

CS162 Operating Systems and Systems Programming Lecture 20

Reliability, Transactions Distributed Systems

November 9th, 2015

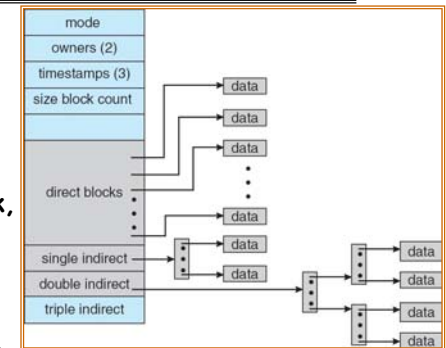
Prof. John Kubiatowicz

<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
Very large files must read many indirect block (four I/Os per block!)

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.2

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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.3

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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.4

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!)

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.5

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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.6

How to make file system durable?

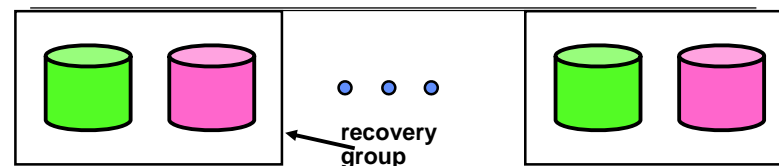
- Disk blocks contain Reed-Solomon error correcting codes (ECC) to deal with small defects in disk drive
 - Can allow recovery of data from small media defects
- Make sure writes survive in short term
 - Either abandon delayed writes or
 - use special, battery-backed RAM (called non-volatile RAM or **NVRAM**) 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
 - Data stored on multiple disks (redundancy)
 - 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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.7

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 ⇒ replace disk and copy data to new disk
 - **Hot Spare:** idle disk already attached to system to be used for immediate replacement

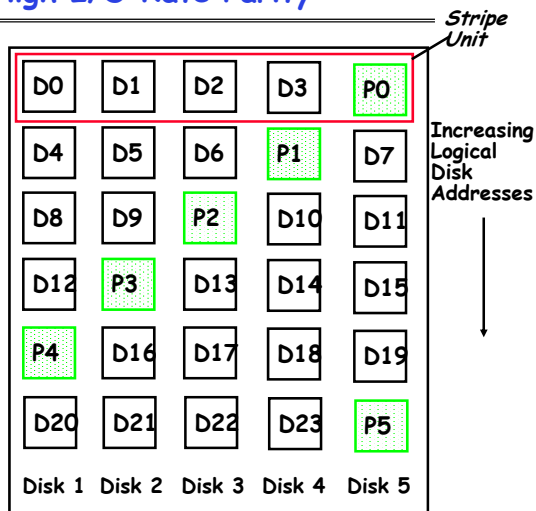
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.8

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 D3 fails, then can reconstruct: $D_3 = D_0 \oplus D_1 \oplus D_2 \oplus P_0$



- Later in term: talk about spreading information widely across internet for durability.

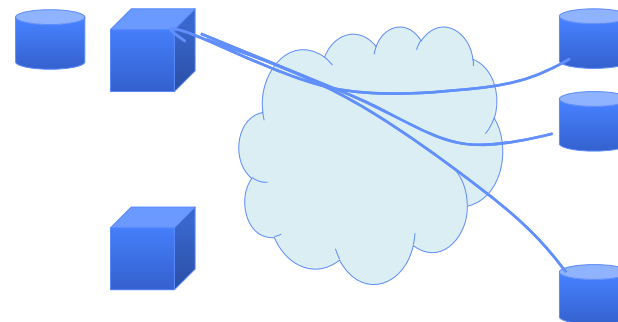
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.9

Higher Durability/Reliability through Geographic Replication

- Highly durable - hard to destroy bits
- Highly available for reads
- Low availability for writes
 - Can't write if any one is not up
 - Or - need relaxed consistency model
- Reliability?



11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.10

File System Reliability

- What can happen if disk loses power or machine 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
 - Bit-for-bit protection of bad state?
 - What if one disk of RAID group not written?
- File system wants durability (as a minimum!)
 - Data previously stored can be retrieved (maybe after some recovery step), regardless of failure

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.11

Storage Reliability Problem

- Single logical file operation can involve updates to multiple physical disk blocks
 - inode, indirect block, data block, bitmap, ...
 - With remapping, single update to physical disk block can require multiple (even lower level) updates
- 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?

