

CS162

Operating Systems and Systems Programming

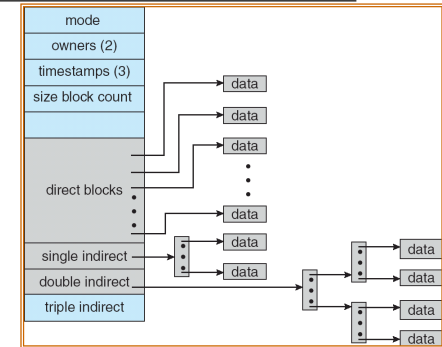
Lecture 20

Reliability, Transactions Distributed Systems

April 16th, 2019
 Prof. John Kubiatowicz
<http://cs162.eecs.Berkeley.edu>

Recall: Multilevel Indexed Files (Original 4.1 BSD)

- Sample file in multilevel indexed format:
 - 10 direct ptrs, 1K blocks
 - How many accesses for block #23? (assume file header accessed on open)?
 - » Two: One for indirect block, one for data
 - How about block #5?
 - » One: One for data
 - Block #340?
 - » Three: double indirect block, indirect block, and data
- UNIX 4.1 Pros and cons
 - Pros: Simple (more or less)
Files can easily expand (up to a point)
Small files particularly cheap and easy
 - Cons: Lots of seeks (lead to 4.2 Fast File System Optimizations)
- Ext2/3 (Linux):
 - 12 direct ptrs, triply-indirect blocks, settable block size (4K is common)



4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.2

Recall: File System Caching

- Key Idea: Exploit locality by caching data in memory
 - Name translations: Mapping from paths→inodes
 - Disk blocks: Mapping from block address→disk content
- **Buffer Cache:** Memory used to cache kernel resources, including disk blocks and name translations
 - Can contain “dirty” blocks (blocks yet on disk)
- Replacement policy? LRU
 - Can afford overhead full LRU implementation
 - Advantages:
 - » Works very well for name translation
 - » Works well in general as long as memory is big enough to accommodate a host’s working set of files.
 - Disadvantages:
 - » Fails when some application scans through file system, thereby flushing the cache with data used only once
 - » Example: `find . -exec grep foo {} \;`
- Other Replacement Policies?
 - Some systems allow applications to request other policies
 - Example, ‘Use Once’:
 - » File system can discard blocks as soon as they are used

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.3

File System Caching (con’t)

- Cache Size: How much memory should the OS allocate to the buffer cache vs virtual memory?
 - Too much memory to the file system cache ⇒ won’t be able to run many applications at once
 - Too little memory to file system cache ⇒ many applications may run slowly (disk caching not effective)
 - Solution: adjust boundary dynamically so that the disk access rates for paging and file access are balanced
- **Read Ahead Prefetching:** fetch sequential blocks early
 - Key Idea: exploit fact that most common file access is sequential by prefetching subsequent disk blocks ahead of current read request (if they are not already in memory)
 - Elevator algorithm can efficiently interleave groups of prefetches from concurrent applications
 - How much to prefetch?
 - » Too many imposes delays on requests by other applications
 - » Too few causes many seeks (and rotational delays) among concurrent file requests

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.4

File System Caching (con't)

- **Delayed Writes:** Writes to files not immediately sent out to disk
 - Instead, `write()` copies data from user space buffer to kernel buffer (in cache)
 - » Enabled by presence of buffer cache: can leave written file blocks in cache for a while
 - » If some other application tries to read data before written to disk, file system will read from cache
 - Flushed to disk periodically (e.g. in UNIX, every 30 sec)
 - Advantages:
 - » Disk scheduler can efficiently order lots of requests
 - » Disk allocation algorithm can be run with correct size value for a file
 - » Some files need never get written to disk! (e..g temporary scratch files written /tmp often don't exist for 30 sec)
 - Disadvantages
 - » What if system crashes before file has been written out?
 - » Worse yet, what if system crashes before a directory file has been written out? (lose pointer to inode!)

Important “ilities”

- **Availability:** the probability that the system can accept and process requests
 - Often measured in “nines” of probability. So, a 99.9% probability is considered “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

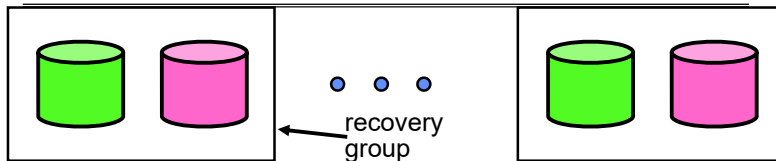
How to Make File System Durable?

