OS and is accessed via the standard *VFS\_xxx* macros.

**.vfsops** The "original" vfs operations vector supplied by the LFS. This vector is only used by DFS and is accessed via the *VFSX\_xxx* macros.

**.vfsgetvolume**

This pointer is set to point at an LFS routine which is capable of converting between a vfs structure and a held DFS volume pointer. See the detailed description of *vfs\_getvolume()* in Section 16.6.1 on page 424.

This procedure is used by the *xglue\_unmount()*, *xglue\_root()* and *xglue\_vget()* glue routines mentioned above. (They need a fileset structure in order to determine if the fileset is exported and writeable.)

For native UFS filesystems, this pointer is set to a DFS routine which is able to "simulate" the needed functionality.

The easiest way to construct a "converted" vfs ops vector is by means of the DFS procedure *xvfs\_InitFromVFSOps()* as follows:

```
xvfs_InitFromVFSOps(&lfs_vfsops, &enh_vfsops, lfs_getvolfunc)
```

The first argument is the address of the standard vfs operations vector provided by the LFS.

The second argument is the address of an enhanced **xvfs\_vfsops** vector which is, as described above, built by the call.

The third argument is the *vfs\_getvolume()* function provided by the LFS and described in Section 16.6.1 on page 424.

There are several options for how the vfs operations vector for a fileset is "converted" as described above. A few are listed below.

1. At initialization time, if it is known that DFS is present, the vfs operations vector can be constructed as shown above and inserted into the `vfssw[]` table so that it will work its way into all vfs structures.

If it is unknown at initialization whether DFS will be installed (dynamically) later, this `vfssw[]` can be performed at a later time when DFS "announces" itself to the LFS for the first time. One possibility would be to do this when an aggregate or fileset was first attached.

To avoid having to write into the `vfssw[]` table, it can initially point at a full sized `xvfs_vfsops` structure within which only a standard `vfsops` vector at the front is present. If DFS is present, the full `xvfs_vfsops` structure can be converted "in place" by passing that structure to the first two arguments of `xvfs_InitFromVFSOps()`.

**Note:** `xvfs_InitFromVFSOps()` is carefully coded to allow this in-place conversion to occur.

Once again, this conversion can be done either at initialization or a later time. (Assuming, as appears to be the case, that the in-place conversion of the vfs ops vector is safe with regard to other processes and processors.

**Note:** Episode behaves this way, performing the in-place conversion at initialization time.

2. Schemes that involve changing, at fileset attach and detach time, the `.vfs_ops` pointer in a vfs structure are also possible. The question of what to do if the fileset isn't mounted (and hence, there is no vfs structure) at attach time must be addressed, however.

**Note:** Although DFS currently "mounts" LFS aggregates that are attached in order to make `vol_sync()` work properly, this doesn't occur until after they are attached.

## 16.2 Administrative Rights

Several of the operations listed here require that the caller have administrative rights to the file being operated upon. Under standard UNIX, these rights would be granted for ROOT as well as the owner of the file.

Within the LFS, administrative rights to a file are granted under the following situations (see *Access Control List Overview*, Chapter 8 on page 155).

1. If the caller is the local ROOT.
2. If the caller is in the special DFS administrative group.
3. If the caller has **perm\_control** rights to the file in question. If the file is not protected by an **ACL**, this right will be granted to the file's owner.

Vnode operations, both standard and extended, which have traditionally required the caller to be the owner of a file (or ROOT) should be changed to work as shown above instead.

### 16.3 Copy-on-Write Impacts

As discussed in the volume operations specification, it is expected that some sort of copy-on-write technology is used in order to make clone operations space and time efficient. Depending on the particular implementation, a number of the **vnode** operations listed below may unexpectedly fail with [ENOSPC] or errors if the copy-on-write sharing cannot be broken. Specifically:

- a. the fileset is currently cloned and the aggregate in question is full (ENOSPC) or
- b. the fileset quota limits prevent the fileset in question (which one depends on the implementation) from physically growing in size (EDQUOT)

As one might expect, any operation that creates or grows something on disk (file, directory, ACL, and so on) can fail. Additional operations that MIGHT fail as well are listed below.

- Over-writing an existing piece of a file.
- Any operations that modify the contents of a directory. For example: *unlink()*, *rename()*, *rmdir()*, and so forth.
- Truncating a file to a smaller size. The copy-on-write algorithm should be designed so that the following operations never fail (with [ENOSPC] or [EDQUOT] errors, at least).
- The *vn\_setattr()* operation. This implies that either anodes are copied or space for them is reserved during a clone operation. Although it is not a DFS requirement, it is suggested that the following operations likewise never fail (with (with [ENOSPC] or [EDQUOT] errors).
- Truncating a file to 0 size. Some implementations may only be able to do this if the file has not been modified since the time of the clone.

### 16.4 Swap Files

Some operating systems allow swapping to ordinary files in the file system. Since there's no benefit in cloning such files, it is suggested that an LFS avoid performance and unexpected error (see ENOSPC and EDQUOT discussion earlier) problems by "ignoring" such files at clone time. Although there are any number of possible approaches here (so long as the OS and LFS agree), the following is one possibility:

- At clone (or reclone) time, any swap files that are encountered are treated as zero-length files. Although the file IS cloned in order to maintain consistent link counts, the contents of the swap file are ignored and not subject to any copy-on-write behavior.
- Any attempt to enable swapping on a currently cloned file (perhaps: with greater than zero length) fails.
- Obviously, any attempt to enable swapping to a file on a read-only clone should fail (presumably, with [EROFS]).

## 16.5 Synchronization Between Vnode and Fileset Operations

For the most part, DFS protects against undesired interactions between volume (`vol_XXX`) administrative operations and ordinary (local or exporter) **vnode** operations.

- a. Both the vnode glue layer and the DFS protocol exporter ensure that vnode operations are properly synchronized with respect to any fileset operations in progress. If the fileset in question has been opened with a "cuncurrency mode" that conflicts with a vnode operation, that **vnode** operation blocks until fileset is available again. See (c) below.
- b. While an administrative application has a fileset open (`vol_open()`), any other open requests that arrive for that fileset are refused.
- c. When a fileset is opened, the caller indicates the "type" of operations that will be performed on it. The open blocks until any "conflicting" **vnode** operations currently in progress finish.

For the duration of the open, step (a) protects against additional **vnode** operations and step (b) protects against additional fileset (volume) operations. Consult the `vol_open()` and `vol_concurr()` operations.

- d. When a fileset is closed, blocked vnode operations are allowed to proceed. There are a number of vnode operations for which the glue does not perform any synchronization, however. The two reasons for these are: (1) the calls are not allowed to block just because a fileset operation is in progress, or (2) there's no harm in allowing the call to proceed concurrently with a volume operation.

In these cases, the LFS appears to be left "on its own" to protect against undesirable (if any are possible) interactions. A typical set of these glue functions is listed below:

- `vn_inactive()`

This call cannot block waiting for a series of volume operations to complete. The LFS is on its own to ensure that the correct thing eventually happens (like, ZLC file deletion) if the fileset in question is open (via `vol_open()`).

- `vn_bmap()`
- `vn_strategy()`
- `vn_bread()`
- `vn_brelse()`
- `vn_select()`
- `vn_ioctl()`
- `vn_pathconf()`
- `vn_fpathconf()`
- `vn_fid()`

It appears this operation left unprotected for performance reasons since it doesn't do anything that could get it into trouble with respect to concurrent volume operations.

- `vn_lockctl()`

Due to the blocking nature of this call, the proper synchronization (which itself can block out other users) is not performed. Depending upon the implementation, no synchronization may be necessary.



### 16.5.1 Directory Offsets

The `vn_readdir()` operation returns (when called by the DFS file exporter, at least) directory entries that contain an `.offset` field. DFS assumes that these offsets have the following properties.

- The offset of the first entry in a directory is 0.
- Each directory entry is assumed to have an offset that never changes for that entry. Specifically: the offset for an entry never changes as the result of the addition or removal of other entries.
- These offsets increase for directory entries as they are returned by `vn_readdir()`.
- Within a returned directory entry, the `.offset` field is set to the offset of the NEXT entry to be returned or, if empty space follows that entry, the offset that would be associated with any entry to be placed within that space in the future.
- As discussed under `vn_readdir()`, these offsets are consistent with the `.uio_offset` value input to and returned from the call.

Some file systems provide a mechanism for performing online reorganization or compaction of directories. If such an operation would cause directory entry offsets to change, the following must hold.

- There must be a way of dynamically disabling this functionality.
- Such reorganization must NEVER be attempted on a readonly fileset (either a `.readonly` replica or a `.backup` clone). Since DFS can switch between replicas at ANY time, these replicas must appear identical, even down to these directory offsets.
- Such reorganization must advance the directory's `1.dataVersion` attribute to allow DFS clients to recognize that any directory offsets being cached might be invalid.

Many file systems, including the standard UNIX UFS, will automatically de-fragment (compact) a directory block during an insertion if there is insufficient space to add a desired entry. Although there may be no way to inhibit this behavior, the file system should at least advance the directory's `.dataVersion` attribute as mentioned above.

## 16.6 Vfs Operations

Since DFS doesn't generally use the `vfs` operations, they have not been extended or standardized the way that the `vnode` ones have. The DFS requirements that do exist on these operations are listed below:

`vfs_vget()` Obtain a `vnode` given an NFS-style `fid` that was originally handed out by `vn_fid()`.

As discussed under `vn_fid()` (see Section 16.7 on page 459), the `anode` index and generations within the `fid` are represented in network (as in `hton32()`) order.

`vfs_statfs()` Obtain statistics for a file system. See Section 15.14 on page 340. More specifically, see Section 13.3.5 on page 279 to see the dynamic status portion of the fileset's status.

The returned `.f_blocks` and `.f_free` fields should be computed from the fileset visible quota fields (`.visQuotaLimit` and

Since filesets do not possess fixed size anode tables, something "appropriate" should be set for the returned `.f_files` and `.f_ffree` fields.

**Note:** `Episode` returns `-1` for these.

### 16.6.1 `vfs_getvolume()`

The `vfs_getvolume()` function below is not part of the standard vfs operations vector. Instead, it is pointed at by the `vfsgetvolume` pointer in the extended `xvfs_vfsops` structure established when a fileset is exported. See *Define the Enhanced Operations Vector*, in Section 13.6.8 on page 299 for information about its fields. There, expect to see the pointer just mentioned.

The `vfs_getvolume()` function has the following signature:

```
vfs_getvolume(
    /* IN */ struct osi_vfs   *vfsp,
    /* OUT */ struct volume  **volpp
);
```

It's description is as follows:

**DESCRIPTION** Obtain a volume structure from a vfs pointer.

**ARGUMENTS**

`vfsp` A pointer to a vfs structure.

`volpp` A pointer to a held volume structure is returned here.

**DISCUSSION** This call should obtain a fileset ID from the supplied vfs structure and then call `volreg_Lookup()` with that in order to obtain a held fileset structure to return.

See Section 15.15 on page 399 for a description of `volreg_Lookup()`

**RETURN VALUES** This call should return whatever value `volreg_Lookup()` returned to it. If all is successful, a 0 will be returned. If the fileset does not reside locally, [ENODEV] will be returned instead.

## 16.7 Base Vnode Interface

This chapter describes the **X-ops**. As discussed in Section 11.4 on page 238, this set of functions is one of two extensions to the array of vnode operations (or equivalent array of operations in some kernels). The **X-ops** do not perform their own synchronization with clients, but assume that it has been done by the caller. Also, the **X-ops** have an interface that is nearly identical across platforms, making them suitable for use by the Protocol Exporter and other components of the DFS File Server. Most of the **X-op** function specifications are derived from the original Sun Microsystems **vnode** architecture, as described in a paper by Kleiman appearing in the proceedings of the 1986 Summer Usenix Conference<sup>11</sup>. Those specifications are given in this chapter. The next section, *Extended Vnode Interface*, Section 16.8 on page 460, gives specifications for some additional **X-ops**, added to meet the needs of DFS.

A description of the operations available within the vnode function array appears below. The array is an exported data type. (The associated declaration is found in Section 13.6.4 on page

11. S. R. Kleinman. *Vnodes: An Architecture for Multiple File System Types in Sun UNIX*, Proceedings Summer Usenix Technical Conference & Exhibition, pp. 238-247, Atlanta, Georgia, June 1986.