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.12

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*.
 - e.g.: transfer funds from BofA to Schwab. 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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.13

Administrivia

- **Midterm II: Coming up in 2 weeks! (11/23)**
 - 7-10PM, "here" (2040, 2050, 2060 VLSB)
 - Topics up to and including previous Wednesday
 - 1 page of hand-written notes, both sides
- **Moved HW4 forward 1 week (hand out next Monday)**
- **No class on Wednesday (it is a holiday)**
- **Only 5 official lectures left (including this one!)**
- **Final (optional) lecture**
 - Monday of RRR week (12/07)
 - Whatever topics you would like!
 - Let me know what you want to hear about
 - Examples: IoT, security hardware, quantum computing.....

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.14

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 in FAT, FFS (fsck), and many app-level recovery schemes (e.g., Word)**

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.15

FFS: Create a File

Normal operation:

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

Recovery:

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

Time proportional to size of disk

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.16

Reliability Approach #2: Copy on Write File Layout

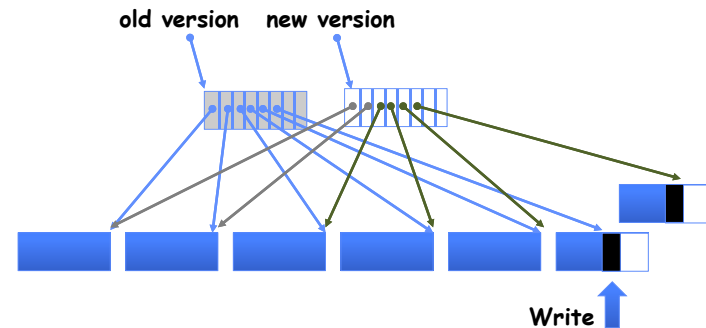
- To update file system, write a new version of the file system containing the update
 - Never update in place
 - Reuse existing unchanged disk blocks
- Seems expensive! But
 - Updates can be batched
 - Almost all disk writes can occur in parallel
- Approach taken in network file server appliances (WAFL, ZFS)

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.17

COW integrated with file system



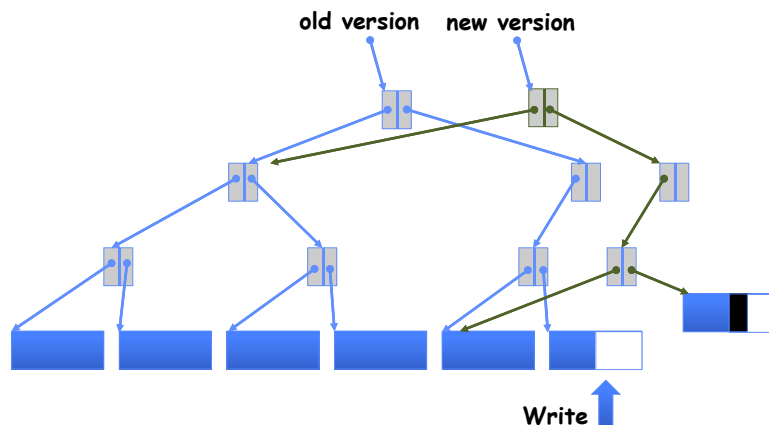
- If file represented as a tree of blocks, just need to update the leading fringe

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.18

COW with smaller-radix blocks



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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.19

ZFS

- 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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.20

More General Solutions

- 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 the many pieces
 - Many applications implement their own transactions
- Redundancy for media failures
 - Redundant representation (error correcting codes)
 - Replication
 - E.g., RAID disks

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.21

Transactions

- Closely related to critical sections in manipulating shared data structures
- Extend concept of atomic update from memory to stable storage
 - Atomically update multiple persistent data structures
- Like flags for threads, 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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.22

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



11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.23

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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.24

"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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.25

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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.26

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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.27

Logging 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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.28

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

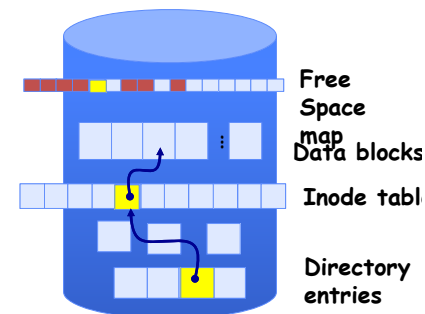
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.29

