

Motivation for OSD

- Improved device and data sharing
 - Platform-dependent metadata moved to device
 - Systems need only agree on naming
- Improved scalability & security
 - Devices directly handle client requests
 - Object security w/ application-level granularity
 - Finer granularity than LUN-based security
- Improved performance
 - Hints, QoS, Differentiated Services
- Improved storage management
 - Self-managed, policy-driven storage
 - Storage devices more autonomous





Objects

Volumes



DIOCKS





OSD Interface





OSD T10 Overview November 2003



OSD overview

Basic Protocol

- READ
- WRITE
- CREATE
- REMOVE
- GET ATTR
- SET ATTR

- Very Basic
- Space Mgmt

Attributes

- timestamps
- vendor-specific
 - shared, opaque

Security

- Authorization on each request
- Integrity for args & data
- SET KEY 1 shared
- SET MASTER KEY ^J secrets

<u>Groups</u>

- CREATE COLLECTION
- REMOVE COLLECTION
- LIST COLLECTION

Specialized

- APPEND write w/o offset
- CREATE & WRITE save msg
- FLUSH OBJ force to media
- LIST recovery of objects

Management

- FORMAT OSD
- CREATE PARTITION
- REMOVE PARTITION



Read (8805h) – parameters

Bit Byte	7	6	5	4	3	2	1	0
8	(MSB)				ON (8905b)			
9				SERVICE ACT	014 (000011)			(LSB)
10	OPTIONS BYTE							
11	Rese	Reserved GET/SET CDBFMT PRIORITY						
12	(MSB)				61 h	ite		
19		_		PARITION_ID 04 DILS				(LSB)
20	(MSB)				6/	hite		
27				USER_OBJEC	U U	5115		(LSB)
28				Decerved	_			
31				Reserved	by	/te add	ressable	e
32	(MSB)							
39		-		LENGTH				(LSB)
40	(MSB)			STADTING BY		/		
47		-		STARTING BY	TE ADDRESS			(LSB)



OSD T10 Overview November 2003



List (8803h) – parameters

Table 38 — LIST service action

Bit Byte	7	6	5	4	3	2	1	0		
8	(MSB)	(MSB) SERVICE ACTION (8903b)								
9				SERVICE ACT				(LSB)		
10	Reserved									
11	Reserved GET/SET CDBFMT SORT ORDER									
12	(MSB)									
19			PARTITION_ID				only one option –			
20	(MSB)					ascendi	ng obje	ect id 🗌		
27				DZEK_ORIEC	I_ID TAG			(LSB)		
28				Becorved						
31				Reserved						
32	(MSB)		(h. ffar		- ila bla		
39				ALLOCATION	LENGIH	putter s	size ava			
40	(MSB)		C		oontir	unation a				
43				LIST TAG	conur	iuation	101055	comman		





Objects



OSD T10 Overview lovember 2003



Object names

Table 8 — Partition_ID and User_Object_ID value assignments

Partition_ID	User_Object_ID	Description
Oh	Oh	Root object
Oh	1h - FFFF FFFF FFFF FFFFh	Reserved
1h - FFFFh	0h - FFFF FFFF FFFF FFFFh	Reserved
10000h - FFFF FFFF FFFF FFFFh	0h	Partition (assigned by OSD)
10000h - FFFF FFFF FFFF FFFFh	1h - FFFFh	Reserved
10000h - FFFF FFFF FFFF FFFFh	10000h - FFFF FFFF FFFF FFFFh	User object (assigned by OSD)

Partition IDs assigned by device

- primary usage case assumes one manager per partition
- Object IDs assigned by device OR by host
 - collection IDs share namespace with objects





Attributes



OSD T10 Overview November 2003



Table 9 — Object attribute page numbers

Attributes	range for each	Page number	OSD objec	object type wit t attributes pag	h which the ge is associated
	object type	0h - 2FFF FFFFh	User		
		3000 0000h - 5FFF FFFFh	Partiti	on	
	Table 10 - Object	6000 0000h - 8FFF FFFFh	Collec	ction	
		9000 0000h - BFFF FFFFh	Root		
		C000 0000h - FFFF FFFEh	Reser	ved	
Base Northan		FFFF FFFFh	All att	ributes pages	
within range	Description			Number	List
0h - 7Fh	Defined by this standa	rd		Yes	Yes
80h - 7FFFh	Reserved				
8000h - EFFFh	Defined by other standards			Yes	Yes
F000h - FFFFh	Defined by OSD manufacturer product specifications			Yes	Yes
1 0000h - 2FFF FFFFh	Assigned by the OSD logical unit ephemeral			No	Yes
2000 0000h - 2FFF FFFFh	Vendor specific			n/a	n/a