- Disk blocks contain Reed-Solomon error correcting codes (ECC) to deal with small defects in disk drive
 - Can allow recovery of data from small media defects
- Make sure writes survive in short term
 - Either abandon delayed writes or
 - Use special, battery-backed RAM (called non-volatile RAM or **NVRAM**) for dirty blocks in buffer cache
- Make sure that data survives in long term
 - Need to replicate! More than one copy of data!
 - Important element: **independence of failure**
 - » Could put copies on one disk, but if disk head fails...
 - » Could put copies on different disks, but if server fails...
 - » Could put copies on different servers, but if building is struck by lightning....
 - » Could put copies on servers in different continents...

RAID: Redundant Arrays of Inexpensive Disks

- Classified by David Patterson, Garth A. Gibson, and Randy Katz here at UCB in 1987
 - Classic paper was first to evaluate multiple schemes
- Data stored on multiple disks (redundancy)
 - Berkeley researchers were looking for alternatives to big expensive disks
 - Redundancy necessary because cheap disks were more error prone
- Either in software or hardware
 - In hardware case, done by disk controller; file system may not even know that there is more than one disk in use
- Initially, five levels of RAID (more now)

RAID 1: Disk Mirroring/Shadowing



- Each disk is fully duplicated onto its “shadow”
 - For high I/O rate, high availability environments
 - Most expensive solution: 100% capacity overhead
- Bandwidth sacrificed on write:
 - Logical write = two physical writes
 - Highest bandwidth when disk heads and rotation fully synchronized (hard to do exactly)
- Reads may be optimized
 - Can have two independent reads to same data
- Recovery:
 - Disk failure \Rightarrow replace disk and copy data to new disk
 - **Hot Spare**: idle disk already attached to system to be used for immediate replacement

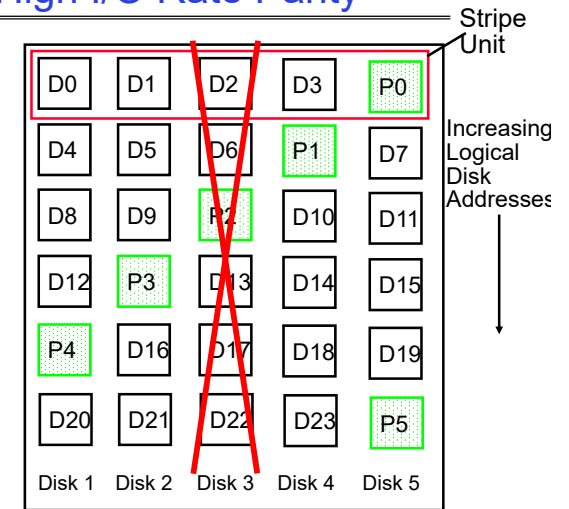
4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.9

RAID 5+: High I/O Rate Parity

- Data striped across multiple disks
 - Successive blocks stored on successive (non-parity) disks
 - Increased bandwidth over single disk
- Parity block (in green) constructed by XORing data blocks in stripe
 - $P_0 = D_0 \oplus D_1 \oplus D_2 \oplus D_3$
 - Can destroy any one disk and still reconstruct data
 - Suppose Disk 3 fails, then can reconstruct: $D_2 = D_0 \oplus D_1 \oplus D_3 \oplus P_0$



- Can spread information widely across internet for durability
 - RAID algorithms work over geographic scale

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.10

Allow more disks to fail!

- 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
- But – must do something more complex than just XORing together blocks!
 - Already used up the simple XOR operation across disks
- Simple option: Check out **EVENODD** code in readings
 - Will generate one additional check disks to support RAID 6
- More general option for general erasure code: **Reed-Solomon codes**
 - Based on polynomials in $GF(2^k)$ (i.e. k-bit symbols)
 - » Galois Field is finite version of real numbers
 - Data as coefficients (a_i), code space as values of polynomial:
 - » $P(x) = a_0 + a_1x + \dots + a_{m-1}x^{m-1}$
 - » Coded: $P(0), P(1), P(2), \dots, P(n-1)$
 - Can recover polynomial (i.e. data) as long as get any m of n ; allows $n-m$ failures!