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



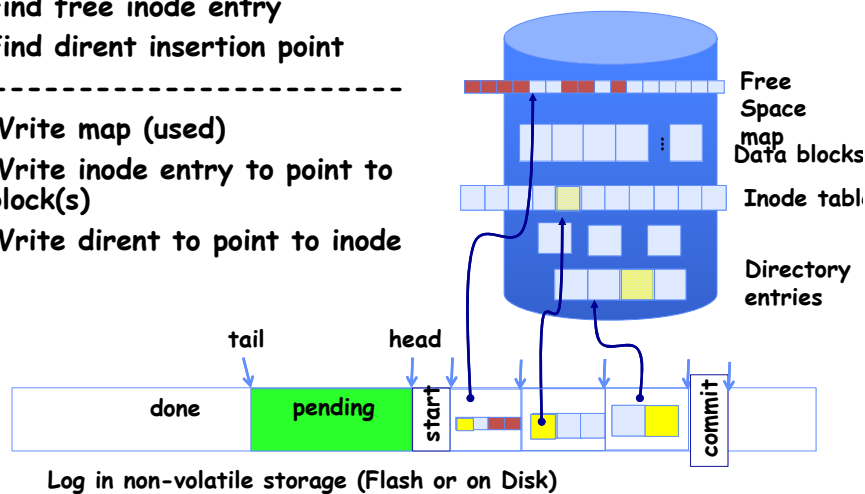
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.30

Ex: Creating a file (as a transaction)

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



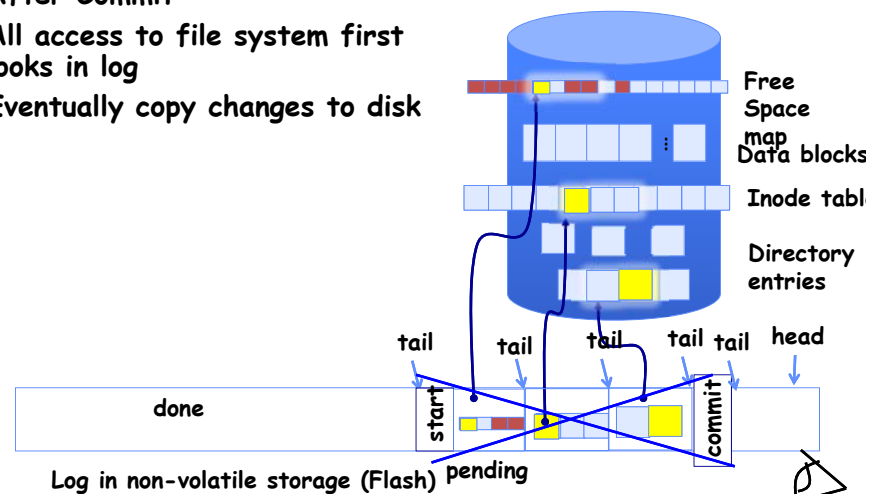
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.31

ReDo log

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



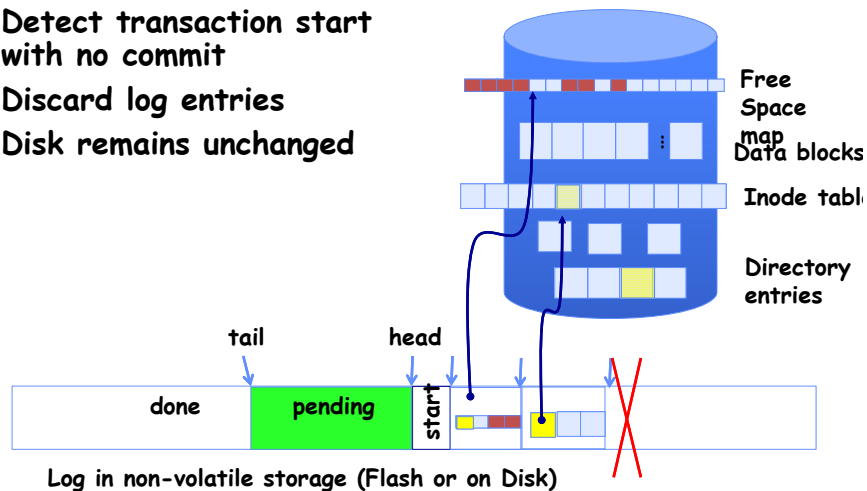
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.32

