

### Recall: Important "ilities"

Availability: the probability that the system can accept and process requests

 Measured in "nines" of probability: e.g. 99.9% probability is "3-nines of availability"

- Key idea here is independence of failures
- Durability: the ability of a system to recover data despite faults
  - This idea is fault tolerance applied to data
  - Doesn't necessarily imply availability: information on pyramids was very durable, but could not be accessed until discovery of Rosetta Stone
- Reliability: the ability of a system or component to perform its required functions under stated conditions for a specified period of time (IEEE definition)
  - Usually stronger than simply availability: means that the system is not only "up", but also working correctly
  - Includes availability, security, fault tolerance/durability
  - Must make sure data survives system crashes, disk crashes, other problems

### Recall: RAID 6 and other Erasure Codes

- In general: RAIDX is an "erasure code"
  - Must have ability to know which disks are bad
  - Treat missing disk as an "Erasure"
- Today, disks so big that: RAID 5 not sufficient!
  - Time to repair disk sooooo long, another disk might fail in process!
  - "RAID 6" allow 2 disks in replication stripe to fail
  - Requires more complex erasure code, such as EVENODD code (see readings)
- More general option for general erasure code: Reed-Solomon codes
  - Based on polynomials in GF(2<sup>k</sup>) (I.e. k-bit symbols)
  - m data points define a degree m polynomial; encoding is n points on the polynomial
    - » P(x)=a<sub>0</sub>+a<sub>1</sub>x<sup>1</sup>+... a<sub>m-1</sub>x<sup>m-1</sup>
    - » Coded: P(0),P(1),P(2)....,P(n-1)
  - Any m points can be used to recover the polynomial; n m failures tolerated
- Erasure codes not just for disk arrays. For example, geographic replication
  - E.g., split data into m = 4 chunks, generate n = 16 fragments and distribute across the Internet
  - Any 4 fragments can be used to recover the original data --- very durable!

Lec 22.3

### Recall: COW with Smaller-Radix Blocks



 If file represented as a free of blocks, just need to update the leading fringe

Kubiatowicz CS162 © UCB Fall 2020

### More General Reliability Solutions

- · Use Transactions for atomic updates
  - Ensure that multiple related updates are performed atomically
  - i.e., if a crash occurs in the middle, the state of the systems reflects either all or none of the updates
  - Most modern file systems use transactions internally to update filesystem structures and metadata
  - Many applications implement their own transactions
- Provide Redundancy for media failures
  - Redundant representation on media (Error Correcting Codes)

Kubiatowicz CS162 © UCB Fall 2020

- Replication across media (e.g., RAID disk array)

Tr	ar	Isa	cti	O	ns

- · Closely related to critical sections for manipulating shared data structures
- They extend concept of atomic update from memory to stable storage
   Atomically update multiple persistent data structures
- · Many ad-hoc approaches
  - FFS carefully ordered the sequence of updates so that if a crash occurred while manipulating directory or inodes the disk scan on reboot would detect and recover the error (fsck)
  - Applications use temporary files and rename

# Key Concept: Transaction

• A *transaction* is an atomic sequence of reads and writes that takes the system from consistent state to another.



- · Recall: Code in a critical section appears atomic to other threads
- Transactions extend the concept of atomic updates from *memory* to *persistent storage*

11/16/20

Lec 22.7

Lec 22.5

11/16/20

Lec 22.6



#### Creating a File (No Journaling Yet) Journaling File Systems · Don't modify data structures on disk directly · Write each update as transaction recorded in a log Find free data block(s) - Commonly called a journal or intention list · Find free inode entry Free Also maintained on disk (allocate blocks for it when formatting) • Find dirent insertion point space Once changes are in the log, they can be safely applied to file system map Data blocks - e.g. modify inode pointers and directory mapping Write map (i.e., mark used) Inode table · Garbage collection: once a change is applied, remove its entry from the log Write inode entry to point to block(s) Directory · Write dirent to point to inode entries • Linux took original FFS-like file system (ext2) and added a journal to get ext3! - Some options: whether or not to write all data to journal or just metadata Other examples: NTFS, Apple HFS+, Linux XFS, JFS, ext4 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.13 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.14 Creating a File (With Journaling) After Commit, Eventually Replay Transaction Find free data block(s) · All accesses to the file system first looks in • Find free inode entry the loa Free Free · Find dirent insertion point space space - Actual on-disk data structure might be stale map Data blocks map Data blocks • [log] Write map (i.e., mark used) Inode table Inode table · Eventually, copy changes to disk and discard transaction from the log [log] Write inode entry to point to block(s) Directory Directory • [log] Write dirent to point to inode entries entries tail head tail head tail tail tail tail commit commit start start pending done done pending Log: in non-volatile storage (Flash or on Disk) Log: in non-volatile storage (Flash or on Disk) 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.15 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.16