296). It is reproduced here for convenience.

```

    /* vnode operations */
struct xvfs_xops {
    int      (*vn_open)();
    int      (*vn_close)();
    int      (*vn_rdwr)();
    int      (*vn_ioctl)();
    int      (*vn_select)();
    int      (*vn_getattr)();
    int      (*vn_setattr)();
    int      (*vn_access)();
    int      (*vn_lookup)();
    int      (*vn_create)();
    int      (*vn_remove)();
    int      (*vn_link)();
    int      (*vn_rename)();
    int      (*vn_mkdir)();
    int      (*vn_rmdir)();
    int      (*vn_readdir)();
    int      (*vn_symlink)();
    int      (*vn_readlink)();
    int      (*vn_fsync)();
    int      (*vn_inactive)();
    int      (*vn_bmap)();
    int      (*vn_strategy)();
    int      (*vn_ustrategy)(); /* assuming stuff already mapped in */
    int      (*vn_bread)();
    int      (*vn_brelse)();
    int      (*vn_lockctl)();
    int      (*vn_fid)();      /* op for old style fid op */
    int      (*vn_hold)();     /* maybe don't need these; revisit */
    int      (*vn_rele)();
    /*
     * new ones for us to provide, rather than just existing to
     * make writing the 0 functions easier (i.e. porting).
     */
    int      (*vn_setacl)();
    int      (*vn_getacl)();
    int      (*vn_afsfid)();
    int      (*vn_getvolume)();
    int      (*vn_getlength)();
    /*
     * Some new ops for AIX 3
     */
    int      (*vn_map)();      /* also used for SunOS 5 */
    int      (*vn_unmap)();
    /*
     * A new op for OSF/1
     */
    int      (*vn_reclaim)();
    /*
     * Some new ops for SunOS 5

```

```

    */
    int (*vn_read)();
    int (*vn_write)();
    int (*vn_realvp)();
    void (*vn_rwlock)();
    void (*vn_rwunlock)();
    int (*vn_seek)();
    int (*vn_space)();
    int (*vn_getpage)();
    int (*vn_putpage)();
    int (*vn_addmap)();
    int (*vn_delmap)();
    int (*vn_pageio)();
#define vn_frlock vn_lockctl          /* overlay equivalent ops */
    /*
     * Ops for HP/UX
     */
    int (*vn_pagein)();
    int (*vn_pageout)();
};

```

The following descriptions do not list all possible error cases and return values. The primary intention has been to identify DFS specific behavior that differs from what would be expected in an *ordinary vnode*-based file system implementation. Some operations have a more comprehensive set of error returns defined than others.

For *vn\_xxx()* operations that are not expected to be called (and, hence, not implemented), it is recommended that the operations vector point to a function which returns [ENOSYS].

**NAME**

vn\_open — Open a file associated with a vnode

**SYNOPSIS**

```
int vn_open(  
    /* INOUT */ struct vnode **avpp,  
    /* IN */     long          aflags,  
    /* IN */     struct ucred  *acred  
);
```

**ARGUMENTS**

<i>avpp</i>	Pointer to the vnode to open.
<i>aflags</i>	Open flags.
<i>acred</i>	Pointer to caller's credential structure.

**DESCRIPTION**

Perform an open protocol on the given *avpp* vnode pointer. The open flags are provide by *aflags*, and the caller's credentials are passed in *acred*. A new vnode may be created, and is passed back in *avpp*.

**DISCUSSION**

This operation is for local access ony. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

None.

**NAME**

vn\_close — Close a file associated with a vnode.

**SYNOPSIS**

```
int vn_close(  
    /* IN */ struct vnode *avp,  
    /* IN */ long         aflags,  
    /* IN */ struct ucred *acred  
);
```

**ARGUMENTS**

*avp*                    Pointer to the vnode to close.  
*aflags*                Flags used with the open.  
*acred*                 Pointer to caller's credential structure.

**DESCRIPTION**

Perform a close protocol on the given *avp* vnode pointer. The flags passed to the original *vn\_open()* are passed in *aflags*, and the caller's credentials are passed in *acred*.

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

None.

**NAME**

vn\_rdwr — Read or Write data from or to a vnode

**SYNOPSIS**

```
int vn_rdwr(
    /* IN */ struct vnode *avp,
    /* IN */ struct uio *auio,
    /* IN */ enum uio_rw arw,
    /* IN */ int aio,
    /* IN */ struct ucred *acred
);
```

**ARGUMENTS**

<i>avp</i>	Pointer to the vnode.
<i>auio</i>	Pointer to the user structure which supplies the I/O arguments, including the address in user space where the data to be read (written) are.
<i>arw</i>	Specifies the direction of the I/O. For a read, specify (UIO_READ); for write, (UIO_WRITE).
<i>aio</i>	The associated I/O flags such as IO_APPEND or (and) IO_SYNC.
<i>acred</i>	Pointer to caller's credential structure.

**DESCRIPTION**

Read or write data from/to the vnode pointed to by *avp*. The user structure which supplies the I/O arguments is pointed to by *auio*, and includes the address in user space where the data is to be read from or written to. The I/O flags contained in *aio* may specify that, if the direction of I/O is output, the output is to be done synchronously. Other flags may be present but will generally be ignored for regular disk files (they may be used by device special files, FIFO's, and the like). If input is done, the file's atime will be modified. If output is done, the file's mtime and ctime will be modified. Under some obscure circumstances, the file's mode may be modified, depending on the caller's credentials.

**DISCUSSION**

The specified file is either read or written, based on the *arw* argument. The actual transfer is controlled by parameters supplied in the **uio** structure and *aio* word.

**Note:** A symlink (its link text) can be accessed via this call.

For a write (UIO\_WRITE):

The usual tests for maximum file size (ulimit) are performed.

IF IO\_SYNC is present, a synchronous write is performed. Unless this call is being issued by ROOT, the file's set-uid and set-group (VSUID, VSGID) flags are cleared.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL]	Returned for a FIFO, BLK or CHR special device.
[EROFS]	Returned if the fileset is read-only.

**NAME**

vn\_ioctl — Handle I/O control on an open file descriptor

**SYNOPSIS**

```
int vn_ioctl(  
    /* IN */ struct vnode *avp,  
    /* IN */ long         acmd,  
    /* IN */ char         *adata,  
    /* IN */ long         aflags,  
    /* IN */ struct ucred *acred  
);
```

**ARGUMENTS**

<i>avp</i>	Pointer to vnode to affect.
<i>acmd</i>	Command to perform.
<i>adata</i>	Pointer to the data involved.
<i>aflags</i>	The associated open flags.
<i>acred</i>	Pointer to the caller's credential structure.

**DESCRIPTION**

Perform an *ioctl()* on the vnode pointed to by *avp*. The command to perform is found in *acmd*, the data involved in the command is found in *adata*, the open flags associated with the file are in *aflags*, and *acred* points to the caller's credential structure.

This function is more or less irrelevant for files. Both the Cache Manager and Episode functions return [EINVAL] without doing anything.

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

None.



**NAME**

vn\_select — Perform a select on a vnode

**SYNOPSIS**

```
int vn_select(  
    /* IN */ struct vnode *avp,  
    /* IN */ long          awhich,  
    /* IN */ struct ucred *acred  
);
```

**ARGUMENTS**

*avp* Pointer to the vnode on which to do the select.

*awhich* Specifies the I/O direction.

*acred* Pointer to the caller's credential structure.

**DESCRIPTION**

Perform a select on the vnode pointed to by *avp*. The direction of the I/O is specified in *awhich*, and the caller's credential structure is pointed to by *acred*.

This function is more or less irrelevant for files. Both the Cache Manager and Episode functions do nothing. On most platforms they return [EINVAL], but on some platforms they return 0.

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

None.

## NAME

vn\_getattr — Get the attributes for a vnode

## SYNOPSIS

```
int vn_getattr(
    /* IN */ struct vnode *avp,
    /* IN */ struct vattr *aattrs,
    /* IN */ int aflag,
    /* IN */ struct ucred *acred
);
```

## ARGUMENTS

*avp* Pointer to the vnode whose attributes are desired.

*aattrs* Pointer to the buffer in which to place the given vnode's attributes. If *aflag* is zero, this is a **vattr** structure, else a **xfvs\_vattr** structure.

*aflag* Flag indicating whether *aattrs* is an extended attribute structure.

*acred* Pointer to the caller's credential structure.

## DESCRIPTION

Get the attributes associated with the vnode pointed to by *avp*. The caller provides *aattrs*, a pointer to a **struct vattr** in which to place the attributes. Also passed is a pointer to the caller's credential structure, *acred*. Vnode attributes include much of the information reported by the *stat* system call, such as the owner's **uid** and **gid**, the file size, and access and modify times.

## DISCUSSION

The standard **vnode** attributes at *aattrs->vattr* are returned in an obvious manner. Fields that might require special attention are listed below.

*.va\_fsid* The value returned here should be chosen with care since it needs to be unique (per fileset) and a given aggregate (device) can hold several filesets. Although an individual LFS is free (seemingly) to use any appropriate algorithm to construct this field, the following is done by Episode:

```
.va_fsid = (volumeID.low << 16) |
            (device major # << 8) |
            (device minor #)
```

*.va\_rdev* The device number of the aggregate on which this fileset resides is returned.

*.va\_blocks* The blocks used value is based on the logical amount of disk space used by the file (as if there were no copy-on-write (clone) sharing taking place). This corresponds to the visible (versus actual allocated) quota on a fileset.

If the flag argument is non-zero, the extended attributes at *aattr->xvattr* are returned as well. Definitions of the fields in this structure and the values to be returned for them can be found earlier in this document. Fields requiring special attention are listed below.

*.anonAccess* If the object in question does not have an **ACL**, 0 is returned in this field. Likewise, if the *any\_other* entry in the **ACL** is invalid, a 0 is returned. Otherwise, it is set to the logical intersection of EVERY valid entry in its **ACL**.

The *acred* argument is used in the computation of the *.callerAccess field*.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[non-zero] This function returns an error.

## NAME

vn\_setattr — Set the attributes for a vnode

## SYNOPSIS

```
int vn_setattr(
    /* IN */ struct vnode *avp,
    /* IN */ struct vattr *aattrs,
    /* IN */ int aflag,
    /* IN */ struct ucred *acred
);
```

## ARGUMENTS

*avp* Pointer to the vnode whose attributes are to be set.

*aattrs* Pointer to the buffer from which to get the given vnode's attributes.

*aflag* Flag indicating whether *aattrs* is an extended attribute structure.

*acred* Pointer to the caller's credential structure.

## DESCRIPTION

Set the attributes of the vnode pointed to by *avp*. The caller provides *aattrs*, a pointer to a **struct vattr** from which the attributes are obtained. In this structure, a value of  $-1$  in any field indicates that the corresponding attribute is not to be changed. Also passed is a pointer to the caller's credential structure, *acred*. This function is called on behalf of various system calls, such as **chmod**, **chown**, **utimes**, and **ftruncate**.

## DISCUSSION

The standard vnode attributes at *aattrs*->*vattr* are applied to the designated file as follows.

<i>.va_ctime</i>	If the <i>.tv_sec</i> sub-field is not equal to $-1$ . This field is <b>ONLY</b> settable if the flags argument is non-zero.
<i>.va_nlink</i>	If not equal to $-1$ . If the link count is being set to zero, the file will be deleted by <i>vn_inactive()</i> when its vnode reference count goes to zero. As described earlier in <i>Zero Link Count Files</i> , Section 15.6 on page 332, such files are preserved across reboots and FSCKs until <i>vn_inactive()</i> explicitly deletes them.
<i>.va_mode</i>	If not equal to $-1$ .
<i>.va_uid</i>	If not equal to $-1$ .
<i>.va_gid</i>	If not equal to $-1$ .
<i>.va_size</i>	If not equal to $-1$ . The file's modified time (mtime) is set to the current time <b>UNLESS</b> it is being explicitly set ( <i>.va_mtime</i> ). If a directory is being operated on, an [EISDIR] error is returned.
<i>.va_mtime</i>	If the <i>.tv_sec</i> sub-field is 0 and the <i>.tv_usec</i> sub-field is $-1$ , the file's atime and mtime are both set to the current time. (System V style <i>utime()</i> call.)  Alternatively, if the <i>.tv_sec</i> field is not equal to $-1$ , the file's mtime is set.
<i>.va_atime</i>	If the <i>.tv_sec</i> sub-field is not equal to $-1$ and the above <i>.va_mtime</i> case does not supercede this one, the file's atime is set.

If the flag argument is non-zero, the extended attributes at *aattrs*->*xvattr* are additionally set as follows.

<i>.volVersion</i>	If the <i>.high</i> or <i>.low</i> sub-field is not equal to -1.
<i>.dataVersion</i>	If the <i>.high</i> or <i>.low</i> sub-field is not equal to -1.
<i>.fileID</i>	The anode generation value, from the <i>.low</i> sub-field, can be changed if it is not equal to -1. The <i>.high</i> field, the anode index, is ignored.
<i>.clientOnlyAttrs</i>	If not equal to -1.

Unless the file's ctime is explicitly set, it is set to the current time by this call.

If setting *.va\_mode*, the set-uid (VSVTX) flag is cleared if the caller is not ROOT and the set-gid (VSGID) flag is cleared if the caller is not a member of the file's group. Other actions (regarding the sticky (VTEXT) flag) might be required as well, depending on the OS.

Administrative rights to the object are required for the caller if any of the following fields are being set: *.va\_mode*, *.va\_uid*, *.va\_gid*, *.va\_atime*, *.va\_mtime*. If *.va\_mtime* and *.va\_atime* are being set to the current time, either administrative rights or write access to the file are sufficient.