Crash during logging - Recover

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



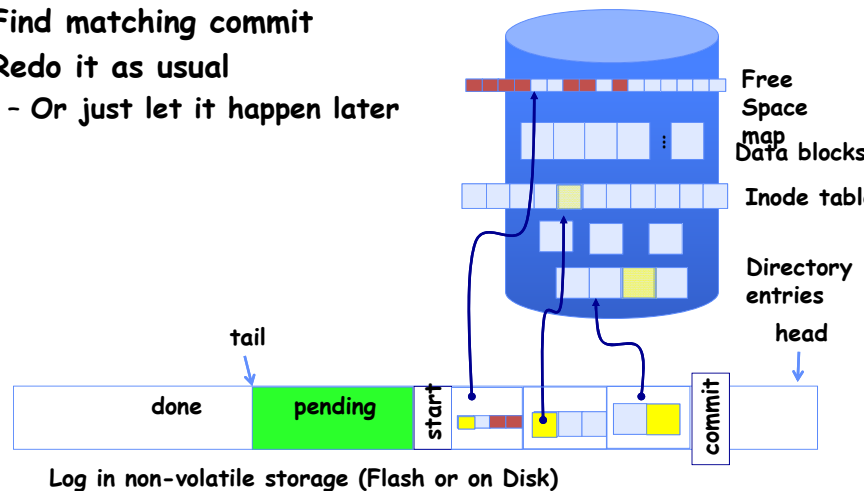
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.33

Recovery After Commit

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

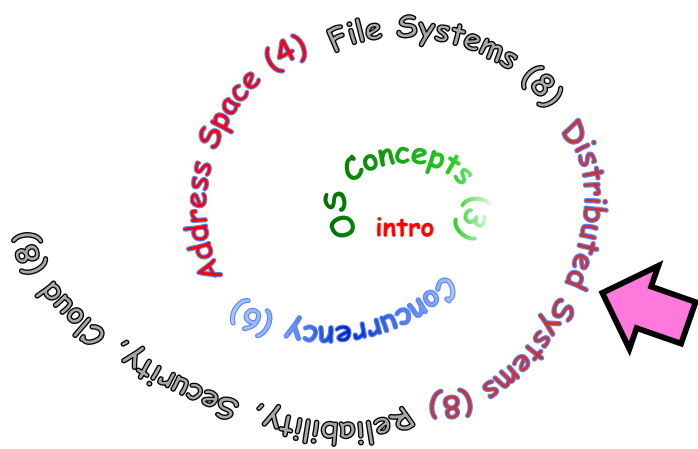


11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.34

Course Structure: Spiral



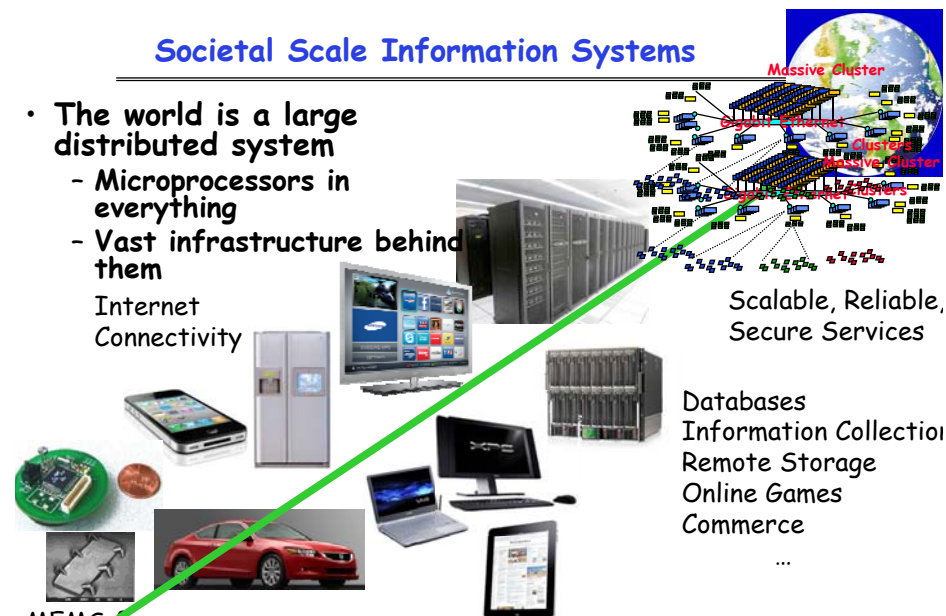
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.35