# Crash Recovery: Discard Partial Transactions

· Upon recovery, scan the log Free · Detect transaction start with no commit space map Data blocks · Discard log entries Inode table Disk remains unchanged Directory entries tail head pending start done Log: in non-volatile storage (Flash or on Disk) 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.17

## Crash Recovery: Keep Complete Transactions



### **Journaling Summary**

#### Why go through all this trouble?

- · Updates atomic, even if we crash:
  - Update either gets fully applied or discarded
  - All physical operations treated as a logical unit

#### Isn't this expensive?

- Yes! We're now writing all data twice (once to log, once to actual data blocks in target file)
- Modern filesystems journal metadata updates only
  - Record modifications to file system data structures
  - But apply updates to a file's contents directly

# Recall: CS 162 Collaboration Policy



Explaining a concept to someone in another group Discussing algorithms/testing strategies with other groups Discussing debugging approaches with other groups Searching online for generic algorithms (e.g., hash table)

#### Sharing code or test cases with another group



Copying OR reading another group's code or test cases Copying OR reading online code or test cases from prior years Helping someone in another group to debug their code

- We compare all project submissions against prior year submissions and online solutions and will take actions (described on the course overview page) against offenders
- Don't put a friend in a bad position by asking for help that they shouldn't give!

# The Log Structured File System (LFS)



# Going Further – Log Structured File Systems

- · The log IS what is recorded on disk
  - File system operations logically replay log to get result
  - Create data structures to make this fast
  - On recovery, replay the log
- · Index (inodes) and directories are written into the log too
- Large, important portion of the log is cached in memory
   Relies on Buffer Cache to make reading fast
- · Do everything in bulk: log is collection of large segments
- Each segment contains a summary of all the operations within the segment - Fast to determine if segment is relevant or not
- Free space is approached as continual cleaning process of segments
   Detect what is live or not within a segment
  - Copy live portion to new segment being formed (replay)
  - Garbage collection entire segment
  - No bit map

11/16/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 22.22

# What about Flash Filesystems?

- Cannot overwrite pages!
  - Must move contents to an erased page
  - Small changes  $\Rightarrow$  lots of rewriting of data/wear out!
- Program/Erase (PE) Wear
  - Permanent damage to gate oxide at each flash cell
  - Caused by high program/erase voltages
  - Issues: trapped charges, premature leakage of charge
  - Need to balance how frequently cells written: "Wear Leveling"
- Flash Translation Layer (FTL)
  - Translates between Logical Block Addresses (at OS level) and Physical Flash Page Addresses
  - Manages the wear and erasure state of blocks and pages
  - Tracks which blocks are garbage but not erased
- Management Process (Firmware)
  - Keep freelist full, Manage mapping,
  - Track wear state of pages
  - Copy good pages out of mostly empty blocks before erasure

Data written in 4 KB Pages	•	4 KB	4 KB	4 KB		
Data erased		4 KB	4 KB	4 KB		
Blocks	-					
4 writable Pages n 1 erasable Bloc	s k	4 KB	4 KB	4 KB		
Tunical NAND Flash Bases and Blas						

Typical NAND Flash Pages and Block



Single FLASH storage bit

# Example Use of LFS: F2FS: A Flash File System

- File system used on many mobile devices
  - Including the Pixel 3 from Google
  - Latest version supports block-encryption for security
  - Has been "mainstream" in linux for several years now
- Assumes standard SSD interface
  - With built-in Flash Translation Layer (FTL)
  - Random reads are as fast as sequential reads
  - Random writes are bad for flash storage
    - » Forces FTL to keep moving/coalescing pages and erasing blocks
       » Sustained write performance degrades/lifetime reduced
- Minimize Writes/updates and otherwise keep writes "sequential"
  - Start with Log-structured file systems/copy-on-write file systems
  - Keep writes as sequential as possible
  - Node Translation Table (NAT) for "logical" to "physical" translation » Independent of FTL
- For more details, check out paper in Readings section of website
  - "F2FS: A New File System for Flash Storage" (from 2015)
  - Design of file system to leverage and optimize NAND flash solutions
  - Comparison with Ext4, Btrfs, Nilfs2, etc

Kubiatowicz CS162 © UCB Fall 2020

Lec 22.23

11/16/20

# e FLASH storage bit







11/16/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 22.27



## **Distributed Systems: Reality**

- · Reality has been disappointing
  - Worse availability: depend on every machine being up
     » Lamport: "A distributed system is one in which the failure of a compute you didn't even know existed can render your own computer unusable
  - Worse reliability: can lose data if any machine crashes
  - Worse security: anyone in world can break into system
- Coordination is more difficult
  - Must coordinate multiple copies of shared state information
  - What would be easy in a centralized system becomes a lot more difficult
- Trust/Security/Privacy/Denial of Service
  - Many new variants of problems arise as a result of distribution
  - Can you trust the other members of a distributed application enough to even perform a protocol correctly?
  - Corollary of Lamport's quote: "A distributed system is one where you can't do work because some computer you didn't even know existed is successfully coordinating an attack on my system!"