#### Notes:

1. If flags is non-zero, an extended set-attribute attribute is being performed and DFS (at the client side) has performed its own permissions check.
2. Episode returns EROFS if the *.va\_uid* or *.va\_gid* is being changed and the file resides on a readonly fileset

Any permissions checks performed by this operation must be performed against the original *mode*, *uid* and *gid* fields (as opposed to the new ones set by this call).

DFS expects to be able to change or set the size of a file even if the caller does not currently own or have write access to the file.

#### RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

#### ERRORS

[EISDIR]	An attempt was made to change the length of a directory with this call.
[EPERM]	The caller possesses insufficient rights to perform this operation.
[EROFS]	The fileset is read-only.

**NAME**

`vn_access` — Check access permissions for a vnode

**SYNOPSIS**

```
int vn_access(  
    /* IN */ struct vnode *avp,  
    /* IN */ long         amode,  
    /* IN */ struct ucred *acred  
);
```

**ARGUMENTS**

*avp* Pointer to the vnode to examine.

*amode* Mode to check for access. Values are S\_IREAD, S\_IWRITE, S\_IEXEC.

*acred* Pointer to the caller's credential structure.

**DESCRIPTION**

Check the access permission for the vnode pointed to by *avp*. The access mode to be checked (read, write, execute) is specified by *amode*, and a pointer to the caller's credential structure is provided by *acred*. If the specified access is denied, `vn_access()` returns a non-zero value.

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success. The caller has the rights specified in the mode argument for the object in question.

**ERRORS**

[EROFS] The fileset is read-only and the caller is asking for writth access. This is either a *.backup* clone or a *.readonly* replica.

[EACCES] The caller does not have the specified rights to the object.

**NAME**

vn\_lookup — Look up a component name in a directory vnode

**SYNOPSIS**

```
int vn_lookup(
    /* IN */ struct vnode *adp,
    /* IN */ char *aname,
    /* OUT */ struct vnode **avpp,
    /* IN */ struct ucred *acred
);
```

**ARGUMENTS**

<i>adp</i>	Pointer to the directory vnode to be searched.
<i>aname</i>	Character string name being looked for.
<i>avpp</i>	A place to store the address of the vnode corresponding to the name being looked up (if one exists).
<i>acred</i>	Pointer to the caller's credential structure.

**DESCRIPTION**

Look up a pathname component name *aname* in the directory represented by the *adp* vnode pointer. The vnode for the corresponding object (if one exists in the given directory) is stored in *avpp*. The caller's credentials are passed in through the *acred* pointer.

**DISCUSSION**

If this call succeeds, it returns a pointer to a held (as in *vn\_hold()*), for the object that was found.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success. The caller has the rights specified in the mode argument for the object in question.

**ERRORS**

[EACCES]	The caller does not have the search rights to the specified directory.
[ENOENT]	The name is not found or the directory has been deleted.

## NAME

vn\_create — Create a new file in the given directory

## SYNOPSIS

```
int vn_create(
    /* IN */ struct vnode *adp,
    /* IN */ char *aname,
    /* IN */ struct vattr *aattrs,
    /* IN */ enum vcexcl aexcl,
    /* IN */ int amode,
    /* OUT */ struct vnode **avpp,
    /* IN */ struct ucred *acred
);
```

## ARGUMENTS

<i>adp</i>	Pointer to the directory vnode in which the new file is to be created.
<i>aname</i>	Character string name for the file to create.
<i>aattrs</i>	Pointer to the attribute structure to give to the new file.
<i>aexcl</i>	Exclusive or non-exclusive creation flag. If exclusive, the specified name ( <i>aname</i> ) must not already reside within the directory.
<i>amode</i>	Access rights required if file exists in a non-exclusive creation (S_IREAD, S_IWRITE, S_IEXEC).
<i>avpp</i>	Set to the address of the vnode corresponding to the newly-created file, if successful.
<i>acred</i>	Pointer to the caller's credential structure.

## DESCRIPTION

Create a file named *aname* in the directory associated with the *adp* vnode. The new file's information is passed in via the following parameters: *aattrs* contains the initial attribute structure, and *aexcl* specifies whether the create is exclusive or non-exclusive. *amode* is not the mode of the new file (that is in the appropriate field in *aattrs*), but specifies access rights that the existing file must allow if there is an existing file and the create is non-exclusive. It is derived from the flags argument to the *open* system call. A pointer to the caller's credentials are provided via *acred*. If *vn\_create()* is successful, *avpp* is set to the address of the vnode pointer describing the new file.

## DISCUSSION

- **If the desired object already exists:**

If the exclusive argument is non-zero, the call fails with [EEXIST].

If the existing object is a directory and the mode argument indicates that the caller wants write access (S\_ISWRITE), the call fails with [EISDIR].

If the caller does not have the rights given by the mode argument to the file, the call fails with [EACCES].

If the existing object is an ordinary file and the supplied *aattrs->va\_size* field is 0, it is truncated to zero length.

- **If the desired object does not already exist:**



The new file's user field (*va\_uid*) is taken from the DFS authentication structure which itself is obtained from the supplied credentials. Depending on whether SYS5 or BSD semantics are being adhered to, the group field (*va\_gid*) is either taken from this PAC or from the parent directory. Consult *Obtaining the Identity of the Principal*, Section 11.8.10 on page 243. If the PAC indicates that this is an unauthenticated RPC operation, the *user* and *group* fields are set instead to -2 and -1, respectively.

**Note:** If BSD behavior is in effect (*va\_gid* from parent directory), then it seems equally valid to form the *va\_gid* this way even for an unauthenticated access. (Episode does not do this.)

The supplied mode bits (in the attributes structure) are used as supplied if (a) a link (VLNK) is being created or (b) the parent directory possesses an *Initial File ACL*. Otherwise, the mode bits are weakened by the process *umask*.

**Note:** For the creation of directories, files and symlinks, DFS sets the file exporter's *umask* to that of the remote client.

Ordinary files, FIFOs and device files can be created via this call. For the latter, the supplied *va\_rdev* device number is applied.

The new file's link count is set to 1. Its *atime*, *mtime* and *ctime* are all set to the current time.

The extended object UUID (as in the *.objid* field of a **Txvattr** structure) is set to zero.

Any user or system property lists (PLIST) on the parent directory are inherited by the newly created file.

Consult *Initial ACL and File Creation - Algorithm*, Section 12.11 on page 258 for additional **ACL** processing that might be required for the new object.

The parent directory's *ctime* and *mtime* fields are set to the current time.

#### RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

#### ERRORS

- |           |   |
|-----------|---|
| [EISDIR]  | An attempt is made to create a directory with this call.<br><br>Also returned if the desired object already exists, is a directory and the <i>amode</i> argument indicates that the caller desired write access (S_IWRITE).   |
| [ENOTDIR] | The object indicated by the <i>adp vnode</i> is not a directory.  |
| [ENOENT]  | The name is not found or the directory has been deleted.  |
| [EEXIST]  | An attempt is mad to create a "." or "..", or if the name already exists and the exclusive flag is non-zero.  |
| [EACCES]  | The caller does not have the search, write and insert rights to the specified directory ( <i>perm_execute</i> , <i>perm_write</i> , <i>perm_insert</i> ).<br><br>Also returned if the desired file already exists and the called does not have the rights specified in the mode argument to the file. |

#### SEE ALSO

*Initial ACL and File Creation - Algorithm*, Section 12.11 on page 258. Also, *Obtaining the Identity of the Principal*, Section 11.8.10 on page 243.

## NAME

vn\_remove — Delete a file in the given directory

## SYNOPSIS

```
int vn_remove(
    /* IN */ struct vnode *adp,
    /* IN */ char *name,
    /* IN */ struct ucred *acred
);
```

## ARGUMENTS

*adp* Pointer to the directory vnode from which file is to be deleted.

*aname* Character string name for the file to delete.

*acred* Pointer to the caller's credential structure.

## DESCRIPTION

Delete a file named *aname* from the directory corresponding to the vnode pointed to by *adp*. A pointer to the caller's credential structure is found in *acred*.

In some kernels, the corresponding vnode operation is passed vnodes for both the parent directory and the file to be deleted. That is, the kernel is expected to look up the file before calling the vnode operation. When this *X-op* is called on behalf of such a vnode operation, the lookup done by the kernel is wasted, as the *X-op* does the lookup itself.

## DISCUSSION

The specified name is removed from the parent directory and the link count on the object is decremented. The parent's mtime and ctime fields are set to the current time. The ctime of the object being removed is set to the current time its link count is being decremented.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

## ERRORS

[ENOTDIR] The object indicated by the *adp* **vnode** is not a directory.

[EPERM] The specified object is a directory.  
**Note:** The local OS might dictate that ROOT be allowed to *vop\_remove()* a directory. In such cases, DFS access should fail.

[ENOENT] The name (by *aname*) is not found in the parent directory.

[EINVAL] An attempt was made to remove "." or ".." by this call.

[EACCES] The caller does not have the search, write and delete rights to the specified directory (perm\_execute, perm\_write, perm\_delete).

## SEE ALSO

*Zero Link Count Files*, Section 15.6 on page 332.

**NAME**

vn\_link — Create a hard link

**SYNOPSIS**

```
int vn_link(
    /* IN */ struct vnode *avp,
    /* IN */ struct vnode *adp,
    /* IN */ char          *aname,
    /* IN */ struct ucred *acred
);
```

**ARGUMENTS**

<i>avp</i>	Pointer to the vnode to be hard-linked.
<i>adp</i>	Pointer to the directory vnode in which the hard link is to be made.
<i>aname</i>	Target name for the hard link.
<i>acred</i>	Pointer to the caller's credential structure.

**DESCRIPTION**

Link the vnode *avp* to the target name *aname* in the target directory associated with vnode pointer *adp*.

**DISCUSSION**

An entry is made within the target directory which links to the supplied source object. The link count on the source object is incremented. The ctime on the source object along with the ctime and mtime on the target directory are set to the current time.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[ENOTDIR]	The object indicated by the <i>adp</i> <b>vnode</b> is not a directory.
[EACCES]	The caller does not have the search, write and insert rights to the specified directory (perm_execute, perm_write, perm_insert).
[EXDEV]	The source object and target directory are on different filesets.
[EPERM]	The specified object is a directory.
	<b>Note:</b> The local OS might dictate that ROOT be allowed to <i>vop_remove()</i> a directory. In such cases, DFS access should fail.
[EEXIST]	The name (by <i>aname</i> ) already exists in the target directory.

**NAME**

vn\_rename — Rename a file

**SYNOPSIS**

```
int vn_rename(
    /* IN */ struct vnode *aodp,
    /* IN */ char          *aname1,
    /* IN */ struct vnode *andp,
    /* IN */ char          *aname2,
    /* IN */ struct ucred *acred
);
```

**ARGUMENTS**

*aodp* Pointer to the directory vnode where the file currently exists.

*aname1* Name of the current file within the above directory.

*andp* Pointer to the directory vnode to which the file will move.

*aname2* The new name of the file.

*acred* Pointer to the caller's credential structure.

**DESCRIPTION**

Rename a file currently named *aname1* in the directory associated with vnode pointer *aodp*. The file's new name will be *aname2*, and it will be moved to the directory associated with vnode pointer *andp*. A pointer to the caller's credential structure is provided by *acred*.

In some kernels, the corresponding vnode operation is passed vnodes for the parent directories and the subfiles (or subdirectories). That is, the kernel is expected to look up the files before calling the vnode operation. When this **X-op** is called on behalf of such a vnode operation, the lookups done by the kernel are wasted, as the **X-op** does the lookups itself.

**DISCUSSION**

If an object with the target's name (*aname2*) already exists:

- If it is a directory, it must be empty.
- If it is a directory, the link count on the target directory is decremented by 1 to account for the removed "." reference.
- Its link count is decremented by either 1 if it is a file, or 2 if a directory.

If the source object (*aname1*, object being renamed) is a directory, the source directory has its link count decremented by 1 and the target directory has its link count incremented by 1 to account for the moved "." reference.

The mtime and ctime of both the source and target directories are set to the current time. The ctime of the source object (*aname1*) is set to the current time.

If the target object (*aname2*) originally existed, its ctime is set to the current time as well.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[ENOENT] The source name does not exist within the source directory.

- [EXDEV] The source object and target directory are on different filesets.
- [ENOTDIR] The object indicated by the *andp* (target) **vnode** is not a directory.  
The object indicated by the *aodp* (source) **vnode** is not a directory.  
The source object (*aname1*) is a directory, the target object (*aname2*) exists and that target object is not a directory.
- [EISDIR] The source object (*aname1*) is not a directory, the target object (*aname2*) exists and that target object is a directory.
- [ENOTEMPTY] The target object (*aname2*) exists and is a not-empty directory.
- [EINVAL] The source object (*aname1*) is a directory, and the target directory is a descendant of it.  
The source object (*aname1*) is either "." or "..".
- [EACCES] The caller does not have the search, write and delete rights to the source directory (perm\_execute, perm\_write, perm\_delete).  
The caller does not have the search, write and insert rights to the target directory (perm\_execute, perm\_write, perm\_insert).  
The directories are different, the source object (*aname1*) is a directory and the caller does not have write access (perm\_write) to that source object.