Limited number defined by standard

• length, size, timestamps

Vendor extensions

- opaque for application use only
- shared device-interpreted (impacts behavior)



OSD T10 Overview November 2003

03-394r0

Also used to do device-level params

- security level
- capacity



Table 26 — Page format get and set attributes CDB parameters format

Bit Byte	7	6	5	4	3	2	1	0
				:	Other CDB :	fields		
79				•				
80	(MSB)					which	attrib	
83				GELATIRIDU	ES PAGE	WIIICH	αιιπρ	(LSB)
84	how n	now much buffer						
87	host h	nas availa	ble	GET ATTRIBUTES ALLOCATION LENGTH				(LSB)
88								
91				RETRIEVED A	TRIBUTES OF	FSET		
92	(MSB)							
95				SELATIRIBUT	E PAGE	which	əttrih	(LSB)
96	(MSB)	_				which		
99				SELATIRIBUT	E NOMBER			(LSB)
100	how m	nuch attri	bs		TICNOTU			
103	l am s	ending		SELATIRIBUT	E LENGTH			(LSB)
104		_		OFT ATTOINT				
107				SELATIRIBUT	ES OFFSET			
108					01 000	<u> </u>		
				:	Other CDB	rields		



Object attributes

Table 64 — User Object Information attributes page contents

Attribute Number	Bytes	Attribute	May be set	OSD Provided
0h	40	Page identification	No	Yes
1h	8	Partition_ID	No	Yes
2h	8	User_Object_ID	No	Yes
3h - 7h		Reserved	No	
8h	8	Used capacity size	No	Yes
9h - Bh		Reserved	No	
Ch	variable	Username	Yes	No
Dh - Eh		Reserved length	No	
Fh	8	Object logical length	No	Yes
10h - FFFF FFFFh		Reserved	No	

Table 69 - User Object Resources attributes page contents

Attribute Number	Bytes	Attribute	May be set	OSD Provided
0h	40	Page identification	No	Yes
1h	8	Capacity quota quota	Yes	No
2h - 7h		Reserved	No	
8h	8	Starting byte address of write or append	No	Yes
9h - FFFF FFFFh		Reserved	No	



OSD T10 Overview November 2003

Object attributes (2)

Table 77 - User Object Timestamps attributes page contents

Attribute Number	Bytes	Attribute	May be set	OSD Provided
0h	40	Page identification	No	Yes
1h	6	Created time	No	Yes
2h	6	Attributes accessed time	No	Yes
3h	6	Attributes modified time	No	Yes
4h	6	Data accessed time	No	Yes
5h	6	Data modified time	No	Yes
6h - FFFF FFFFh		Reserved	No	

Table 79 — Collections attributes page contents

Attribute Number	Bytes	Attribute	May be set	OSD Provided
Oh	40	Page identification	No	Yes
1h - FFFF FF00h	0 or 8	Collection pointer	Yes	No
FFFF FF01h - FFFF FFFFh		Reserved	No	

set of collections an object belongs to





Security



OSD T10 Overview November 2003



Read – security

48	- protect arguments		
59	protect arguments	REQUEST INTEGRITY CHECK VALUE	(LSB)
60	nrotoot roplovo	DEQUEST NONCE	
71	protect replays	REQUEST NONCE	
72			
75	protect attributes	DATA-IN INTEGRITY CHECK VALUE OFFSET	
76	and data		
79		DATA-OUT INTEGRITY CHECK VALUE OFFSET	
80		Cat and act attributes parameters (cap 5, 1, 1, 2)	
107		Get and set attributes parameters (see 5.1.1.5)	
108		Capability (see 4 6 4 4 2)	
163		Capability (see 4.6.4.4.5)	