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.11

Allow more disks to fail! (Con't)

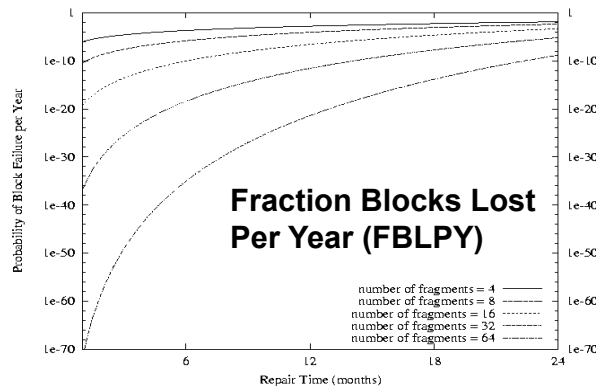
- How to use **Reed-Solomon** code in practice?
 - Each coefficient has a fixed (k) number of bits. So, must encode with symbols that size
 - Example: $k=16$ bit symbols, $m=4$, encoding 16×4 bits at a time
 - » Take original data, split into 4 chunks. On each encoding step, grab 16 bits from each chunk to use as coefficients
 - » Each data point yields a 16-bit symbol, which you distributed to final encoded chunks
 - (better version of Reed-Solomon code for erasure channels is the “Cauchy Reed-Solomon” code; it is isomorphic to the version here)
- Examples (with $k=16$):
 - Suppose have 6 disks, want to tolerate 2 failures
 - » Split data into 4 chunks, encode 16 bits from each chunk at a time, by generating 6 points (of 16 bits) on 3rd-degree polynomial
 - » Distribute data from polynomial to 6 disks – each disk will ultimately hold data that is $\frac{1}{4}$ size of original data
 - » Can handle 2 lost disks for 50% overhead
 - More interesting extreme for Internet-level replication:
 - » Split data into 4 chunks, produce 16 chunks
 - » Each chunk is $\frac{1}{4}$ total size of original data, Overhead = factor of 4
 - » But – only need 4 of 16 fragments! **REALLY DURABLE!**

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.12

Use of Erasure Coding in general: High Durability/overhead ratio!



- Exploit law of large numbers for durability!
- 6 month repair, FBLPY with 4x increase in total size of data:
 - Replication (4 copies): 0.03
 - Fragmentation (16 of 64 fragments needed): 10^{-35}

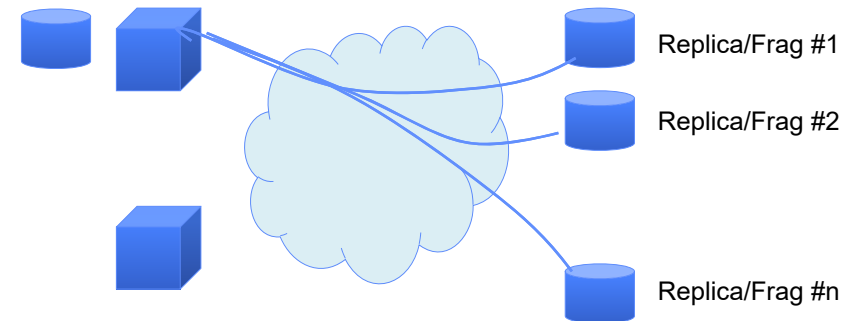
4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.13

Higher Durability/Reliability through Geographic Replication

- Highly durable – hard to destroy all copies
 - Simple replication: read any copy
 - Erasure coded: read m of n
- Highly available for reads
 - Simple replication: read any copy
 - Erasure coded: read m of n
- Low availability for writes
 - Can't write if any one replica is not up
 - Or – need relaxed consistency model
- Reliability? – availability, security, durability, fault-tolerance



4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.14

Administrivia

- Last Midterm: 5/2
 - Can have 3 handwritten sheets of notes – both sides
 - Focus on material from lecture 17-24, but all topics fair game!
- Don't forget to do your group evaluations!
 - Very important to help us understand your group dynamics
- Optional HW4 will come out soon
 - Will give you a chance to try out using the language "Go" to build a two-phase commit protocol
 - You will be testing it out for next term
 - » Not sure that we will be giving out points for it. Stay tuned!

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.15

File System Reliability: (Difference from Block-level reliability)

- What can happen if disk loses power or software crashes?
 - Some operations in progress may complete
 - Some operations in progress may be lost
 - Overwrite of a block may only partially complete
- Having RAID doesn't necessarily protect against all such failures
 - No protection against writing bad state
 - What if one disk of RAID group not written?
- File system needs durability (as a minimum!)
 - Data previously stored can be retrieved (maybe after some recovery step), regardless of failure

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.16

Storage Reliability Problem

- Single logical file operation can involve updates to multiple physical disk blocks
 - inode, indirect block, data block, bitmap, ...
 - With sector remapping, single update to physical disk block can require multiple (even lower level) updates to sectors
- At a physical level, operations complete one at a time
 - Want concurrent operations for performance
- How do we guarantee consistency regardless of when crash occurs?