## NAME

vn\_mkdir — Create a directory

## SYNOPSIS

```
int vn_mkdir(
    /* IN */ struct vnode *adp,
    /* IN */ char *aname,
    /* IN */ struct vattr *aattrs,
    /* OUT */ struct vnode **avpp,
    /* IN */ struct ucred *acred
);
```

## ARGUMENTS

<i>adp</i>	Pointer to the directory vnode in which the new directory is to be created.
<i>aname</i>	Character string name for the directory to create.
<i>aattrs</i>	Pointer to the attribute structure to give to the new directory.
<i>avpp</i>	If successful, set to the pointer of the newly-created directory's vnode.
<i>acred</i>	Pointer to the caller's credential structure.

## DESCRIPTION

Create a directory named *aname* in the directory associated with the *adp* vnode. The new directory's initial attribute structure is passed in *aattrs*, and a pointer to the caller's credential structure appears in *acred*. If *vn\_mkdir()* is successful, *avpp* is set to the address of the vnode pointer describing the new directory.

## DISCUSSION

An empty directory, containing just the "." and ".." entries, is created.

The new directory's user field (*va\_uid*) is taken from the DFS authentication structure which itself is obtained from the supplied credentials. Depending on whether SYS5 or BSD semantics are being adhered to, the group field (*va\_gid*) is either taken from this PAC or from the parent directory. Consult *Obtaining the Identity of the Principal*, Section 11.8.10 on page 243. If the PAC indicates that this is an unauthenticated RPC operation, the *user* and *group* fields are set instead to -2 and -1, respectively.

**Note:** If BSD behavior is in effect (*va\_gid* from parent directory), then it seems equally valid to form the *va\_gid* this way even for an unauthenticated access. (Episode does not do this.)

The supplied mode bits (in the attributes structure) are used as supplied if the parent directory possesses an *Initial File ACL*. Otherwise, the mode bits are weakened by the process *umask*.

**Note:** For the creation of directories, files and symlinks, DFS sets the file exporter's *umask* to that of the remote client.

The new directory's link count is set to 2. Its *atime*, *mtime* and *ctime* are all set to the current time.

The extended object UUID (as the *.objid* field of a **Txvattr** structure) is set to zero.

Any user or system property lists (PLIST) on the parent directory are inherited by the newly created directory.

The parent directory's *ctime* and *mtime* fields are set to the current time and its link count is incremented to account for the new ".." reference.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[ENOTDIR] The object indicated by the *theadp* (target) **vnode** is not a directory.

[EEXIST] The directory name already exists.

[EACCES] The caller does not have search, write and insert rights to the target directory (perm\_execute, perm\_write, perm\_insert).

**SEE ALSO**

Also, *Obtaining the Identity of the Principal*, Section 11.8.10 on page 243.

## NAME

vn\_rmdir — Delete a directory

## SYNOPSIS

```
int vn_rmdir(
    /* IN */ struct vnode *adp,
    /* IN */ char *aname,
    /* IN */ struct ucred *acred
);
```

## ARGUMENTS

*adp* Pointer to the directory vnode in which the directory appears.

*aname* Character string name for the directory to delete.

*acred* Pointer to the caller's credential structure.

## DESCRIPTION

Delete a directory named *aname*, which appears within the directory whose vnode is pointed to by *adp*. The caller's credentials are passed in via *acred*.

In some kernels, the corresponding vnode operation is passed vnodes for both the parent directory and the file to be deleted. That is, the kernel is expected to look up the file before calling the vnode operation. When this **X-op** is called on behalf of such a vnode operation, the lookup done by the kernel is wasted, as the **X-op** does the lookup itself.

## DISCUSSION

If the specified directory (*aname*) is empty, its name is removed, its link count is decremented by two (name plus "." reference) and its ctime is set to the current time.

The containing directory (*adp*) has its link count decremented by 1 (to account for the "." reference) and its mtime and ctime set to the current time.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

## ERRORS

[ENOENT] The specified directory (*aname*) does not exist.

[EINVAL] The specified directory (*aname*) is either "." or "..".

[ENOTDIR] The object specified by the *adp* **vnode** is not a directory.

[EACCES] The caller does not have the search, write and delete rights to the containing directory (perm\_execute, perm\_write, perm\_delete).

[EBUSY] The directory being deleted is mounted on top of.

[ENOTEMPTY] The directory being deleted is not empty.

## SEE ALSO

*Zero Link Count Files*, Section 15.6 on page 332.



**NAME**

vn\_readdir — Read entries from a directory

**SYNOPSIS**

```
int vn_readdir(
    /* IN */ struct vnode *adp,
    /* IN */ struct uio   *auio,
    /* IN */ struct ucred *acred,
    /* IN */ int          *eofp,
    /* IN */ int          isPX
);
```

**ARGUMENTS**

<i>adp</i>	Pointer to the directory vnode being read from.
<i>auio</i>	Pointer to user structure which supplies the I/O arguments, including the address in user space into which the directory data are to be read.
<i>acred</i>	Pointer to the caller's credential structure.
<i>eofp</i>	Optional. A pointer to a returned EOF (end of file) flag.
<i>isPX</i>	Optional. If non-zero, this call is being made by the file exporter on behalf of a remote DFS client.

**DESCRIPTION**

Read entries from directory *adp*. The user structure pointed to by *auio* supplies the I/O arguments. The **uio** offset is set to a file system-dependent number which represents the logical offset in the directory when the reading is done. This is necessary because the number of bytes returned by *vn\_readdir()* is not necessarily the number of bytes in the equivalent part of the on-disk directory.

**DISCUSSION**

Directory entries are read into the supplied buffer. The format of these returned entries is controlled by the *isPX* argument. If it is zero, entries are returned in whatever native format the local OS expects. If this argument is non-zero, the DFS file exporter expects to see entries in **struct dirent** format as defined in *Directory Entry Formats*, Section 13.6.11 on page 301. Fields should be filled in as follows:

<i>.offset</i>	The directory offset of whatever (either a valid entry or empty space) FOLLOWS this entry in the directory. This field must be consistent with the <b>uio</b> <i>.uio_offset</i> argument to this call.  Specifically: calling <i>vn_readdir()</i> with a <i>.uio_offset</i> set to the <i>.offset</i> from some directory entry results in any and all entries AFTER that one being returned.
<i>.inode</i>	An LFS-dependent anode number.
<i>.recordlen</i>	The offset (in bytes) from this entry in the output buffer to the next entry. Since directory entries should be aligned on modulo-4 boundaries, this field will be a multiple of 4.
<i>.namelen</i>	The length of the file name field, not including the terminating null character.
<i>.dir_name</i>	A variable sized file name, including a terminating null character.

**Note:** The *isPX* argument is present because the DFS file exporter requires an *.offset* field within directory entries while HPUX (along with OSF1, etc.) does not provide such a field in its native entry format. For AIX and SunOS-5 platforms, this call behaves the

same regardless of the setting of this flag.

As many entries as will fit in the supplied buffer, or as remain in the directory, are actually returned. Entries are returned on mod-4 byte boundaries. (The supplied return buffer has this alignment.)

Fields within the `uio` structure are used as follows.

<code>.uio_iov</code>	On entry: a pointer to an array of io vectors describing the buffer to be read (into).
<code>.uio_iovec</code>	On entry: the number of io vectors pointed at by <code>.uio_iov</code> . Episode fails with [EINVAL] if this number is other than 1.
<code>.uio_offset</code>	On entry: the offset, initially 0, at which to begin "reading" directory entries from.  At exit: the offset (LFS-dependent) at which to continue reading directory entries from on the next call.  This field should be treated consistently with the <code>.offset</code> field in returned entries.
<code>.uio_segflg</code>	UIO_<whatever>
<code>.uio_resid</code>	On entry: the total size of the user supplied buffer.  At exit: the number of bytes in the user's buffer that were NOT read into by this call.  If no directory entries are actually returned, the <code>.uio_resid</code> field should be left unchanged on return.

The `.uio_offset` argument will initially be set to zero if the caller wishes to read the directory from the beginning. On return, it should be updated by the LFS so that a subsequent call with it will continue at the next entry in the directory. As mentioned above, if on input it is set to the `.offset` field from an entry read earlier, the directory read will commence immediately FOLLOWING that entry.

DFS makes no assumptions regarding the ordering of entries within a directory.

The atime of the directory is set to the current time.

**Note:** Episode rounds the supplied output buffer size down to a multiple of 512 bytes. The LAST directory entry returned has its `.recordlen` field incremented to cover the remaining space within the final 512 byte block. (it "points" to the next 512 byte boundary)

This behavior exists only to allow the file exporter to manufacture meaningful directory offsets on platforms whose native entries do not contain such a field. It turns out that the higher level exporter cannot manufacture correct offsets. There is no reason for an LFS to mimic this 512-byte block behavior.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

## ERRORS

[EINVAL] (Optional - not done by Episode) The supplied buffer is not aligned on a 0-modulo-4 byte boundry.

[ENOTDIR]      The object indicated by the *adv* **vnode** is not a directory.

[EACCES]      The caller does not have read rights to the directory.

**SEE ALSO**

*Directory Offsets*, Section 16.5.1 on page 423, and also, *Directory Entry Formats*, Section 13.6.11 on page 301.

## NAME

vn\_symlink — Create a symbolic link

## SYNOPSIS

```
int vn_symlink(
    /* IN */ struct vnode *adp,
    /* IN */ char *aname,
    /* IN */ struct vattr *aattrs,
    /* IN */ char *atargetName,
    /* IN */ struct ucred *acred
);
```

## ARGUMENTS

<i>adp</i>	Pointer to the directory vnode which will contain the new symlink.
<i>aname</i>	Name of the new link, to be inserted into the directory associated with the <i>adp</i> directory vnode pointer.
<i>aattrs</i>	Pointer to the attribute structure to give the new symlink.
<i>atargetName</i>	String name of the path which is the target of the new symlink.
<i>acred</i>	Pointer to the caller's credential structure.

## DESCRIPTION

Create a symlink named *aname* within the directory associated with directory vnode pointer *adp*. This symlink's initial attributes are taken from *aattrs*, and its target pathname is taken from *atargetName*. A pointer to the caller's credential structure is present in *acred*.

## DISCUSSION

The terminating null at the end of *atargetName* is *not* written into the link object. The link file is created as described under *vn\_create()*, with the following exceptions:

- The target name *aname* cannot already exist.
- The supplied *.va\_type* field is ignored and VLNK is used instead.
- Since the permissions on a symlink are never checked, there is no reason for it to inherit an *Initial ACL* from the parent directory.
- The parent directory's *ctime* and *mtime* fields are set to the current time.

## RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

## ERRORS

[ENOTDIR]	The object indicated by the <i>adp</i> <b>vnode</b> is not a directory.
[EACCES]	The caller does not have search, write and insert rights to the target directory.
[EEXIST]	The target name ( <i>aname</i> ) already exists within the target directory.

**NAME**

vn\_readlink — Read a symbolic link's value

**SYNOPSIS**

```
int vn_readlink(
    /* IN */ struct vnode *avp,
    /* IN */ struct uio *auio,
    /* IN */ struct ucred *acred
);
```

**ARGUMENTS**

*avp* Pointer to the symlink vnode being read from.

*auio* Pointer to user structure which supplies the I/O arguments, including the address in user space into which the symlink value is to be read.

*acred* Pointer to the caller's credential structure.

**DESCRIPTION**

Read the value (contents) associated with the symlink associated with the *avp* vnode pointer. The user structure pointed to by *auio* supplies the I/O arguments, including where the symlink's value is to be placed. A pointer to the caller's credential structure is passed via *acred*.

**DISCUSSION**

The symbolic link link is read from the supplied object and deposited into the supplied buffer. A terminating null is *not* appended to the returned link text.

At the time the call is made, the *.uio\_offset* field should have been set to 0 and the *.uio\_resid* field has been set to the total size of the supplied buffer. When the call returns, the *.uio\_resid* field has been decremented by the actual size of the link text copied into the buffer.

The atime of the symbolic link is set to the current time.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL] The object indicated by the *avp* **vnode** is not a symbolic link (VLNK).

**NAME**

`vn_fsync` — Write out all in-memory information for a file

**SYNOPSIS**

```
int vn_fsync(  
    /* IN */ struct vnode *avp,  
    /* IN */ struct ucred *acred  
);
```

**ARGUMENTS**

*avp*                    Pointer to the file to be flushed to disk.  
*acred*                  Pointer to the caller's credential structure.

**DESCRIPTION**

Flush the state of the file associated with vnode pointer *avp* to disk. The credential structure found in *acred* is ignored.