Societal Scale Information Systems

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



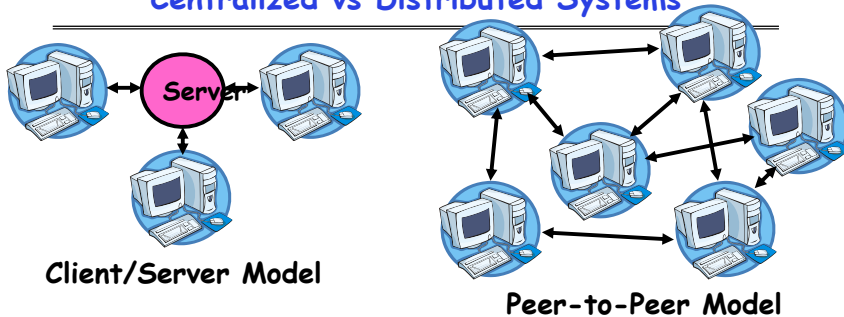
MEMS for Sensor Nets

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.36

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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.37

Distributed Systems: Motivation/Issues

- 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
- 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

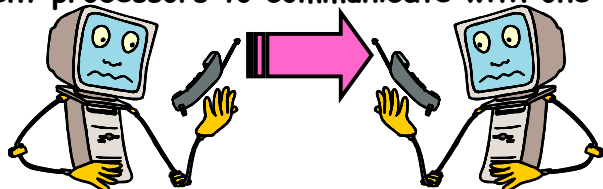
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.38

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 in the system
- Transparency and collaboration require some way for different processors to communicate with one another



11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.39

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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.40

Examples of Protocols in Human Interactions

• Telephone

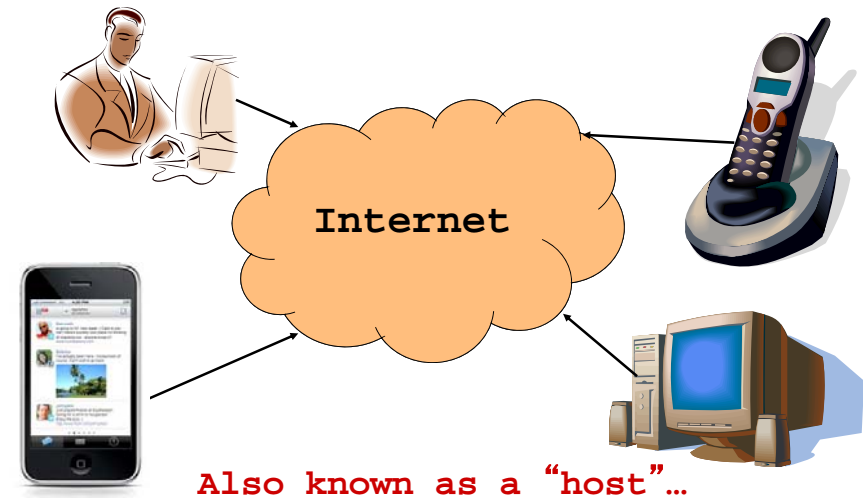
1. (Pick up / open up the phone)
2. Listen for a dial tone / see that you have service
3. Dial
4. Should hear ringing ...
5. **Callee: "Hello?"**
6. Caller: "Hi, it's John...."
Or: "Hi, it's me" (← what's *that* about?)
7. Caller: "Hey, do you think ... blah blah blah ..." pause
1. **Callee: "Yeah, blah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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.41

End System: Computer on the 'Net



11/9/15

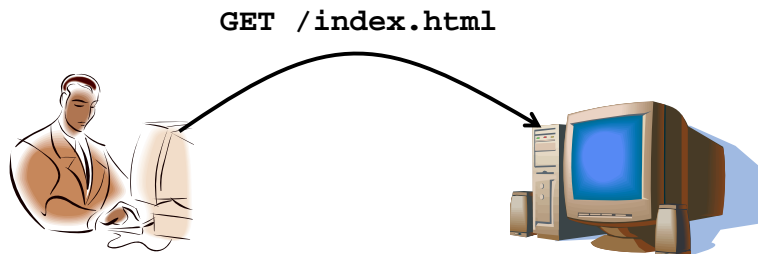
Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.42

