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
- 43 Log-structured File Systems
- 44 Flash-based SSD
- 45 Data Integrity and Protection

398

Reading a file from disk

- Issue an open("/foo/bar", O_RDONLY),
- Traverse the pathname
 - begin at the root of the file system (/)
 - root inode number often 2 (specified in superblock)
 - read in block containing inode 2.
 - use "/" pointer blocks to get "/" directory contents
 - recurse on "/foo"
 - check permissions, memory for metadata, file descriptor
- when read() issued
 - · consult inode, find and read in first block
 - update open file table, file offset
- When file closed
 - dealloc file descriptor, logically the file may be deallocated, but not usually done here

Open and read / foo/bar timeline

	data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data [0]	bar data [1]	bar data [2]	
open(bar)			read	read		read	read				1
					read		icuu				
read()					read						t
					write			read			
read()					read						¥
					write				read		
read()					read						
	read()				write	te			read	ad	
	open(bar) read() read() read()	data bitmapopen(bar)read()read()read()	data bitmapinode bitmapopen(bar)	data bitmapinode inodeopen(bar)readread()readread()readread()read	data inode bitmaproot inodefoo inodeopen(bar)readreadread()readreadread()readreadread()readread	$ \begin{array}{ c c c c c } \begin{tabular}{ c c c c c c } \begin{tabular}{ c c c c c c c } \begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{ c c c c c } \begin{tabular}{ c c c c c } \begin{tabular}{ c c c c c } \begin{tabular}{ c c c c c c } \begin{tabular}{ c c c c c c } \begin{tabular}{ c c c c c c c } \begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{ c c c c c c } \begin{tabular}{ c c c c c } \begin{tabular}{ c c c c c } \begin{tabular}{ c c c c c } \begin{tabular}{ c c c c } \begin{tabular}{ c c c c c c } \begin{tabular}{ c c c c c c c } \begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{ c c c c c } \hline \begin{tabular}{ c c c } \hline \begin{tabular}{ c c } \hline \hline \begin{tabular}{ c c } \hline \begin{tabular}{ c c }$	$ \begin{array}{ c c c c c } \hline \begin{tabular}{c c c c } \hline \begin{tabular}{c c } \hline \hline \begin{tabular}{c c } \hline \hline \ \begin{tabular}{c c } \hline \hline \ \begin{tabular}{c c } \hline \hline \begin{tabular}{c c } \hline \hline \ \begin{tabular}{c c } \hline \hline \ \ \begin{tabular}{c c } \hline \hline \ \ \begin{tabular}{c c } \hline \hline \ \ \b$	data inode bitmaproot inodefoo inodebar inoderoot inodefoo databar databar databar databar databar databar databar databar databar databar databar databar databar databar databar databar databar databar databar

400

Create / foo/bar timeline

		data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data [0]	bar data [1]	bar data [2]	
				read	road		read					
5 I/Os 10 I/Os	create (/foo/bar)		read		Icau			read				
			write			read		write				
					write	write						· .
	write()	read write				read						
						write			write			+
	write()	read write				read				write		
						write						
	write()	read write				icau					write	
						write						

Unified Buffer Cache

- buffer cache
 - recently used pages/disk-blocks held in memory
 - writes *buffered* (delayed)
 - consecutive writes batched
 - scheduled more efficiently
 - dynamically partitioned (not fixed size)
- some apps (databases) ignore the cache and file system
 - call rsync()
 - use *direct I/O* interfaces to disk
 - write to raw disks

402

Locality and the Fast File System FFS

• "old" file system





- FFS uses block groups
 - disk maps onto cylinder groups
 - each has superblock, bitmaps, inodes, data

s	ib db	Inodes	Data
---	-------	--------	------

- "Keep related stuff together" single group
 - most directories
 - file data and related inodes
 - large files have chunks sprayed across multiple groups





iNodes FFS

- Use inode structure
 - direct links and blocks same group
 - indirect blocks, and blocks pointed to, different group
 - each 1024 blocks (4MB) in different group (4K pages, 4-byte ptr)



Errata

- more internal fragmentation
- used subblocks for space efficiency (small disks in stone ages)
 - copy to regular blocks when full
 - libc buffers, so most large files never subblock'd
- parameterization
 - blocks laid out so that OS has time to request block i + 1 after reading block it, *before* i + 1 rotates past
- track buffer
- long file names
- symbolic links