**DISCUSSION**

Any modified data or status for the specified file is written to disk along with data from the VM system.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[`error_status_ok`] This function always returns success.

**ERRORS**

None.

**NAME**

vn\_inactive — The given vnode is no longer referenced

**SYNOPSIS**

```
int vn_inactive(  
    /* IN */ struct vnode *avp,  
    /* IN */ struct ucred *acred  
);
```

**ARGUMENTS**

*avp* Pointer to the vnode no longer being actively referenced by the system.  
*acred* Pointer to the caller's credential structure.

**DESCRIPTION**

The vnode pointed to by *avp* is no longer being referenced by any agent in the vnode layer, and may be safely deallocated. This function is called when a vnode's reference count drops to zero. If the corresponding file also has a zero link count, any space allocated for it on disk may also be freed at this point. The credential structure is ignored.

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

This call is made by the **vnode** (fs independent) layer to indicate that a **vnode** is no longer referenced because its *.v\_count* field has been decremented to 0. The deletion does *not* occur on a readonly fileset (a *.readonly* replica or a *.backup* clone) or a fileset mounted readonly for local access. Files on a *.readonly* replica are only destroyed by an explicit *vol\_depletd()* operation. Files on a *.backup* clone are only destroyed when their fileset is deleted or as the result of a *vol\_reclone()* operation.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**NAME**

vn\_bmap — Logical to physical block number mapping

**SYNOPSIS**

```
int vn_bmap(  
    /* IN */   struct vnode *avp,  
    /* IN */   long         abn,  
    /* OUT */  struct vnode *anvp,  
    /* OUT */  long         anbn  
);
```

**ARGUMENTS**

<i>avp</i>	Pointer to the vnode to consider.
<i>abn</i>	Logical block number to map.
<i>anvp</i>	Pointer to the vnode representing the associated physical device.
<i>anbn</i>	Set to the associated physical block number.

**DESCRIPTION**

Map logical block number *abn* in the file associated with vnode pointer *avp* to its physical device and block number. A vnode pointer describing the corresponding physical device is returned in *anvp*, and *anbn* is set to the physical block number mapped to by *abn*.

This X methods of mapping files to virtual memory (see Section 16.8 on page 460).

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.



**NAME**

vn\_strategy — Asynchronous read/write interface between file blocks and buffers

**SYNOPSIS**

```
int vn_strategy(  
    /* IN */ struct buf *abp  
);
```

**ARGUMENTS**

*abp*                    Pointer to buffer descriptor.

**DESCRIPTION**

This function provides an asynchronous, block-oriented interface to read or write a logical block from a file into or out of a buffer. The only parameter, *bp*, is a pointer to a buffer header which contains a pointer to the affected vnode. This function does *not* write through the buffer cache, and is used by the kernel's virtual memory manager (if any) to perform I/O to and from virtual memory.

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**NAME**

vn\_bread — Read a logical block from a file into a buffer

**SYNOPSIS**

```
int vn_bread(
    /* IN */   struct vnode *avp,
    /* IN */   daddr_t      albn,
    /* OUT */  struct buf    **abpp,
    /* OUT */  long          *asizep
);
```

**ARGUMENTS**

<i>avp</i>	Pointer to the vnode for the file being read from.
<i>albn</i>	Logical block number to read.
<i>abpp</i>	Set to the pointer to a buffer header containing the desired block.
<i>asizep</i>	Pointer to location at which to store the buffer size.

**DESCRIPTION**

Read logical block *albn* from the file associated with the vnode pointed to by *avp*. Parameter *abpp* is set to point to a buffer header pointer which describes the buffer just read in.

This *X-op* is not used in kernels that have more general methods of mapping files to virtual memory (see Section 16.8 on page 460).

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

*Extended Vnode Interface*, Section 16.8 on page 460.

**NAME**

vn\_brelse — Release a buffer header pointer

**SYNOPSIS**

```
int vn_brelse(  
    /* IN */ struct vnode *avp,  
    /* IN */ struct buf *abp  
);
```

**ARGUMENTS**

*avp* Pointer to the vnode to which the buffer header belongs.  
*abp* Pointer to the buffer header to be released.

**DESCRIPTION**

Release the buffer header pointed to by *abp*, which is associated with the vnode pointed to by *avp*. The buffer header pointer was originally generated by *vn\_bread()*.

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**NAME**

vn\_lockctl — File and record locking interface

**SYNOPSIS**

```
int vn_lockctl(  
    /* IN */   struct vnode *avp,  
    /* IN */   struct flock *af,  
    /* IN */   int      acmd,  
    /* IN */   struct ucred *acred  
);
```

**ARGUMENTS**

<i>avp</i>	Pointer to vnode to be locked.
<i>af</i>	Pointer to lock structure.
<i>acmd</i>	Locking command.
<i>acred</i>	Pointer to the caller's credential structure.

**DESCRIPTION**

This function provides System V-style locking on vnodes. A pointer to the vnode to be operated upon is found in *avp*, a pointer to the desired lock structure in *af*, the lock command to be executed in *acmd*, and a pointer to the caller's credential structure in *acred*.

**DISCUSSION**

This operation is for local access only. It is NEVER called by the file exporter.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**NAME**

vn\_fid — Return the file ID associated with a vnode

**SYNOPSIS**

```
int vn_fid(
    /* IN */ struct vnode *avp,
    /* OUT */ struct fid **afidpp
);
```

**ARGUMENTS**

*avp* Pointer to the vnode whose file ID is to be discovered.

*afidpp* Set to the address of a pointer containing the file ID associated with vnode *avp*.

**DESCRIPTION**

Given a pointer to a vnode *avp*, set *afidpp* to point to a file ID pointer describing the file ID associated with the vnode.

This *X-op* is called only on behalf of the NFS server. It is not called on behalf of the Protocol Exporter, which uses a different interface (see Section 16.8 on page 460).

In most kernels, this *X-op* allocates the *struct fid* that is requested; that structure is later freed by the kernel, rather than by another *X-op*. In IBM's AIX kernel, the *struct fid* is allocated by the caller. Accordingly, the parameter *afidpp* is a **struct fid \*** rather than a **struct fid \*\***.

**DISCUSSION**

An NFS-style **fid** (file ID) is allocated, filled in and returned to the caller of this operation. The storage for this **fid** should be allocated via *kmem\_alloc()* since it will be freed by the FS-independent kernel at a later time.

Any information stored within the **fid** (for example, an LFS-dependent anode index and generation number) should be stored in network byte order. Information within these **fid**s must be stored in network byte order to allow the movement of a fileset to a server with a different byte ordering while the **fid** is held at some client.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

## 16.8 Extended Vnode Interface

### 16.8.1 Generalized Differences

A number of *X-ops*, in addition to those described in the *Base Vnode Interface* section, have been added to meet the needs of various components of DFS. There are the following categories:

1. **Support for Protocol Exporter.** The *vn\_afsfid()* *X-op* is used by the Protocol Exporter analogously to the use of the *vn\_fid()* *X-op* by the NFS server.
2. **Support for O-op wrapper functions.** As described in earlier chapters, the glue functions (“*O-ops*”) must perform fileset synchronization to avoid conflicts between individual file references and whole-fileset operations such as cloning, dumping, or restoring. The *vnx\_getvolume()* *X-op* is used to identify a file’s fileset. Referring to the pseudocode of Section C.3 on page 246, this *X-op* would be called from *ReferenceCorrespondingVolume()*.
3. **Support for vnode Operations in newer kernels.** The *vn\_cmp()* and *\_vnctl* *X-ops* were added to support more recent versions of Sun’s SunOS kernels. The *vn\_map()* and *vn\_unmap* *X-ops* were added to support IBM’s AIX kernel. The *vn\_reclaim* *X-op* was added to support OSF’s OSF/1 kernel. Additional additions for SunOS 5 include *vn\_read()*, *vn\_write()*, *vn\_realvp()*, *vn\_rwlock()*, *vn\_rwunlock()*, *vn\_seek()*, *vn\_space()*, *vn\_getpage()*, *vn\_putpage()*, *vn\_addmap()*, *vn\_delmap()*, *vn\_pageio()* and *vn\_frlock()*, which is equivalent to *vn\_lockctl()*. Also, the following were added for HP/UX: *vn\_pagein()*, *vn\_pageout()*.

As noted in *Overview of Interfaces to the LFS*, in Section 11.3 on page 235, these are LFS-specific facilities. They are thus not documented here.

4. **Support for ACLs.** The *X-ops* *vn\_getacl()* and *vn\_setacl()* *X-ops* are used by the Protocol Exporter to read and write ACLs. There is also a system call interface, implemented using these *X-ops*, which makes ACLs available to local users.
5. **Support for VM. Hold, Rele, Lock, Unlock, Pvnlock, Getlength, Setlength, Pin, Unpin, Nlrw, Cached, Update, Readahead, Delaywrite.**

The idea of supporting VM by calling VM-related *X-ops* from the *N-op* wrapper functions is no longer being used in new code, but not all old code has been revised at this writing, and so the VM-related *X-ops* are still included in the *X-op* specifications.

**NAME**

vn\_getvolume — Get the vnode's fileset

**SYNOPSIS**

```
int vn_getvolume(
    /* IN */ struct vnode *avp,
    /* OUT */ struct volume *avolpp
);
```

**ARGUMENTS**

*avp* Pointer to the vnode whose fileset is desired.  
*avolpp* Pointer to place to put fileset pointer.

**DESCRIPTION**

If the **vnode**'s fileset is in the *Fileset Registry*, put a pointer to the fileset's entry in the registry in the location indicated by *avolpp*. Otherwise, return an error code.

**DISCUSSION**

This operation extracts a fileset ID from the **vnode** and uses it to obtain and return the volume structure for the fileset on which it resides.

**Note:** This implies that LFS vnodes either contain or point to a fileset ID.

The actual fileset structure (**struct volume**) should be obtained by an up-level call to the *Fileset Registry* routine for that purpose, *volreg\_Lookup()* defined in Section 15.15 on page 399.

**RETURN VALUE**

This function will return the value that *volreg\_Lookup()* returned to it. If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[ENODEV] The object does not reside on an exported fileset (there is no fileset for the **vnode**).

**SEE ALSO**

*volreg\_Lookup()* defined in Section 15.15 on page 399.

**NAME**

vn\_afsfid — Return file ID suitable for exporter

**SYNOPSIS**

```
int vn_afsfid(
    /* IN */ struct vnode *avp,
    /* OUT */ struct afsFid *af,
    /* IN */ int wantv
);
```

**ARGUMENTS**

*avp* Pointer to the vnode whose file ID is desired.

*aafsfidp* Pointer to the file ID structure as used by the Protocol Exporter.

*wantv* Flag indicating whether the *volume* and *cell* fields in the file ID structure should be filled in.

**DESCRIPTION**

*vn\_fid()* and *vn\_afsfid()* have to be separate **X-ops**. Given a pointer to a vnode *avp*, set *af* to point to a file ID pointer describing the file ID associated with the vnode.

This **X-op** is called only on behalf of Protocol Exporter.

*vn\_fid()* is called by the NFS server.

**DISCUSSION**

Fields within the returned **afsFid** are constructed as follows:

*.Vnode* Set to the LFS-dependent (non-canonical) anode number.

*.Unique* Set to the anode generation number.

*.Volume* If *wantv* is non-zero, set to the fileset ID of the fileset on which this object resides.

*.Cell* If *wantv* is non-zero, the *.high* and *.low* subfields should be set to 1 and 0, respectively.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[error\_status\_ok] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

*vn\_fid()* which does a similar function.



**NAME**

vn\_getacl — Get the vnode's access control list

**SYNOPSIS**

```
int vn_getacl(
    /* IN */ struct vnode *avp,
    /* OUT */ struct dfs_acl *aaclp,
    /* IN */ long aw,
    /* IN */ struct ucred *acred
);
```

**ARGUMENTS**

*avp* Pointer to the vnode whose **ACL** is desired.

*aaclp* Pointer to the buffer to which the contents and length of the chosen **ACL** are written.

*aw* Specifies which of a file's **ACLs** is desired.

*acred* Pointer to the caller's credential structure.

**DESCRIPTION**

Read the **ACL** associated with the file whose **vnode** pointer is *avp*. A pointer to the caller's credentials is contained in *acred*. Argument *aw* describes which of a file's **ACLs** is desired. For regular files, there is only one possibility, but for directories, either the directory's own **ACL**, or the default **ACL** for its subdirectories, or the default **ACL** for its subfiles may be specified.

The format of the **ACL** is described in the topic, *ACL Storage Format*, Chapter 9 on page 165.

**DISCUSSION**

The **ACL** designated by the *aw* argument is encoded as described in the above mentioned document and returned. The *aw* argument can take the following values.

**VNX\_ACL\_REGULAR\_ACL**

The **ACL** for the object itself is returned.

**VNX\_ACL\_DEFAULT\_ACL**

The *Initial Container ACL* for a directory object is returned. If the object in question is not a directory, an [EINVAL] error is returned.