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.17

Threats to Reliability

- Interrupted Operation
 - Crash or power failure in the middle of a series of related updates may leave stored data in an inconsistent state
 - Example: transfer funds from one bank account to another
 - What if transfer is interrupted after withdrawal and before deposit?
- Loss of stored data
 - Failure of non-volatile storage media may cause previously stored data to disappear or be corrupted

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.18

Reliability Approach #1: Careful Ordering

- Sequence operations in a specific order
 - Careful design to allow sequence to be interrupted safely
- Post-crash recovery
 - Read data structures to see if there were any operations in progress
 - Clean up/finish as needed
- Approach taken by
 - FAT and FFS (fsck) to protect filesystem structure/metadata
 - Many app-level recovery schemes (e.g., Word, emacs autosaves)

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.19

FFS: Create a File

Normal operation:

- Allocate data block
- Write data block
- Allocate inode
- Write inode block
- Update bitmap of free blocks and inodes
- Update directory with file name → inode number
- Update modify time for directory

Recovery:

- Scan inode table
- If any unlinked files (not in any directory), delete or put in lost & found dir
- Compare free block bitmap against inode trees
- Scan directories for missing update/access times

Time proportional to disk size

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.20

Reliability Approach #2: Copy on Write File Layout

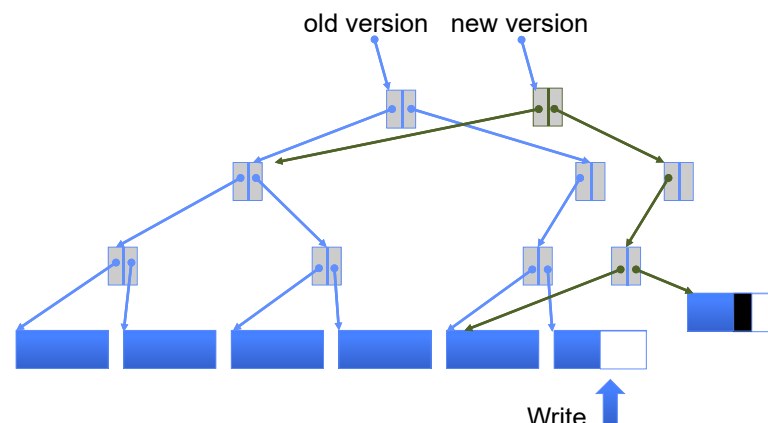
- To update file system, write a new version of the file system containing the update
 - Never update in place
 - Reuse existing unchanged disk blocks
- Seems expensive! But
 - Updates can be batched
 - Almost all disk writes can occur in parallel
- Approach taken in network file server appliances
 - NetApp's Write Anywhere File Layout (WAFL)
 - ZFS (Sun/Oracle) and OpenZFS

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.21

COW with Smaller-Radix Blocks



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

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.22

ZFS and OpenZFS

- Variable sized blocks: 512 B – 128 KB
- Symmetric tree
 - Know if it is large or small when we make the copy
- Store version number with pointers
 - Can create new version by adding blocks and new pointers
- Buffers a collection of writes before creating a new version with them
- Free space represented as tree of extents in each block group
 - Delay updates to freespace (in log) and do them all when block group is activated

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.23

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

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

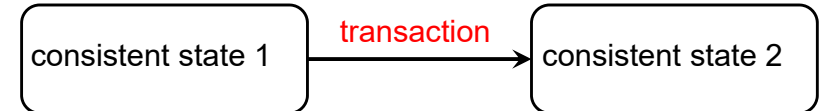
Lec 19.24

Transactions

- 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

- An **atomic sequence** of actions (reads/writes) on a storage system (or database)
- That takes it from one **consistent state** to another



Typical Structure

- **Begin** a transaction – get transaction id
- Do a bunch of updates
 - If any fail along the way, **roll-back**
 - Or, if any conflicts with other transactions, **roll-back**
- **Commit** the transaction

“Classic” Example: Transaction

```
BEGIN;    --BEGIN TRANSACTION
UPDATE accounts SET balance = balance - 100.00 WHERE
  name = 'Alice';

UPDATE branches SET balance = balance - 100.00 WHERE
  name = (SELECT branch_name FROM accounts WHERE name
    = 'Alice');

UPDATE accounts SET balance = balance + 100.00 WHERE
  name = 'Bob';

UPDATE branches SET balance = balance + 100.00 WHERE
  name = (SELECT branch_name FROM accounts WHERE name
    = 'Bob');

COMMIT;   --COMMIT WORK
```

Transfer \$100 from Alice's account to Bob's account

The ACID properties of Transactions

- **Atomicity:** all actions in the transaction happen, or none happen
- **Consistency:** transactions maintain data integrity, e.g.,
 - Balance cannot be negative
 - Cannot reschedule meeting on February 30
- **Isolation:** execution of one transaction is isolated from that of all others; no problems from concurrency
- **Durability:** if a transaction commits, its effects persist despite crashes