How to get integrity values

Table 81 — Current Command attributes page contents

Attribute Number	Bytes	Attribute	May be set	OSD Provided
0h	40	Page identification	No	Yes
1h		Reserved	No	
2h	8	Created User_Object_ID	No	Yes
3h	12	Response integrity check value	No	Yes
4h - FFFF FFFFh		Reserved	No	

Special attribute to read the integrity value

Bit 7 6 5 4 3 2 1 0 Byte 0 Traditional command or parameter m-1 m Meta data n-1 n Integrity check value n+11

overall 0

structure

E	10		1 1	
-		1.44	7/	
1	-			

OSD T10 Overview November 2003

03-394r0



Table 1 — OSD Data-In Buffer and Data-Out Buffer model

OSD Security – Illustrated



Security levels

Table 13 — Security level threats thwarted

	Threat Thwarted by Security Level							
Threat	0	1	2	3				
Forgery of credential	No	Yes	Yes	Yes				
Alteration of capabilities	No	Yes	Yes	Yes				
Replay of command or status	No	No	Yes	Yes				
Alteration of command or status	No	No	Yes	Yes				
Replay of data	No	No	No	Yes				
Alteration of data	No	No	No	Yes				
Inspection of command, status or data	No	▼ No	No	No				
Level 1 manda to be levered								

Level 1 needs to be layered

Level 3 needs streaming SHA-1



Credentials

Table 14 - Credential format

Bit Byte	7	6	5	4	3	2	1	0
0								
55			Capability (see 4.6.4.4.3)					
56								
73		-		OSD SYSTEM ID				
74	(MSB)			PARTITION_ID Uniquely i			elv ide	ntify
81							t in tim	
82	(MSB)		\sim	object in t				e
87		-		OBJECT CREA	TION TIME			(LSB)
88				Decerved				
91				Reserved				
92	(MSB)							
111				CREDENTIAL INTEGRITY CHECK VALUE			(LSB)	



OSD T10 Overview November 2003



Credential format (2)

Table 15 — Capability format







Bit	Allowed Operation	Bit	Allowed Operation	
APP	APPEND	RD	READ	
CRE	CREATE or CREATE AND WRITE	RMV	REMOVE	
CRE_COL	CREATE COLLECTION	RMV_COL	REMOVE COLLECTION	
CRE_PART	CREATE PARTITION	RMV_PART	REMOVE PARTITION	
FLS_O	FLUSH OBJECT	ST_ATTR	SET ATTRIBUTES	no hit fo
GT_ATTR	GET ATTRIBUTES	ST_KY	SET KEY	
LST	LIST	ST_MKY	SET MASTER KEY	each
LST_COL	LIST COLLECTION	WRT	WRITE O	peratior
OPNCLS	OPEN or CLOSE			

Table 18 – Capabilities Permission Bits

The USER OBJECT DESCRIPTOR TYPE field (see table 19) specifies the format of information that appears in the USER OBJECT DESCRIPTOR field.



User Object Descriptor Type	Description
0h	The USER OBJECT DESCRIPTOR field shall be ignored
1h	Single user object single object
2h - Fh	Reserved multi-object in future







Nonces – replay protection

Table 21 — Request nonce format

Bit Byte	7	6	5	4	3	2	1	0
0	(MSB)	(MSB)						
5			TIMESTAMP					
6	(MSB)							
11			RANDOM NUMBER				(LSB)	

Prevent requests from being captured and re-sent

- at a later point
- overwrite file data
- denial of service

Nonce management

- rough window of time can age old nonces
- must remember future nonces
- under attack change key version; or "cut off" a host via audit nonce





Key hierarchy

Table 22 - OSD secret key hierarchy

Key Name	Key Shared Using	Key Used To	Key Update Frequency				
Keys Shared Between the Security Manager and the OSD Device Server							
Master	SET MASTER KEY service action	Update Drive key	Change of OSD device owner				
Drive	SET KEY service action	Update Partition key	When Partition key may have been compromised (i.e., very infrequently)				
Partition ^a	SET KEY service action	Update Working keys	When Working key updates may have been compromised (i.e., infrequently)				
Working ^a	SET KEY service action	Create Capability keys	When normal key use affords to much chance that the working key might be reverse engineered (i.e., regularly)				
Keys Shared E	Between the Security Mana	ger and the Applicatior	n Client ^c				
Capability ^d	Capability ^d Credentials and mecha- nisms not specified in this standard Secure commands, responses, and data						
^d As dual purpose number, the capability key is different from other keys in the hierarchy. The capability key is the credential integrity check value, meaning that even though the security manager computes its the computation is based on values beyond the security manager's control (e.g., the user object to which the credential allows access). The time interval during which the capability key is very short. While changing the working key used to construct the credential integrity check value invalidates the capability key, the credential may expire long before that and thus invalidate the capability key.							