Clients and Servers

• Client program

- Running on end host
- Requests service
- E.g., Web browser



11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.43

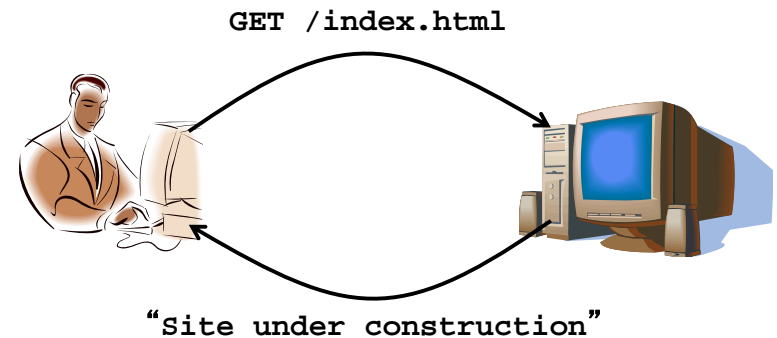
Clients and Servers

• Client program

- Running on end host
- Requests service
- E.g., Web browser

• Server program

- Running on end host
- Provides service
- E.g., Web server



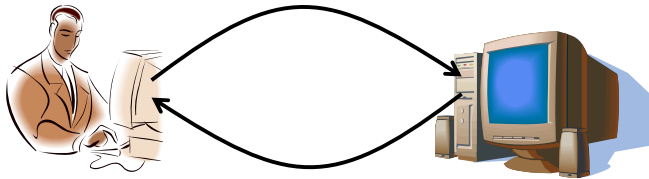
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.44

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



11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.45

Peer-to-Peer Communication

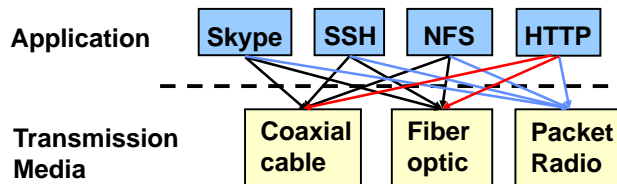
- No always-on server at the center of it all
 - Hosts can come and go, and change addresses
 - Hosts may have a different address each time
- Example: peer-to-peer file sharing (e.g., BitTorrent)
 - Any host can request files, send files, query to find where a file is located, respond to queries, and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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.46

Global Communication: The Problem



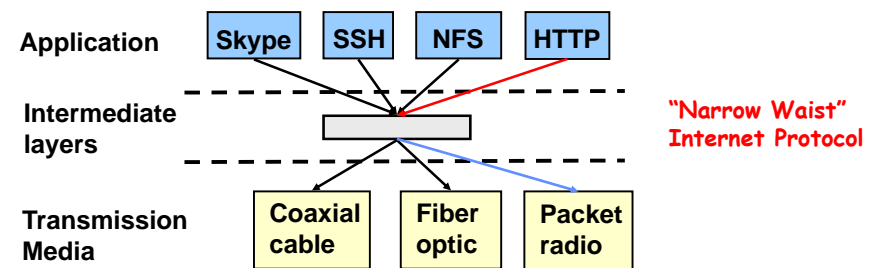
- Many different applications
 - email, web, P2P, etc.
- Many different network styles and technologies
 - Wireless vs. wired vs. optical, etc.
- How do we organize this mess?
 - Re-implement every application for every technology?
- No! But how does the Internet design avoid this?

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.47

Solution: Intermediate Layers



- Introduce intermediate layers that provide **set of abstractions** for various network functionality & technologies
 - A new app/media implemented only once
 - Variation on “add another level of indirection”

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.48

Distributed Applications

- How do you actually program a distributed application?
 - Need to synchronize multiple threads, running on different machines
 - » No shared memory, so cannot use test&set



- One Abstraction: send/receive messages
 - » Already atomic: no receiver gets portion of a message and two receivers cannot get same message
- Interface:
 - Mailbox (mbox): temporary holding area for messages
 - » Includes both destination location and queue
 - Send(message, mbox)
 - » Send message to remote mailbox identified by mbox
 - Receive(buffer, mbox)
 - » Wait until mbox has message, copy into buffer, and return
 - » If threads sleeping on this mbox, wake up one of them

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.49

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→T2
 - T1→buffer→T2
 - Very similar to producer/consumer
 - » Send = V, Receive = P
 - » However, can't tell if sender/receiver is local or not!

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.50

Messaging for Producer-Consumer Style

- Using send/receive for producer-consumer style:

```
Producer:
int msg1[1000];
while(1) {
    prepare message;
    send(msg1, mbox);
}
```

Send Message

```
Consumer:
int buffer[1000];
while(1) {
    receive(buffer, mbox);
    process message;
}
```

Receive Message

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

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.51

Messaging for Request/Response communication

- What about two-way communication?
 - Request/Response
 - » Read a file stored on a remote machine
 - » Request a web page from a remote web server
 - Also called: **client-server**
 - » Client ≡ requester, Server ≡ responder
 - » Server provides "service" (file storage) to the client
- Example: File service

```
Client: (requesting the file)
char response[1000];

send("read rutabaga", server_mbox);
receive(response, client_mbox);
```