**VNX\_ACL\_INITIAL\_ACL**

The *Initial File ACL* for a the object in question is not a directory; an [EINVAL] error is returned. If the specified **ACL** exists for the object in question, it is returned. If *aw* equals **VNX\_ACL\_REGULAR\_ACL**, it may be necessary to coerce the "required" (user\_obj, group\_obj, other\_obj) entries to agree with the mode bits on the file. Consult the "Relationship between an **ACL** and the UNIX protection mode bits" section of the LFS **ACL** specification.

If the requested **ACL** does not exist, matters are somewhat more complicated. A minimal **ACL** is constructed as is described in *ACL Creation From Mode Bits Algorithm*, Section 12.10 on page 257. The actual specifics differ slightly depending on the type of **ACL** requested (the *aw* argument):

**VNX\_ACL\_REGULAR\_ACL**

The mode bits used in the **ACL** synthesis are taken from the object's mode.

**VNX\_ACL\_DEFAULT\_ACL**

The default is used.

**VNX\_ACL\_INITIAL\_ACL**

These initial **ACLs** are synthesized from mode bits of 0777 (**Read-Write-Execute** for all).

If the object in question does possess a regular **ACL**, the default realm in the constructed initial **ACL** is set to the value from that regular **ACL** instead of the LOCAL-REALM (local cell, that is).

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[EINVAL] Returned if the *aw* argument is not one of the values VNX\_ACL\_REGULAR\_ACL, VNX\_ACL\_DEFAULT\_ACL or VNX\_ACL\_INITIAL\_ACL.

Also returned if *aw* is either VNX\_ACL\_DEFAULT\_ACL or VNX\_ACL\_INITIAL\_ACL and the object (*avp*) is not a directory.

**SEE ALSO**

*ACL Storage Format*, Chapter 9 on page 165, and *ACL Creation From Mode Bits Algorithm*, Section 12.10 on page 257.

**NAME**

vn\_setacl — Set the vnode's access control list

**SYNOPSIS**

```
int vn_setacl(
    /* IN */ struct vnode *avp,
    /* IN */ char *aacpl,
    /* IN */ struct vnode *asvp,
    /* IN */ long adestw,
    /* IN */ long asrcw,
    /* IN */ struct ucred *acred
);
```

**ARGUMENTS**

<i>avp</i>	Pointer to the vnode whose <b>ACL</b> is to be overwritten.
<i>aacpl</i>	Pointer to the buffer from which to read the new length and contents of the chosen <b>ACL</b> .
<i>asvp</i>	Pointer to the vnode structure (if any) to use as the “shared” <b>ACL</b> object.
<i>adestw</i>	Indication of which of the file's <b>ACLs</b> is to be set.
<i>asrcw</i>	Indication of which of the source file's <b>ACLs</b> is to be copied.
<i>acred</i>	Pointer to the caller's credential structure.

**DESCRIPTION**

Set the contents of the **ACL** associated with the file whose **vnode** pointer is *avp*. If *aacpl* is non-null, write the **ACL** directly from that structure. If *aacpl* is null and *asvp* is non-null, then copy the **ACL** of *avp* from the **ACL** associated with the file whose **vnode** pointer is *asvp*. This copying gives the underlying file system the opportunity to detect sharing of **ACLs**. However, the **ACLs** are not *semantically* shared; if a subsequent call to *vn\_setacl()* modifies the **ACL** of *avp*, the **ACL** of *asvp* is not modified with it. Exactly one of *adestw* and *asrcw* must be non-zero. A pointer to the caller's credentials is passed via *acred*. The format of the information in the **ACL** is identical to that in *vn\_getacl()*

**DISCUSSION**

This operation sets an **ACL** on the target object (*avp*). The *adestw* argument indicates which type of **ACL** is being set and can take the following values:

VNX\_ACL\_REGULAR\_ACL

The **ACL** on the object itself is set.

VNX\_ACL\_DEFAULT\_ACL

The initial directory **ACL** on a directory object is set.

VNX\_ACL\_INITIAL\_ACL

The initial file **ACL** on a directory object is set.

The **ACL** to be applied to the target is specified by either the *asvp* or *aacpl* arguments, exactly 1 of which should be non-NULL:

- **If *aacpl* is non-NULL:**

The **ACL** has been supplied as an argument to this call. It's representation is described in the above mentioned LFS **ACL** specification.

If the supplied **ACL** is "empty" (*dfs\_acl\_len* == 0), any existing target **ACL** is removed.

If the supplied **ACL** is present (*.dfs\_acl\_len* > 0) and the *destw* argument equals `VNX_ACL_REGULAR_ACL`, it may be necessary to update the mode bits of the object.

The *srcw* argument is unused.

- **If *asvp* is non-NULL:**

The **ACL** to be applied to the target is copied from the source object indicated by *asvp*. The *srcw* argument indicates which **ACL** on the source is to be used and will have one of the `VNX_ACL_XXX` values listed above. There is no requirement that the source and target **ACL** types agree.

**Note:** If either *srcw* or *destw* equal `VNX_ACL_REGULAR_ACL`, the internal Episode implementation of **a**CLs makes itself aparent here.

Normally, **ACL** operations under Episode behave as follows. When one sets an **ACL**, the file's mode bits are changed to agree with the required **ACL** entries. When a file's mode bits are changed (**chmod**), its **ACL** is not physically changed on disk. Instead, the mode bits are "merged" into the **ACL** dynamically during **ACL** checks and when **ACL**s are returned to user space.

If *srcw* equals `VNX_ACL_REGULAR_ACL`, the *vn\_setacl()* operation fetches the source **ACL** as it resides on disk, without coercing the **ACL**'s required entries to agree with the file's mode bits. If *destw* equals `VNX_ACL_REGULAR_ACL`, the **ACL** is applied to the target object without changing the file's mode bits.

Finally, the *ctime* on the target object is set to the current time.

#### RETURN VALUE

If this function succeeds, it returns a value of zero. This function succeeds if:

[*error\_status\_ok*] This function returns success.

#### ERRORS

- [EINVAL] Returned if both *asvp* and *aaclp* are non-NULL.
- Returned if both *asvp* and *aaclp* are NULL.
- Returned if *adestw* is not equal to one of the legal values `VNX_ACL_XXX` listed above.
- Returned if *asvp* is non-NULL and *asrcw* is not one of the legal `VNX_ACL_XXX` values listed above.
- Returned if *adestw* is either `VNX_ACL_DEFAULT_ACL` or `VNX_ACL_INITIAL_ACL` and the target object (*avp*) is not a directory.
- Returned if *asrcw* is either `VNX_ACL_DEFAULT_ACL` or `VNX_ACL_INITIAL_ACL`, *asvp* is non-NULL and the source object is not a directory.
- [EINVAL] If *aaclp* is non-NULL and the supplied **ACL** is in an illegal or inconsistent state. Specifically, a valid **ACL** requires that:
  - The **ACL**'s length is consistent with its actual contents.
  - The *.mgr\_type\_field* is correct.
  - The required entries are all present.
  - The *entry\_type\_foreign\_user* and *entry\_type\_foreign\_group* entries must contain a realm UUID that differs from the default one.

- The *user\_obj* entry grants **perm\_control** rights.

[EPERM] Returned if The caller does not have administrative rights to the target object (*avp*).

**SEE ALSO**

See *vn\_getacl()*, **ACL Storage Format**, Chapter 9 on page 165, and **ACL Creation From Mode Bits Algorithm**, Section 12.10 on page 257.



# DCE DFS VFS+ Extended Credential Package

## 17.1 The xcred Package

This chapter describes the interface used for authentication within the VFS+ layer. It presents an extensible and upwards-compatible interpretation of the VFS UNIX **credential** structure. Section 17.2 provides a short overview of this extended credential package, with Section 17.3 on page 470 specifying the exported interface. Important data structures are detailed in *Extended Credential Data Types*, Section 13.4 on page 289. The set of reserved attributes for the user's property list is described in Section 13.4.3 on page 290.

## 17.2 Package Overview

This chapter, describing VFS client context extensions, does not describe new VFS interface functions. Rather, it describes what is essentially a subroutine library allowing VFS functions to interpret the VFS credentials structure in a more general manner.

Specifically, different file systems may have different authentication requirements, and one system may not be satisfied with the authentication guarantees made by a particular file exporter.

### 17.2.1 Interface Overview

To give an idea of the problems must be solved here, consider some of the possible authenticated identities a DCE fileset might have to evaluate. First, a call might enter via a authenticated RPC call, in which case the UUID would be used as the user's identity. A call might also enter via a kernel-uthenticated access on the File Exporter itself. In this mode of authentication, the user is identified by his or her UNIX UID to the DCE fileset. A third user accessing the file system might come in via an essentially unauthenticated NFS call, also specifying a UNIX UID for the accessing user. However, this access path, while it results in a UNIX UID identifying the user, probably should be treated as less secure than a kernel-authenticated access.

So, a client context is created and associated with a typed authentication structure. There are various calls in the interface to access these objects. One call might ask for the string name for the authenticated user. Another might ask for a 32-bit UID from the authentication library (and the interface might call the protection server to map a name to an ID if necessary to provide the information). Since the authentication level for a 32-bit UID might differ considerably (compare getting that UID from the kernel versus getting it from an unprotected NFS RPC call), some hint as to what level of authentication is permitted needs to be provided.

### **17.3 xcred Functions**

This section describes the functions exported by the credential package.



**NAME**

xcred\_Init — Initialize the module

**SYNOPSIS**

```
long xcred_Init(  
    /* No parameters */  
);
```

**ARGUMENTS**

No parameters.

**DESCRIPTION**

This function initializes the internal state of the xcred module. It must be called at least once before any of the other exported routines are invoked. Failure to do so will cause unpredictable results from the xcred module. Only the first call will actually perform initialization, all others will be no-ops.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[XCRED\_EC\_ALLOC\_FAILED]  
There was a memory allocation failure.

**NAME**

`xcred_Create` — Create an extended credential structure

**SYNOPSIS**

```
long xcred_Create(  
    /* OUT */ xcred_t    **anewXCredPP  
);
```

**ARGUMENTS**

*newXCredPP*      Address of xcred pointer to fill.

**DESCRIPTION**

This function returns the location of a new (empty) xcred structure.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[*error\_status\_ok*] This function returns success.

**ERRORS**

[*XCRED\_EC\_ALLOC\_FAILED*]  
                    There was a memory allocation failure.

**NAME**

xcred\_Hold — Increment the reference count on a given xcred structure

**SYNOPSIS**

```
long xcred_Hold(  
    /* IN */ xcred_t  *axcredP  
);
```

**ARGUMENTS**

*axcredP*            Pointer to the xcred structure to update.

**DESCRIPTION**

This function simply increments the reference count on the relevant structure.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[XCRED\_EC\_INVALID\_PARAM]  
The xcred was invalid.

[XXCRED\_EC\_BAD\_REFCOUNT]  
The refcount was already less than or equal to zero.

**NAME**

`xcred_Delete` — Mark an `xcred` structure for deletion

**SYNOPSIS**

```
long xcred_Delete(  
    /* IN */ xcred_t *axP  
);
```

**ARGUMENTS**

`axP` Pointer to the `xcred` structure to delete.

**DESCRIPTION**

This function marks the `xcred` structure for deletion. This deletion will occur when the `xcred_Release()` function is called with this `xcred` structure as an argument.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[`error_status_ok`] This function always returns success.

**ERRORS**

None.

**NAME**

xcred\_Release — Decrement the reference count on a given xcred structure

**SYNOPSIS**

```
long xcred_Release(  
    /* IN */ xcred_t *axcredP  
);
```

**ARGUMENTS**

*axcredP* Pointer to the xcred structure to update.

**DESCRIPTION**

This function decrements the reference count on the relevant structure. If the reference count drops to zero, or a call to *xcred\_Delete* was previously issued against the xcred structure, then the structure is freed.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[XCRED\_EC\_INVALID\_PARAM]  
The xcred was invalid.

[XCRED\_EC\_BAD\_REFCOUNT]  
The refcount was already less than or equal to zero.

**NAME**

`xcred_AssociateCreds` — Associate the given `xcred` and UNIX `ucred`

**SYNOPSIS**

```
long xcred_AssociateCreds(  
    /* IN */ xcred_t      *axcredP,  
    /* OUT */ struct ucred **aucredPP,  
    /* IN */ long         apag  
);
```

**ARGUMENTS**

*axcredP*            Pointer to the `xcred` structure for the association.  
*aucredPP*           Pointer to the related UNIX credential structure.  
*apag*                `pag` to allocate for linkage, if any.

**DESCRIPTION**

This function associates the given `xcred` and `ucred` structures so that having a pointer to the `ucred` structure will permit the location of the `xcred` structure with which it is associated. If *apag* is not-specified, (equal to zero), this function will generate one.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[`error_status_ok`] This function returns success.