Transactional File Systems

- Better reliability through use of log
 - All changes are treated as *transactions*
 - A transaction is *committed* once it is written to the log
 - » Data forced to disk for reliability
 - » Process can be accelerated with NVRAM
 - Although File system may not be updated immediately, data preserved in the log
- Difference between “Log Structured” and “Journaled”
 - In a Log Structured filesystem, data stays in log form
 - In a Journaled filesystem, Log used for recovery
- Journaling File System
 - Applies updates to system metadata using transactions (using logs, etc.)
 - Updates to non-directory files (i.e., user stuff) can be done in place (without logs), full logging optional
 - Ex: NTFS, Apple HFS+, Linux XFS, JFS, ext3, ext4
- Full Logging File System
 - All updates to disk are done in transactions

Journaled File Systems

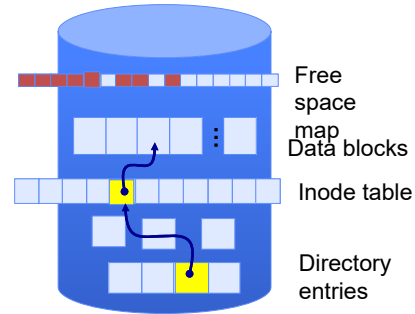
- Instead of modifying data structures on disk directly, write changes to a journal/log
 - Intention list: set of changes we intend to make
 - Log/Journal is **append-only**
 - Single commit record commits transaction
- Once changes are in the log, it is safe to apply changes to data structures on disk
 - Recovery can read log to see what changes were intended
 - Can take our time making the changes
 - » As long as new requests consult the log first
- Once changes are copied, safe to remove log
- But, ...
 - If the last atomic action is not done ... poof ... all gone
- Basic assumption:
 - Updates to sectors are atomic and ordered
 - Not necessarily true unless very careful, but key assumption

Redo Logging

- Prepare
 - Write all changes (in transaction) to log
- Commit
 - Single disk write to make transaction durable
- Redo
 - Copy changes to disk
- Garbage collection
 - Reclaim space in log
- Recovery
 - Read log
 - Redo any operations for committed transactions
 - Garbage collect log

Example: Creating a File

- Find free data block(s)
- Find free inode entry
- Find dirent insertion point
-
- Write map (i.e., mark used)
- Write inode entry to point to block(s)
- Write dirent to point to inode



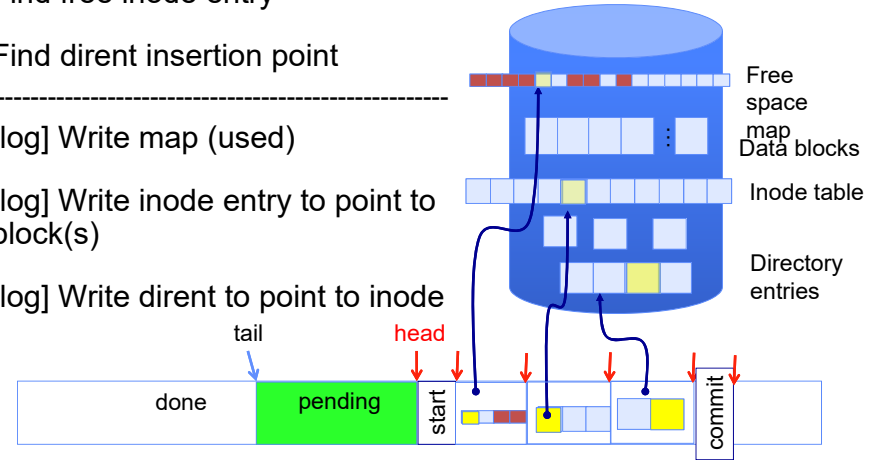
4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.33

Ex: Creating a file (as a transaction)

- Find free data block(s)
- Find free inode entry
- Find dirent insertion point
-
- [log] Write map (used)
- [log] Write inode entry to point to block(s)
- [log] Write dirent to point to inode



Log: in non-volatile storage (Flash or on Disk)