Leslie Lamport

# Distributed Systems: Goals/Requirements

- Transparency: the ability of the system to mask its complexity behind a simple interface
- Possible transparencies:
  - Location: Can't tell where resources are located
  - Migration: Resources may move without the user knowing
  - Replication: Can't tell how many copies of resource exist
  - Concurrency: Can't tell how many users there are
  - Parallelism: System may speed up large jobs by splitting them into smaller pieces
  - Fault Tolerance: System may hide various things that go wrong
- Transparency and collaboration require some way for different processors to communicate with one another



Kubiatowicz CS162 © UCB Fall 2020

Lec 22.31

11/16/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 22.32



Kubiatowicz CS162 © UCB Fall 2020

Lec 22.35

#### The Internet Hourglass Implications of Hourglass Single Internet-layer module (IP): · Allows arbitrary networks to interoperate SMTP HTTP DNS NTP Applications - Any network technology that supports IP can exchange packets Allows applications to function on all networks UDP TCP Transport - Applications that can run on IP can use any network Waist Supports simultaneous innovations above and below IP -But changing IP itself, i.e., IPv6, very involved Data Link SONET 802.11 **Physical** Fiber Radio The Hourglass Model Copper There is just one network-layer protocol, IP. The "narrow waist" facilitates interoperability. 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.37 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.38 Drawbacks of Layering **End-To-End Argument** • Hugely influential paper: "End-to-End Arguments in System Design" by Layer N may duplicate layer N-1 functionality Saltzer, Reed, and Clark ('84) - E.g., error recovery to retransmit lost data "Sacred Text" of the Internet Layers may need same information - Endless disputes about what it means - E.g., timestamps, maximum transmission unit size - Everyone cites it as supporting their position · Layering can hurt performance • Simple Message: Some types of network functionality can only be correctly - E.g., hiding details about what is really going on implemented end-to-end · Some layers are not always cleanly separated - Reliability, security, etc. - Inter-layer dependencies for performance reasons · Because of this, end hosts: - Some dependencies in standards (header checksums) - Can satisfy the requirement without network's help · Headers start to get really big - Will/must do so, since can't rely on network's help - Sometimes header bytes >> actual content Therefore don't go out of your way to implement them in the network 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.39 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.40

#### Example: Reliable File Transfer Discussion • Solution 1 is incomplete Host A Host B - What happens if memory is corrupted? - Receiver has to do the check anyway! Appl. OS • Solution 2 is complete OK - Full functionality can be entirely implemented at application layer with no need for reliability from lower layers Solution 1: make each step reliable, and then concatenate them • Is there any need to implement reliability at lower layers? - Well, it could be more efficient • Solution 2: end-to-end check and try again if necessary 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.41 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.42 **Conservative Interpretation of E2E End-to-End Principle** Implementing complex functionality in the network: · Don't implement a function at the lower levels of the system Doesn't reduce host implementation complexity unless it can be completely implemented at this level · Does increase network complexity · Probably imposes delay and overhead on all applications, even if Or: Unless you can relieve the burden from hosts, don't bother they don't need functionality · However, implementing in network can enhance performance in some cases -e.g., very lossy link 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.43 11/16/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 22.44



### Using Messages: Send/Receive behavior

- When should send (message, mbox) return?
  - When receiver gets message? (i.e. ack received)
  - When message is safely buffered on destination?
  - Right away, if message is buffered on source node?
- · Actually two questions here:
  - When can the sender be sure that receiver actually received the message?
  - When can sender reuse the memory containing message?
- Mailbox provides 1-way communication from T1 $\rightarrow$ T2
  - $-\operatorname{T1} \rightarrow \! \text{buffer} \! \rightarrow \! \text{T2}$
  - Very similar to producer/consumer
    - » Send = V, Receive = P
    - » However, can't tell if sender/receiver is local or not!

# Messaging for Producer-Consumer Style

• Using send/receive for producer-consumer style:



- No need for producer/consumer to keep track of space in mailbox: handled by send/receive
  - Next time: will discuss fact that this is one of the roles the window in TCP: window is size of buffer on far end
  - Restricts sender to forward only what will fit in buffer

Lec 22.47



