Persistence

- 36 I/O Devices
- 37 Hard Disk Drives
- 38 RAID
- 39 File and Directories
- 40 File System Implementation
- 41 Locality and the Fast File System
- 42 Crash Consistency and Journaling
- 43 Log-structured File Systems
- 44 Flash-based SSD
- 45 Data Integrity and Protection

LFS ISSUES need for a cleaner

- no overwrite means
 - files, dirs, etc. become fragmented
 - parts of the log no long active:
 - all but most-recent versions of inodes
 - data that has been modified
 - out-of-date imap chunks
- cleaner process asynchronously copies live data
 - from *full segments* to clean new segments
 - cleaned segments are empty, can be used again
 - might use this opportunity to segregate by age, activity, etc.
 - segment full of rarely-changing data rarely needs cleaning

444

LFS cleaning costs

- Cleaner
 - read some number of live segments
 - copy live data out into fewer new segments
 - old segments are now free.
- But....write amplification! Let:
 - N : num segments to be cleaned
 - u : percent of these segments that is live
 - write cost (wc): write amplification of each new byte

```
write cost = (#readSegs + #writeLive + #writeNew) / (#writeNew)
= (N + N*u + (N*(1-u)) / (N * (1-u))
= 2/(1-u)
```

- if utilization low, say 10%: wc = 2.22
- if utilization high, say 90%: wc = 20.00

446

LFS cleaner notes

- advantages
 - asynchronous
 - can be done in bulk, i.e. fast
- opportunities
 - older data less likely to be modified than new data
 - can segregate data based on age for cleaner writes
- implemented on bare disk
 - log chunk to be written to disk is many pages long
 - LFS can report consistency check of all blocks back to OS
 - LFS guarantees that pg *i* written correctly before pg *i*+1

Persistence

- 36 I/O Devices
- 37 Hard Disk Drives
- 38 RAID
- 39 File and Directories
- 40 File System Implementation
- 41 Locality and the Fast File System
- 42 Crash Consistency and Journaling
- 43 Log-structured (and other) File Systems
- 44 Flash-based SSD
- 45 Data Integrity and Protection

SSDs

- non-volatile storage
 - we will assume NAND flash, though rapidly evolving
- terminology
 - a flash chip implements one or more banks (or planes)
 - a bank contains some number of (erase) blocks
 - might be 128 KB or 256 KB
 - a block contains some number of pages
 - maybe 4 KB

Block:		()			1				2	2	
Page:	00	01	02	03	04	05	06	07	08	09	10	11
Content:												

448

SSDs operations

- reads
 - any page can be read, same cost
 - very fast, low microseconds
- erase
 - before writing, a page's entire block must be erased
 - slow, milliseconds
 - needs to be done in advance, usually asynchronously
- program (write)
 - entire page written
 - slower, 100's of usec
- tech constantly evolving, but generally costs follow:
 - read << write << erase

450

SSDs

• Pages can be in one of several states:

		iiii	Initial: pages in block are invalid (1)
Erase()	\rightarrow	EEEE	State of pages in block set to erased (E)
Program(0)	\rightarrow	VEEE	<i>Program page 0; state set to valid</i> (\vee)
Program(0)	\rightarrow	error	Cannot re-program page after programming
Program(1)	\rightarrow	VVEE	Program page 1
Erase()	\rightarrow	EEEE	Contents erased; all pages programmable

SSDs example

• Unrealistically small for example. All start as valid:

Page 0	Page 1	Page 2	Page 3
00011000	11001110	0000001	00111111
VALID	VALID	VALID	VALID

• If we want to write page 0, must first erase:

Page 0	Page 1	Page 2	Page 3
11111111	11111111	11111111	11111111
ERASED	ERASED	ERASED	ERASED

• Now we can program page 0:

Page 0	Page 1	Page 2	Page 3
00000011	11111111	11111111	11111111
VALID	ERASED	ERASED	ERASED

• But, but...pages 1-3 are gone....

452

SSDs deets

	Read	Program	Erase
Device	(μs)	_(μs)	(μs)
SLC	25	200-300	1500-2000
MLC	50	600-900	~3000
TLC	~75	~900-1350	~4500

• Reliability

- no head crashes
- erasure causes blocks to wear out
- NANDs leak
 - not good for archival storage

SSDs from flash

- SSD contain
 - some amount of RAM for mapping tables
 - FLASH
 - control logic
- flash translation layer (FTL)
 - maps logical blocks to physical pages
 - handles erasures asynchronously
 - modifies mappings as needed
 - because of erasures (we don't write in place)
 - failures
 - wear leveling
- log-structured...

454

Why not direct mapped?