Request File

Get Response

```
Server: (responding with the file)
char command[1000], answer[1000];

receive(command, server_mbox);
decode command;
read file into answer;
send(answer, client_mbox);
```

Receive Request

Send Response

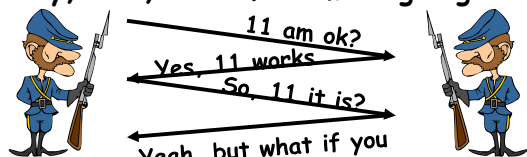
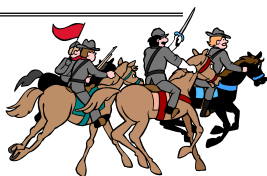
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.52

General's Paradox

- **General's paradox:**
 - Constraints of problem:
 - » Two generals, on separate mountains
 - » Can only communicate via messengers
 - » Messengers can be captured
 - Problem: need to coordinate attack
 - » If they attack at different times, they all die
 - » If they attack at same time, they win
 - Named after Custer, who died at Little Big Horn because he arrived a couple of days too early
- Can messages over an unreliable network be used to guarantee two entities do something simultaneously?
 - Remarkably, "no", even if all messages get through



- No way to be sure last message gets through!

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.53

Two-Phase Commit

- Since we can't solve the General's Paradox (i.e. simultaneous action), let's solve a related problem
 - Distributed transaction: Two machines agree to do something, or not do it, **atomically**
- Two-Phase Commit protocol does this
 - **Persistent stable log on each machine:** keep track of whether commit has happened
 - » If a machine crashes, when it wakes up it first checks its log to recover state of world at time of crash
 - **Prepare Phase:**
 - » The global coordinator requests that all participants will promise to commit or rollback the transaction
 - » Participants record promise in log, then acknowledge
 - » If anyone votes to abort, coordinator writes "Abort" in its log and tells everyone to abort; each records "Abort" in log
 - **Commit Phase:**
 - » After all participants respond that they are prepared, then the coordinator writes "Commit" to its log
 - » Then asks all nodes to commit; they respond with ack
 - » After receive acks, coordinator writes "Got Commit" to log
- Log can be used to complete this process such that all machines either commit or don't commit

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.54

2PC Algorithm

- Developed by Turing award winner Jim Gray (first Berkeley CS PhD, 1969)
- One coordinator
- N workers (replicas)
- High level algorithm description
 - Coordinator asks all workers if they can commit
 - If all workers reply "**VOTE-COMMIT**", then coordinator broadcasts "**GLOBAL-COMMIT**",
Otherwise coordinator broadcasts "**GLOBAL-ABORT**"
 - Workers obey the **GLOBAL** messages
- Use a persistent, stable log on each machine to keep track of what you are doing
 - If a machine crashes, when it wakes up it first checks its log to recover state of world at time of crash

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.55

Detailed Algorithm

Coordinator Algorithm

Coordinator sends **VOTE-REQ** to all workers

- If receive **VOTE-COMMIT** from all N workers, send **GLOBAL-COMMIT** to all workers
- If doesn't receive **VOTE-COMMIT** from all N workers, send **GLOBAL-ABORT** to all workers

Worker Algorithm

- Wait for **VOTE-REQ** from coordinator
- If ready, send **VOTE-COMMIT** to coordinator
- If not ready, send **VOTE-ABORT** to coordinator
 - And immediately abort

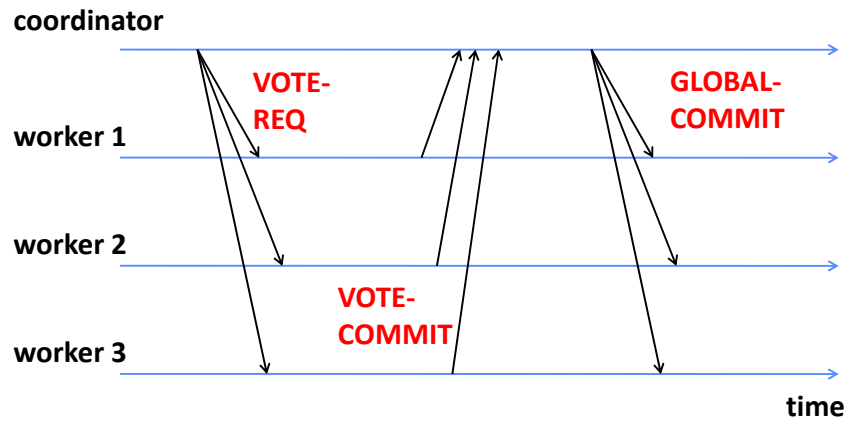
- If receive **GLOBAL-COMMIT** then commit
- If receive **GLOBAL-ABORT** then abort