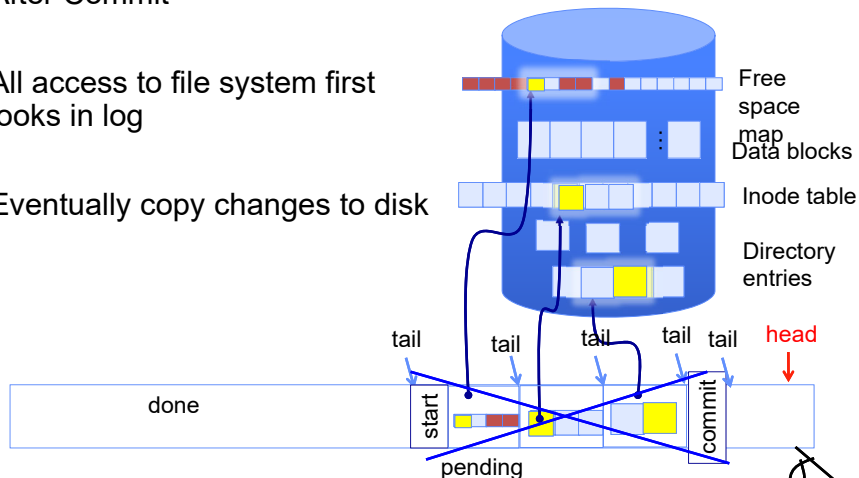
4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.34

ReDo Log

- After Commit
- All access to file system first looks in log
- Eventually copy changes to disk



Log: in non-volatile storage (Flash or Disk)

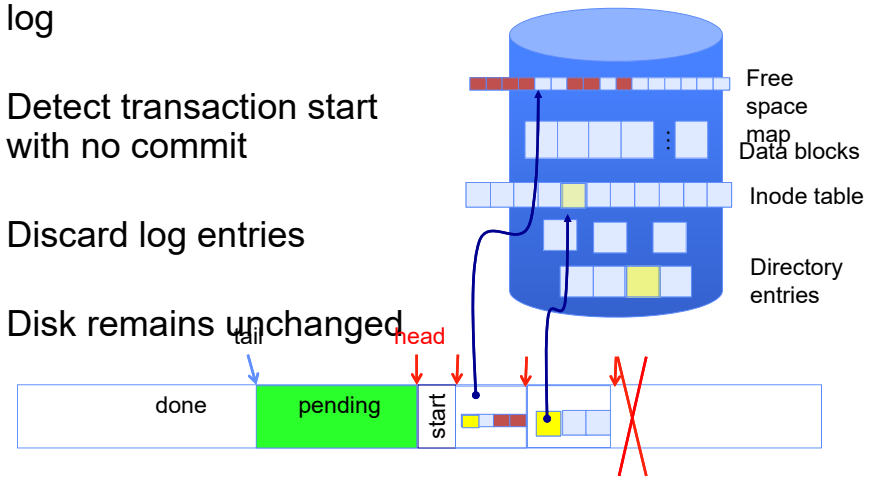
4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.35

Crash During Logging – Recover

- Upon recovery scan the log
- Detect transaction start with no commit
- Discard log entries
- Disk remains unchanged



Log: in non-volatile storage (Flash or on Disk)

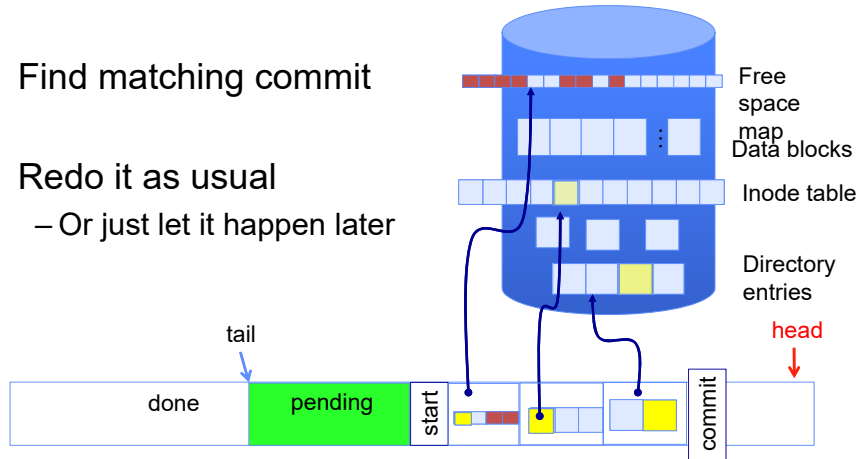
4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.36

Recovery After Commit

- Scan log, find start
- Find matching commit
- Redo it as usual
 - Or just let it happen later



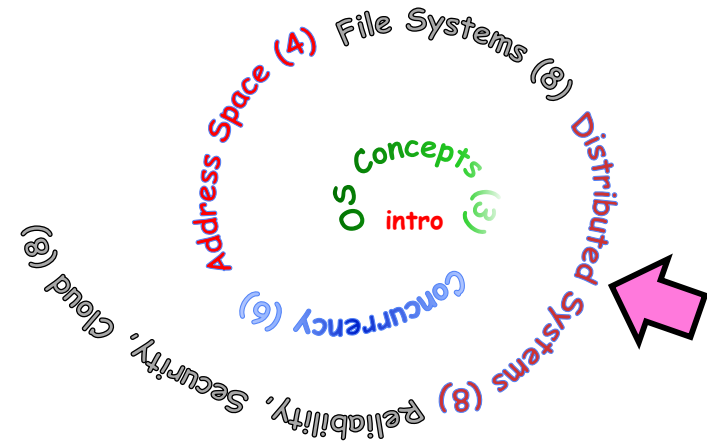
Log: in non-volatile storage (Flash or on Disk)