- Problems if FTL mapping LBA *N* to physical page *N*
 - performance
 - write to *N* requires:
 - read block
 - erase block
 - write block
 - reliability
 - · hot spots in program cause some blocks to fail
 - no wear leveling

SSDs ft/

- log structure
 - in storage device
 - also in file system above
 - keeps mapping table
- Assume:
 - externally a disk w/ 512-byte sectors
 - client is reading/writing 4k blocks
 - SSD has many 16-KB blocks, w/ 4-KB pages

SSDs example

Write al to F	S bl	lock	100	0, a	$2 \rightarrow$	• 10	01, ł	- 1	$\rightarrow 2$	000	, b2	\rightarrow	2001
Block:			0				1			2	2		
Page:	00	01	02	03	04	05	06	07	08	09	10	11	
Content:													
State:	i	i	i	i	i	i	i	i	i	i	i	i	
Block:		()				1			2	2		
Page:	00	01	02	03	04	05	06	07	08	09	10	11	
Content:													
State:	Е	Е	Е	Е	i	i	i	i	i	i	i	i	
Table:	10	00 -	►0										
													_
Block:		()			1				2	2		
Block: Page:	00	01) 02	03	04	1 05	06	07	08	2 09	2 10	11	
Block: Page: Content:	00 a1	01) 02	03	04	1 05	06	07	08	2 09	2 10	11	
Block: Page: Content: State:	00 a1 V	01 E) 02 E	03 E	04 i	1 05 i	06 i	07 i	08 i	2 09 i	2 10 i	11 i	
Block: Page: Content: State:	00 a1 V	01 E) 02 E	03 E	04 i	1 05 i	06 i	07 i	08 i	2 09 i	2 10 i	11 i	
Block: Page: Content: State:	00 a1 V	01 E) 02 E	03 E	04 i	1 05 i	06 i	07 i	08 i	2 09 i	2 10 i	11 i	
Block: Page: Content: State: Table:	00 a1 V	01 E) 02 E ►	03 E 10	04 i	1 05 i	06 i	07 i	08 i	2 09 i 200	2 10 i	11 i	
Block: Page: Content: State: Table: Block:	00 a1 V	01 E) 02 E ►0	03 E 10	04 i	1 05 i ►1	06 i	07 i	08 i	2 09 i 200	2 10 i	11 i	_
Block: Page: Content: State: Table: Block: Page:	00 a1 V 10	01 E D0 -) 02 E ► 0	03 E 10	04 i	1 05 i ►1	06 i 20 06	07 i	08 i ►2	2 09 i 200 200	2 10 i 01→	11 i	_
Block: Page: Content: State: Table: Block: Page: Content:	00 a1 V 10 00 a1	01 E D0 - 01 a2) 02 E ► 0 02 02 b1	03 E 10 03 b2	04 i 01 -	1 i ► 1	06 i 20 06	07 i 00→	08 i ►2	2 09 i 200 2 09	2 10 i 01 →	11 i	_
Block: Page: Content: State: Table: Block: Page: Content: State:	00 a1 V 10 a1 V	(01 E 00 - (01 a2 V) 02 E ► 0 02 b1 V	03 E 10 03 b2 V	04 i 11 -	1 05 i ►1 1 05 i	06 i 20 06 i	07 i 00→ 07	08 i ►2 08	2 09 i 200 2 09 i	2 10 i 01→ 10 i	11 i ►3	-

Rewrite $c1 \rightarrow 100, c1 \rightarrow 101$

Table:	100 →4	101 ->5	2000-2	2001->3
Table.	100 -4	101 - 5	2000 - 2	2001 - 5

Block:		()				1			2	2	
Page:	00	01	02	03	04	05	06	07	08	09	10	11
Content:	a1	a2	b1	b2	c1	c2						
State:	V	۷	۷	۷	V	۷	Е	Е	i	i	i	i

	Ga	ırb	age	со	llec	t						
Table:	10	- 00	►4	10	01 -	►5	20	000-	►6	20	001-	►7
Block:		()				1			2	2	
Page:	00	01	02	03	04	05	06	07	08	09	10	11
Content:					c1	c2	b1	b2				
State:	Е	Е	Е	Е	V	٧	٧	٧	i	i	i	i

logical block (FS) to physical page mapping table

456

SSDS mapping table size assume 1TB drive: if assume 4 bytes / 4k block, 1GB for table, in memory! block approach: reduces size by factor $\frac{size_{block}}{size_{page}}$, but more complicated assume page-level mappings $2000 \rightarrow 4$, $2001 \rightarrow 5$, $2002 \rightarrow 6$, $2003 \rightarrow 7$: all have chunk 500 • offsets 0, 1, 2, 3 LBA translation: get chunk from top-level bits add offset to chunk -> page mapping • Table: **500** → 4 Memory Block: 0 2 1 Page: 00 01 02 03 04 05 06 07 08 09 10 11 Flash b Content: a c d Chip State: i | V V | i i V V i i i i reads easy, but writes require prior reads and erases 458

SSDs hybrid mapping

- direct writes to a few empty blocks (log blocks):
 - log table : per-page mappings (checked first) •
 - data table : per-block mapping (checked next) •

Say: $a \rightarrow 1000, b \rightarrow 1001, c \rightarrow 1002, d \rightarrow 1003$

Log Table: But what if re-write 1000, 1001? Data Table: 250 -8 Block: 0 2 1 Log Table: 1000→0 1001→1 00 01 02 03 04 05 06 07 08 09 10 11 Page: Data Table: 250 -> 8 a b c d V V V V Content: i V ii State: Block: 0 1 2 Page: 00 01 02 03 04 05 06 07 08 09 10 11 Content: a' b' abcd Log Table: 1000-0 1001-1 1002-2 1003-3 State: V V i i i i i i V V V V Data Table: 250 -8 Block: 0 2 1 Page: 00 01 02 03 04 05 06 07 08 09 10 11 Content: a' b' c' d' a b c d V V V V i i i i V V V V State: Log Table: Data Table: 250 →0 Block: 0 Page: 00 01 02 03 04 05 06 07 08 09 10 11 switch merge Content: a' b' c' d' State: V V V V i i i i i