**ERRORS**

[`CRED_EC_ALLOC_FAILED` | `RED_EC_INVALID_PARAM`]  
An invalid structure address was passed in.  
[`XCRED_EC_BAD_REFCOUNT`]  
The refcount was already less than or equal to zero.  
[`XCRED_EC_CRED_FULL`]  
The UNIX cred can't handle additional xcreds.  
[`OSI_NOPAG`] The *apag* value is invalid.

**SEE ALSO**

`xcred_UCredToXCred( )`.

**NAME**

xcred\_UCredToXCred — Obtain pointer to ucred associated with the given xcred

**SYNOPSIS**

```
long xcred_UCredToXCred(  
    /* IN */ struct ucred *aucredP,  
    /* OUT */ xcred_t **axcredPP  
);
```

**ARGUMENTS**

*aucredP* Pointer to the UNIX credential structure to examine.

*axcredPP* Pointer to the xcred structure for the association.

**DESCRIPTION**

This function obtains the corresponding xcred structure from the given UNIX cred, increments its reference count by 1, and returns its location. If an xcred structure has not been associated with the ucred, *axcredPP* will be set to null and the appropriate error code is returned.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[XCRED\_EC\_INVALID\_PARAM]  
An invalid structure address was passed in.

[XCRED\_EC\_BAD\_REFCOUNT]  
The refcount on the xcred was already less than or equal zero.

[XCRED\_EC\_NO\_ENTRY]  
The given UNIX cred had no associated xcred.

**NAME**

xcred\_FindByPag — Obtain pointer to xcred

**SYNOPSIS**

```
long xcred_FindByPag(  
    /* IN */ long apag  
);
```

**ARGUMENTS**

*apag* pag to match against.

**DESCRIPTION**

This function returns a pointer to the xcred structure that contains the pag equal to *apag*. If no match is found then NULL is returned.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[non\_NULL] This function returns a non\_NULL pointer to the xcred structure upon success.

**ERRORS**

[NULL\_Pointer] No match was found. This is a FAILURE.



**NAME**

xcred\_PutProp — Enter a property list attribute-value pair into an xcred

**SYNOPSIS**

```
long xcred_PutProp(
    /* IN */ xcred_t      *axcredP,
    /* IN */ char         *aattributeP,
    /* IN */ long        aattrLength,
    /* IN */ char         *avalueP,
    /* IN */ long        aLength,
    /* IN */ char         *abaseAttrP,
    /* IN */ long        abaseAttrLength
);
```

**ARGUMENTS**

*axcredP* Pointer to the xcred structure to modify.

*aattributeP* Name of the property list entry.

*aattrLength* Number of bytes pointed to by *aattributeP*.

*avalueP* A counted array of bytes comprising the value portion.

*alength* The length of the *avalueP* field in bytes.

*abaseAttrP* Name of the attribute from which this plist entry is derived.  
**Note:** This parameter is now ignored.

*abaseAttrLength* Number of bytes pointed to by *abaseAttrP*.  
**Note:** This parameter is now ignored.

**DESCRIPTION**

This function modifies the contents of the specified attribute associated with the given xcred. If *avalueP* is null, the attribute is deleted from the property list. The *avalueP* parameter's type is a counted array of bytes, not a character string, despite the declaration, and thus may contain null characters that would otherwise terminate a C language character string. This is also true of *aattributeP*.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[XCRED\_EC\_INVALID\_PARAM]  
 An invalid parameter was passed in.

[XCRED\_EC\_BAD\_REFCOUNT]  
 The refcount on the xcred was already less than or equal zero.

[XCRED\_EC\_NO\_BASE\_ATTRIBUTE]  
 The base attribute specified doesn't exist.

**NAME**

`xcred_GetProp` — Retrieve attribute value from an `xcred`'s property list

**SYNOPSIS**

```
long xcred_GetProp(  
    /* IN */ xcred_t      *axcredP,  
    /* IN */ char         *aattributeP,  
    /* IN */ long         aattrLength,  
    /* OUT */ char        **avaluePP,  
    /* OUT */ long        *arealLengthP  
);
```

**ARGUMENTS**

*axcredP*            Pointer to the `xcred` structure to examine.

*aattributeP*        Name of the property to retrieve.

*aattrLength*        Number of bytes pointed to by *aattributeP*.

*avalueP*            Pointer to counted array of bytes in which to store the retrieved value, if any (may be null).

*arealLengthP*       The number of bytes placed into *avaluePP*.

**DESCRIPTION**

This function returns the value from the specified attribute for an `xcred`. The value is placed in *avaluePP*, and its length is placed in *arealLengthP*.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[`error_status_ok`] This function returns success.

**ERRORS**

[`XCRED_EC_BAD_REFCOUNT`]  
The reference count on the `xcred` was less than or equal zero.

[`XCRED_EC_NO_ENTRY`]  
There was no entry named attribute associated with the `xcred`.

[`XCRED_EC_NOT_ENOUGH_ROOM`]  
There wasn't enough room in the caller's buffer to copy out the value component.

[`XCRED_EC_ALLOC_FAILED`]  
There was a memory allocation failure.

**NAME**

xcred\_GetUFlags — Get user flags from an xcred

**SYNOPSIS**

```
long xcred_GetUFlags(  
    /* IN */ xcred_t *xcredP  
);
```

**ARGUMENTS**

*xcredP*                    Pointer to the xcred structure from which the flags are to be obtained.

**DESCRIPTION**

This function returns the value of the user defined flags for the given xcred structure. Currently this is the value of the uflags field.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always returns success. This function succeeds if:

[value]                    This function returns success. The value of the uflags field is returned (here).

**ERRORS**

None.

**NAME**

`xcred_SetUFlags` — Set user flags in an `xcred`

**SYNOPSIS**

```
long xcred_SetUFlags(  
    /* IN */ xcred_t *xcredP,  
    /* IN */ long orflags,  
    /* IN */ long andflags  
);
```

**ARGUMENTS**

<i>xcredP</i>	Pointer to the <code>xcred</code> structure from which the flags are to be set.
<i>orflags</i>	Flags that are to be set (turned on).
<i>andflags</i>	Flags that are to be unset (turned off).

**DESCRIPTION**

This function sets the value of the user defined flags in the `xcred` structure. Currently this is the value of the `uflags` field. Those flag bits set in the *orflags* parameter are set, and the flag bits set in the *andflags* parameter are unset.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

None.

**NAME**

xcred\_EnumerateProp — Apply a function to all elements of the given xcred's property list

**SYNOPSIS**

```
long xcred_EnumerateProp(
    /* IN */ xcred_t *axcredP,
    /* IN */ long (*aprockP)(),
    /* IN */ long *arockP
);
```

**ARGUMENTS**

*axcredP* Pointer to the xcred structure being operated on.

*aprockP* Ptr to function to call on each member of the xcred's property list.

*arockP* The first parameter passed to *\*aprockP()*. The given routine is actually given three parameters, the second being *axcredP* and the third a pointer to the parameter list entry currently being operated upon.

**DESCRIPTION**

This call applies the function given in *aprockP* to all the elements of an xcred's property list. The function will be passed three arguments, the first is *arockP*, the second is *axcredP*, and last is a pointer to the plist entry currently being operated upon. The function should return a long value which is 0 for success, and any non-zero value to indicate failure. Note that the xcred is read-locked for the duration of the call, implying that the function being applied should not try to lock the xcred itself.

This function will halt at the first unsuccessful invocation of the *aprockP* function, and a non-zero value will be returned in this case.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function succeeds if:

[error\_status\_ok] This function returns success.

**ERRORS**

[XCRED\_EC\_INVALID\_PARAM]  
A parameter was illegal.

[XCRED\_EC\_BAD\_REFCOUNT]  
The xcred's reference count was less than or equal zero.

[any\_other\_non\_zero]  
This function did not complete successfully. The function was not applied to all elements of the xcred's property list.

**NAME**

`xcred_DeleteEntry` — Mark a property list entry as deleted

**SYNOPSIS**

```
long xcred_DeleteEntry(  
    /* IN */ xcred_t          *axcredP,  
    /* IN */ struct xcred_PlistEntry *aplp  
);
```

**ARGUMENTS**

*axcredP*            Pointer to the xcred structure being operated on.  
*aplp*                The property list entry to mark as deleted.

**DESCRIPTION**

This function marks a property list entry as deleted. This function is only to be called from within a function that is invoked via the *xcred\_EnumerateProp()* function (the only legal source for *xcred\_PListEntry* objects). It assumes this is the case and does not lock the xcred structure.

**RETURN VALUE**

If this function succeeds, it returns a value of zero. This function always succeeds. This function succeeds if:

[*error\_status\_ok*] This function always returns success.

**ERRORS**

None.

**SEE ALSO**

Section 13.4.1 on page 289.

# Index

ACL .....	6, 413	AFS_HardLink	
ACL Format		IDL Definition .....	34
Access Type .....	166	AFS_HardLink() .....	<b>126</b>
ACL Entry Types .....	167	AFS_Lookup	
ACL Structure .....	169	IDL Definition .....	38
Complex Entry Type .....	168	AFS_Lookup() .....	<b>116</b>
Extended Complex Entry Type .....	169	AFS_LookupRoot	
Foreign Cell Principal ID .....	166	IDL Definition .....	23
Principal ID .....	165	AFS_LookupRoot() .....	<b>104</b>
Simple Entry Type .....	168	AFS_MakeDir	
AFS4Int		IDL Definition .....	35
IDL Definitions .....	22-48	AFS_MakeDir() .....	<b>128</b>
AFS4Int Interface		AFS_MakeMountPoint	
End .....	49	IDL Definition .....	42
Start .....	13	AFS_MakeMountPoint() .....	<b>124</b>
AFS_BulkFetchVV		AFS_ProcessQuota	
IDL Definition .....	44	IDL Definition .....	46
AFS_BulkFetchVV() .....	<b>139</b>	AFS_ProcessQuota() .....	<b>140</b>
AFS_BulkKeepAlive		AFS_Readdir	
IDL Definition .....	45	IDL Definition .....	37
AFS_BulkKeepAlive() .....	<b>141</b>	AFS_Readdir() .....	<b>132</b>
AFS_CreateFile		AFS_ReleaseTokens	
IDL Definition .....	31	IDL Definition .....	40
AFS_CreateFile() .....	<b>118</b>	AFS_ReleaseTokens() .....	<b>137</b>
AFS_FetchACL		AFS_RemoveDir	
IDL Definition .....	25	IDL Definition .....	36
AFS_FetchACL() .....	<b>107</b>	AFS_RemoveDir() .....	<b>130</b>
AFS_FetchData		AFS_RemoveFile	
IDL Definition .....	24	IDL Definition .....	30
AFS_FetchData() .....	<b>105</b>	AFS_RemoveFile() .....	<b>114</b>
AFS_FetchStatus		AFS_Rename	
IDL Definition .....	26	IDL Definition .....	32
AFS_FetchStatus() .....	<b>108</b>	AFS_Rename() .....	<b>120</b>
AFS_GetServerInterfaces		AFS_SetContext	
IDL Definition .....	47	IDL Definition .....	22
AFS_GetServerInterfaces() .....	<b>143</b>	AFS_SetContext() .....	<b>103</b>
AFS_GetStatistics		AFS_SetParams	
IDL Definition .....	43	IDL Definition .....	48
AFS_GetStatistics() .....	<b>136</b>	AFS_SetParams() .....	<b>142</b>
AFS_GetTime		AFS_StoreACL	
IDL Definition .....	41	IDL Definition .....	28
AFS_GetTime Constants .....	14	AFS_StoreACL() .....	<b>111</b>
AFS_GetTime() .....	<b>138</b>	AFS_StoreData	
AFS_GetToken		IDL Definition .....	27
IDL Definition .....	39	AFS_StoreData() .....	<b>109</b>
AFS_GetToken() .....	<b>134</b>		