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.37

Course Structure: Spiral



4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.38

Societal Scale Information Systems

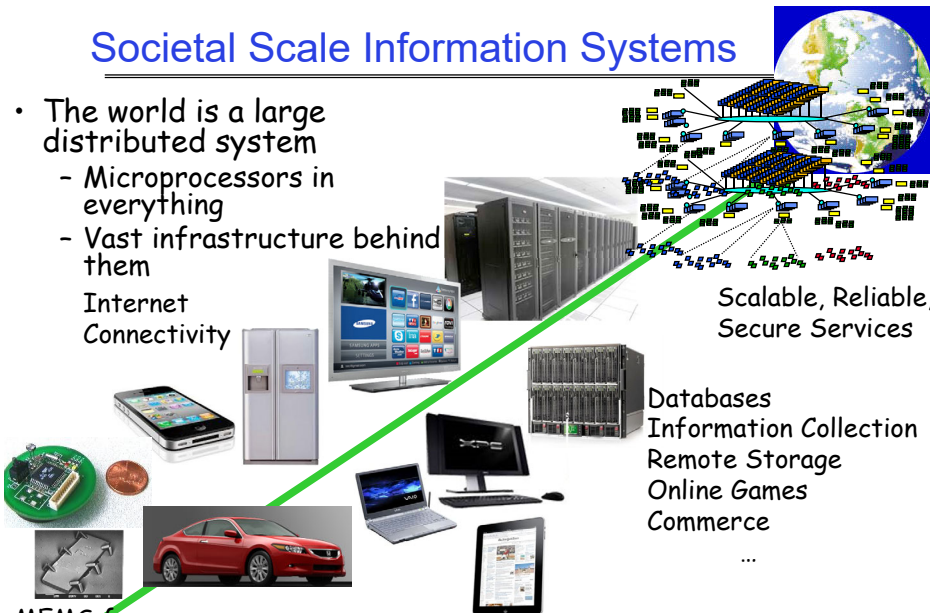
- The world is a large distributed system
 - Microprocessors in everything
 - Vast infrastructure behind them

Internet Connectivity

Scalable, Reliable, Secure Services

Databases
Information Collection
Remote Storage
Online Games
Commerce
...

MEMS for Sensor Nets

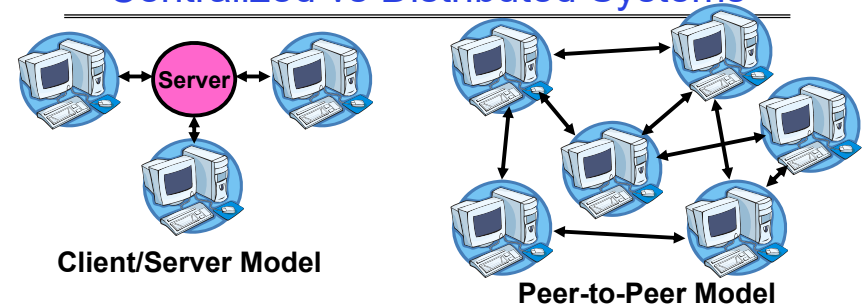


4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.39

Centralized vs Distributed Systems



- **Centralized System:** System in which major functions are performed by a single physical computer
 - Originally, everything on single computer
 - Later: client/server model
- **Distributed System:** physically separate computers working together on some task
 - Early model: multiple servers working together
 - » Probably in the same room or building
 - » Often called a “cluster”
 - Later models: peer-to-peer/wide-spread collaboration

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.40

Distributed Systems: Motivation/Issues/Promise

- Why do we want distributed systems?
 - Cheaper and easier to build lots of simple computers
 - Easier to add power incrementally
 - Users can have complete control over some components
 - Collaboration: much easier for users to collaborate through network resources (such as network file systems)
- The **promise** of distributed systems:
 - Higher availability: one machine goes down, use another
 - Better durability: store data in multiple locations
 - More security: each piece easier to make secure

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.41

Distributed Systems: Reality

- Reality has been disappointing
 - Worse availability: depend on every machine being up
 - » Lamport: “a distributed system is one where I can't do work because some machine I've never heard of isn't working!”
 - 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 (using only a network)
 - What would be easy in a centralized system becomes a lot more difficult