Backup Slides



OSD T10 Overview November 2003



OSD Status

- History
 - Started with NSIC NASD research 1995-1999
 - Carnegie Mellon, HP, IBM, Quantum, STK, Seagate
 - Seagate led NSIC OSD into SNIA in 1999
- Today
 - Intel & IBM leading SNIA OSD effort
 - EMC, HP, Panasas, Seagate, Veritas involved
 - IBM architecting objects into version 2 of StorageTank
 - Lustre CFS/HP/BlueArc open-source OSD for DoE
 - 1,000 node; 225 TB cluster installed October 2002
 - Panasas shipping OSD-based products today
 - scalable NAS; large-scale systems (300+ devices)





OSD Commands

Table 29 — Commands for OSD devices (part 1 of 2)

	Operation	Service		
Command name	code	action ^a	Туре	Reference
APPEND	7Fh	8807h	М	6.2
CHANGE ALIASES	A4h	0Bh	0	SPC-3
CREATE	7Fh	8802h	М	6.3
CREATE AND WRITE	7Fh	8812h	М	6.4
CREATE COLLECTION	7Fh	8815h	0	6.5
CREATE PARTITION	7Fh	880Bh	М	6.6
FLUSH OBJECT	7Fh	8808h	М	6.7
FORMAT OSD	7Fh	8801h	М	6.8
GET ATTRIBUTES	7Fh	880Eh	М	6.9
INQUIRY	12h		М	SPC-3
LIST	7Fh	8803h	М	6.10
LIST COLLECTION	7Fh	8817h	0	6.11
LOG SELECT	4Ch		0	SPC-3
LOG SENSE	4Dh		0	SPC-3
MODE SELECT(10)	55h		0	SPC-3
MODE SENSE(10)	5Ah		0	SPC-3





OSD Commands (2)

PERSISTENT RESERVE IN	5Eh		М	SPC-3
PERSISTENT RESERVE OUT	5Fh		М	SPC-3
PREVENT ALLOW MEDIUM REMOVAL	1Eh		0	SPC-3
READ	7Fh	8805h	М	6.12
READ BUFFER	3Ch		0	SPC-3
RECEIVE COPY RESULTS	84h		0	SPC-3
RECEIVE DIAGNOSTIC RESULTS	1Ch		0	SPC-3
REMOVE	7Fh	880Ah	М	6.13
REMOVE COLLECTION	7Fh	8816h	0	6.14
REMOVE PARTITION	7Fh	880Ch	М	6.15
REPORT ALIASES	A3h	0Bh	0	SPC-3
REPORT LUNS	A0h		0	SPC-3
REQUEST SENSE	O3h		М	SPC-3





OSD Commands (3)

Table 29 — Commands for OSD devices (part 2 of 2)

Command name	Operation code	Service action ^a	Туре	Reference		
SEND DIAGNOSTIC	1Dh		М	SPC-3		
SET ATTRIBUTES	7Fh	880Fh	М	6.16		
SET KEY	7Fh	8818h	М	6.17		
SET MASTER KEY	7Fh	8819h	М	6.18		
START STOP UNIT	1Bh		0	SBC-2		
TEST UNIT READY	00h		М	SPC-3		
WRITE	7Fh	8806h	М	6.19		
WRITE BUFFER	3Bh		0	SPC-3		
Key: M = Command implementation is mandatory. O = Command implementation is optional. OB = Command implementation is defined in a previous standard						
^a No entry in the service action column means that the SERVICE ACTION field does not apply to the command. OSD service action codes not listed in this table (i.e., 8800h and 8810h-88FFh) are reserved for future standardization.						