AFS_StoreStatus	
IDL Definition .....	29
AFS_StoreStatus Mask Values .....	14
AFS_StoreStatus() .....	113
AFS_Symlink	
IDL Definition .....	33
AFS_Symlink() .....	122
Aggregate .....	6
aggregate .....	236
Aggregate Operations Vector .....	240
Aggregate Registry .....	236, 303
Aggregate Structure .....	323
Aggregate Type .....	240
ag_attach() .....	314
ag_detach() .....	313
ag_hold() .....	306
ag_lock() .....	308
ag_rele() .....	307
ag_setops() .....	240
ag_stat() .....	310
ag_sync() .....	315
ag_unlock() .....	309
ag_volCreate() .....	311
ag_volInfo() .....	312
Algorithms	
Access Check .....	161
Access Check Delegation .....	162
Access Rights .....	257, 262-263, 265
ACL Creation .....	257
ACL Exists .....	263
File Creation .....	258
Initial ACL .....	258
mask_obj ACL Impact .....	265
Mode Bits .....	257
No ACL Exists .....	263
PAC From Ucred .....	261
Principals Access .....	161
Principals Identity .....	259
Anode .....	6
Anode Index .....	6
Backing Anode .....	6
Block .....	6
callback .....	234, 236
Cell .....	6
Clone .....	6
Common_data Interface	
End .....	65
Start .....	57
Constants	
General AFS constants .....	57
afsConnParams Mask .....	62
AFS_GetTime .....	14
for afsFStype .....	58
for afsVolumeType .....	58
for AFS_Mount .....	58
for Cell and Hosts .....	57
for RPC Versioning Scheme .....	64
for tn_tag .....	59
General Client Constants .....	76
General Server Constants .....	84
General UPDATE constants .....	73
IS_COMM_ERR .....	84
List of Supported Interfaces .....	76
Quota Opcodes .....	58
Quota Types .....	58
Container .....	6
COW Anode .....	6
credential .....	234, 237
dacl_AclMgrName() .....	212
dacl_AddEntryToAcl() .....	195
dacl_AreObjectEntriesRequired() .....	209
dacl_AreObjectUidsRequiredOnAccessCheck() .....	210
dacl_ArePermBitsRequiredOnAccessCheck() .....	211
dacl_CheckAccessAllowedPac() .....	178
dacl_CheckAccessId() .....	172
dacl_CheckAccessPac() .....	176
dacl_ChmodAcl() .....	188
dacl_CreateAclOnDisk() .....	193
dacl_DeleteAclEntry() .....	197
dacl_DeleteAllEntries() .....	198
dacl_DetermineAccessAllowed() .....	174
dacl_EntryType_FromString() .....	204
dacl_EntryType_ToString() .....	203
dacl_epi_CheckAccessAllowedPac() .....	179
dacl_epi_CheckAccessPac() .....	177
dacl_epi_FlattenAcl() .....	183
dacl_ExtractPermBits() .....	187
dacl_FlattenAcl() .....	182
dacl_FlattenAclWithModeBits() .....	181
dacl_FreeAclEntries() .....	189
dacl_InitAclEntryFromStrings() .....	201
dacl_InitEpiAcl() .....	208
dacl_ModifyAclEntry() .....	196
dacl_NameAndTypeStringsFromEntry() .....	202
dacl_PacFromUcred() .....	180
dacl_ParseAcl() .....	186
dacl_ParseAclDiskOption() .....	184
dacl_ParseSyscallAcl() .....	185
dacl_Permset_FromString() .....	206
dacl_Permset_ToString() .....	205
dacl_PrintAcl() .....	191
dacl_PrintAclEntry() .....	190



## Index

dacl_ReadFromDisk()	194	ag_status_st	269
dacl_ValidateBuffer()	207	ag_sync() Flags	271
dacl_WriteToDisk()	192	asfBulkFEX	20
Data Types	16, 59, 64, 74, 78, 88	astab	273
Access Permissions	166	as_type	274
ACL	165-170, 254	BulkKeepAlive afsFidExp	20
ACL Formats	165	BulkKeepAlive re-check	20
ACL Permission Sets	166	Client	78
afsACL	20	Client Globals	78
afsBulkStats	21	Concurrency Levels	273
afsBulkVolIDs	21	dacl	169
afsBulkVVs	20	dacl_complex_entry	168
afsConnParams	62	dacl_complex_entry_type	168
afsDBCachEntry	62	dacl_entry	168
afsDBLock	63	dacl_entry_type_t	167
afsDBLockDesc	62	dacl_extended_info	169
afsFetchStatus	17	dacl_format_label	169
afsFid	60	dacl_permset_t	166
afsFids	63	dacl_simple_entry	168
afsFidTaggedName	60	dacl_simple_entry_type	168
afsHyper	60	dfstab entries	273
afsNetAddr	60	dfstab File System Type	274
afsNetData	16	dfs_interfaceList	64
afsQuota	20	Directory Entry	301
afsRecordLock	61	dirent	301
afsReturnDesc	62	Enhanced Operations	299
afsReturns	63	Enhanced Vnode Functions	298
afsRevokeDesc	61	epi_sec_id	170
afsRevokes	63	epi_sec_id_foreign	170
afsStatistics	18	epi_uuid	165
afsStoreStatus	17	epi_uuid_foreign	166
afsStrings	64	Exported Registry Items	272
afsTaggedName	59	External ACL	254
afsTaggedPath	60	File System Type Macros	274
afsTimeval	60	Fileset Functions	274
afsToken	61	Fileset NextHole	277
afsTokens	63	Flags	271
afsVolSync	16	for RPC Versioning Scheme	64
aggr	270	General	59
Aggregate Dynamic Status	269	General, for AFS Data Structures	59
Aggregate Fields	271	Locks	240
Aggregate Operations Vector	272	lock_data	240
Aggregate States	271	Manager	88
Aggregate Static Status	269	Manager Error Exit	88
Aggregate Status	270	Manager Locking	88
Aggregate Status Types	270	Manager Unlocking	88
Aggregates	270	NextHole	277
aggrops	272	Physical File System	18
ag_attach() Flags	271	Registry	272
ag_status	270	Txvattr	294
ag_status_dy	269	UFS dfstab entry	274

UPDATE .....	74
updateFileStatS .....	74
vfsops .....	300
Vnode Functions .....	296
Vnode Ops .....	273
volCreate() Flags .....	271
volume .....	285
volumeops .....	274
vol_dirent .....	286
vol_handle .....	288
vol_NextHole .....	277
vol_status .....	282
vol_statusDesc .....	282
vol_stat_dy .....	279
vol_stat_st .....	277
xcred Property List .....	289
xcred Structure .....	290
xcred_PListEntry_t .....	289
xcred_t .....	290
xvfs_attr .....	296
xvfs_vfsops .....	299
xvfs_vnodeops .....	298
xvfs_xops .....	296
Defines	
Macros .....	300
states .....	282
Vnode Operation Classifications .....	300
VOLHOLE_MAX_HOLES .....	277
VOLOP_XXX .....	276
Volume Operations .....	276
VOL_ERR_XXX .....	286
VOL_MAX_XXX .....	289
VOL_OP to VOL_SYS .....	287
VOL_OP_XXX .....	286
VOL_ROOTINO .....	289
VOL_STAT_XXX .....	284
VOL_SYNC_XXX .....	288
VOL_XXX .....	282
VOPX_UPDATE .....	299
VOPX_XXX .....	298
v_* .....	286
DFS .....	7
DFS Locking Functions .....	240
Directory Entry Formats .....	301
EFS .....	7
Episode UUID .....	254
Epi_PrinId_Cmp() .....	215
Epi_PrinId_FromUuid() .....	214
Epi_PrinId_ToUuid() .....	213
External ACL .....	254
f2vfs_getvolume() .....	424
File .....	7
Fileset .....	7
Fileset Dynamic Status .....	279
Fileset Handle .....	288
Fileset Registry .....	236, 239, 325, 327, 352
Fileset Static Status .....	277
Fileset Status .....	282
Fileset Status Description .....	282
Flags	
for afsRevokeDesc flags .....	59
for afsRevokeDesc outFlags .....	59
Flags	
AFS_FetchACL .....	14
AFS_GetToken .....	14
AFS_ReleaseToken .....	14
AFS_SetContext .....	16
AFS_SetParams .....	16
AFS_StoreACL .....	14
for Client-only Attribute .....	16
for Fileset operations .....	15
for Token Recovery .....	15
Getting a Token .....	14
Token Recovery .....	52
Fragment .....	7
glue function .....	246
Group .....	7
hton_epi_principal_id() .....	218
hton_epi_uuid() .....	216
LFS .....	7
locking .....	236, 238
Locking .....	240
Locking Functions .....	240
Locks .....	240
Magic Cookie .....	7
magic cookie .....	237, 294
Masks	
afsStoreStatus .....	14
VOL_SETSTATUS .....	284
VOL_STAT_XXX .....	284
MAX_TOKEN_RELEASE	
in AFS_ReleaseTokens() .....	137
ntoh_epi_principal_id() .....	219
ntoh_epi_uuid() .....	217
on-disk states .....	283
PAC .....	7
Principal .....	8
Privately Attached Storage .....	323
Quota .....	8
Quota Size .....	289
Realm .....	8
Replica .....	8

## Index

Rights .....	8	UPDATE Interface	
Root Anode Index .....	289	End.....	75
sec_acl_FlattenAcl() .....	199	Start .....	73
sec_acl_ParseAcl() .....	200	UPDATE_FetchFile	
States		Exported Procedure .....	90
on-disk .....	283	IDL Definition .....	75
struct aggr .....	323	UPDATE_FetchInfo	
TKN4Int		Exported Procedure .....	91
IDL Definitions .....	52-55	IDL Definition .....	74
TKN4Int Debugging		UPDATE_FetchObjectInfo	
TKN_GetCE Routine .....	54	Exported Procedure .....	93
TKN_GetLock Routine.....	53	IDL Definition .....	75
TKN4Int Interface		UPDATE_GetServerInterfaces	
End.....	55	Exported Procedure .....	89
Recovery Flags.....	52	IDL Definition .....	74
Start .....	51	UPDATE_v4_0_manager_epv	
TKN_AsyncGrant		Entry Point Vector .....	94
IDL Definition .....	55	vfs Operations.....	423
TKN_AsyncGrant() .....	152	VFS Operations Vector .....	300
TKN_GetCE		VFS+ Interface .....	8
IDL Definition .....	54	VFS+ Switch.....	290
TKN_GetCellName		Vnode Preliminary.....	293
IDL Definition .....	53	VNOPS Vector Organization .....	290
TKN_GetCellName() .....	149	vn_access() .....	436
TKN_GetLock		vn_afsfid() .....	462
IDL Definition .....	53	vn_bmap() .....	454
TKN_GetServerInterfaces		vn_bread() .....	456
IDL Definition .....	54	vn_brelse().....	457
TKN_GetServerInterfaces() .....	151	vn_close().....	428
TKN_InitTokenState		vn_create().....	438
IDL Definition .....	52	vn_fid().....	459
TKN_InitTokenState().....	146	vn_fsync() .....	452
TKN_Probe		vn_getacl().....	463
IDL Definition .....	52	vn_getattr().....	432
TKN_Probe().....	147	vn_getvolume() .....	461
TKN_SetParams		vn_inactive() .....	453
IDL Definition .....	54	vn_ioctl().....	430
TKN_SetParams() .....	150	vn_link().....	441
TKN_TokenRevoke		vn_lockctl().....	458
IDL Definition .....	53	vn_lookup().....	437
TKN_TokenRevoke() .....	148	vn_mkdir() .....	444
Token.....	8	vn_open().....	427
token.....	236	vn_rdwr().....	429
Token Manager .....	238, 247-248	vn_readdir() .....	447
Txvattr Structure.....	294	vn_readlink() .....	451
UPDATE		vn_remove().....	440
IDL Definition.....	74-75	vn_rename() .....	442
UFS.....	8	vn_rmdir() .....	446
Uniquifier .....	8	vn_select() .....	431
UPDATE		vn_setacl() .....	465
Exported Procedure.....	89-91, 93-94	vn_setattr() .....	434

vn_strategy()	455	vol_write()	359
vn_symlink()	450	VOPX_UPDATE	299
volreg_Delete()	398	VOPX_XXX	298
volreg_Enter()	397	Xcred	7
volreg_Lookup()	399	xcred_AssociateCreds()	476
Volume Fields	286	xcred_Create()	472
Volume Handle	288	xcred_Delete()	474
volumeops	274	xcred_DeleteEntry()	484
vol_appenddir()	390	xcred_EnumerateProp()	483
vol_attach()	353	xcred_FindByPag()	478
vol_bulksetstatus()	392	xcred_GetProp()	480
vol_clone()	369	xcred_GetUFlags()	481
vol_close()	350	xcred_Hold()	473
vol_concurr()	384	xcred_Init()	471
vol_copyacl()	382	xcred_PListEntry_t	289
vol_create()	357	xcred_PutProp()	479
vol_delete()	361	xcred_Release()	475
vol_deplete()	352	xcred_SetUFlags()	482
vol_destroy()	351	xcred_UCredToXCred()	477
vol_detach()	354	xfvs_attr Structure	296
vol_freedystat()	380	xfvs_vnodeops Structure	298
vol_getacl()	365	xfvs_xops Structure	296
vol_getattr()	362		
vol_getnextholes()	395		
vol_getstatus()	355		
vol_getvv()	378		
vol_getzlc()	394		
vol_hold()	341		
vol_isroot()	377		
vol_lock()	343		
vol_open()	345		
vol_pushstatus()	387		
vol_read()	358		
vol_readdir()	388		
vol_reclone()	371		
vol_rele()	342		
vol_root()	376		
vol_scan()	349		
vol_seek()	347		
vol_setacl()	367		
vol_setattr()	363		
vol_setdystat()	379		
vol_setnewvid()	381		
vol_setstatus()	356		
vol_swapids()	385		
vol_sync()	386		
vol_tell()	348		
vol_truncate()	360		
vol_unclone()	373		
vol_unlock()	344		
vol_vget()	375		