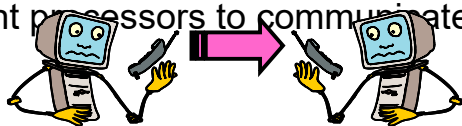
4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.42

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

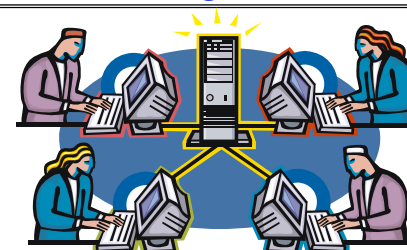


4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.43

Networking Definitions



- **Network**: physical connection that allows two computers to communicate
- **Packet**: unit of transfer, sequence of bits carried over the network
 - Network carries packets from one CPU to another
 - Destination gets interrupt when packet arrives
- **Protocol**: agreement between two parties as to how information is to be transmitted

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.44

What Is A Protocol?

- A protocol is an **agreement on how to communicate**
- Includes
 - **Syntax**: how a communication is specified & structured
 - » Format, order messages are sent and received
 - **Semantics**: what a communication means
 - » Actions taken when transmitting, receiving, or when a timer expires
- Described formally by a state machine
 - Often represented as a message transaction diagram

4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.45

Examples of Protocols in Human Interactions

- Telephone
 1. (Pick up / open up the phone)
 2. Listen for a dial tone / see that you have service
 3. Dial
 4. Should hear ringing ...
 5. Callee: "Hello?"
 6. Caller: "Hi, it's John..."
Or: "Hi, it's me" (← what's *that* about?)
 7. Caller: "Hey, do you think ... blah blah blah ..." **pause**
 1. Callee: "Yeah, blah blah blah ..." **pause**
 2. Caller: Bye
 3. Callee: Bye
 4. Hang up

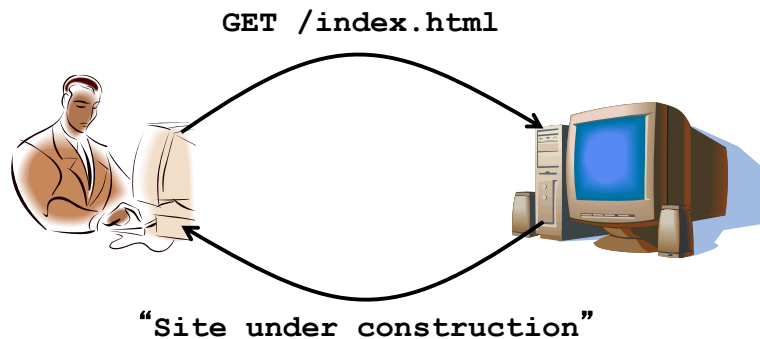
4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.46

Clients and Servers

- Client program
 - Running on end host
 - Requests service
 - E.g., Web browser
- Server program
 - Running on end host
 - Provides service
 - E.g., Web server



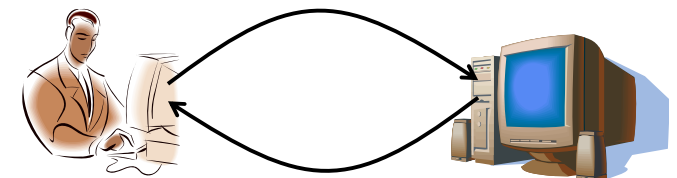
4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.47

Client-Server Communication

- Client "sometimes on"
 - Initiates a request to the server when interested
 - E.g., Web browser on your laptop or cell phone
 - Doesn't communicate directly with other clients
 - Needs to know the server's address
- Server is "always on"
 - Services requests from many client hosts
 - E.g., Web server for the *www.cnn.com* Web site
 - Doesn't initiate contact with the clients
 - Needs a fixed, well-known address



4/16/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.48

Peer-to-Peer Communication

- No always-on server at the center of it all
 - Hosts can come and go, and change addresses
 - Hosts may have a different address each time
- Example: peer-to-peer file sharing (e.g., BitTorrent)
 - Any host can request files, send files, query to find where a file is located, respond to queries, and forward queries
 - Scalability by harnessing millions of peers
 - Each peer acting as **both a client and server**

Summary

- Important system properties
 - **Availability**: how often is the resource available?
 - **Durability**: how well is data preserved against faults?
 - **Reliability**: how often is resource performing correctly?
- **RAID**: Redundant Arrays of Inexpensive Disks
 - RAID1: mirroring, RAID5: Parity block
- Use of Log to improve Reliability
 - Journalled file systems such as ext3, NTFS
- **Transactions**: ACID semantics
 - Atomicity
 - Consistency
 - Isolation
 - Durability