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.56

Failure Free Example Execution



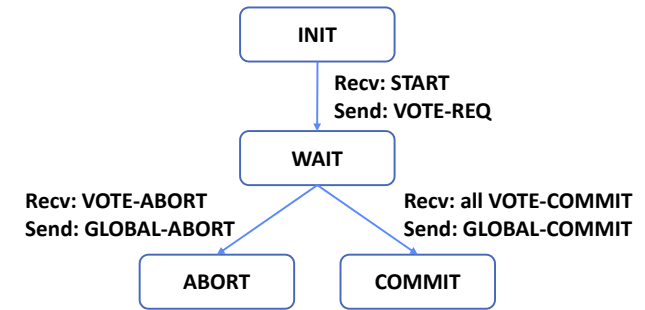
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.57

State Machine of Coordinator

- Coordinator implements simple state machine:

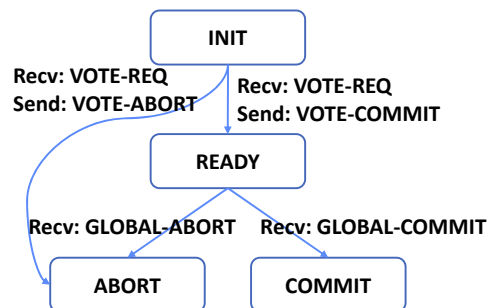


11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.58

State Machine of Workers



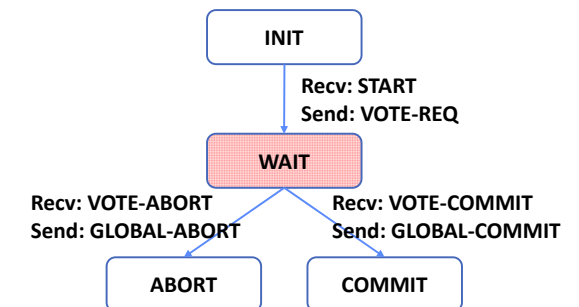
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.59

Dealing with Worker Failures

- How to deal with worker failures?
 - Failure only affects states in which the node is waiting for messages
 - Coordinator only waits for votes in "WAIT" state
 - In WAIT, if doesn't receive N votes, it times out and sends GLOBAL-ABORT

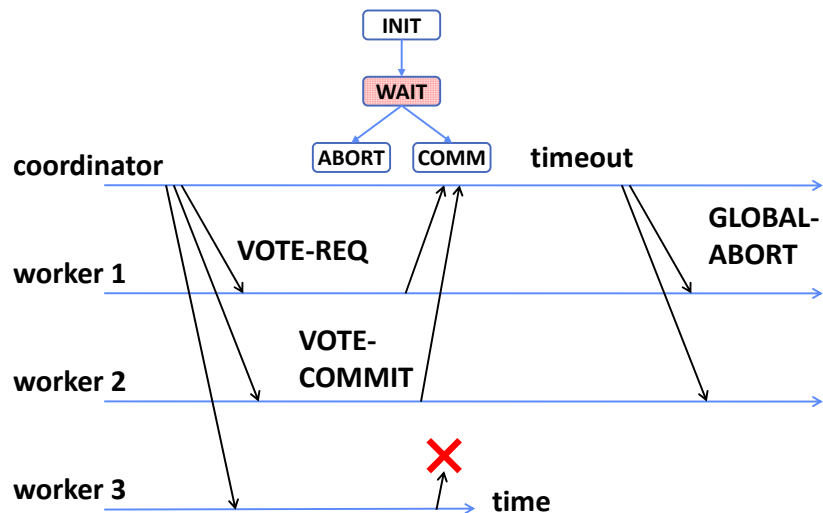


11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.60

Example of Worker Failure



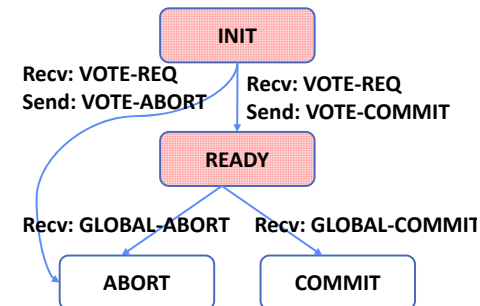
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.61

Dealing with Coordinator Failure

- How to deal with coordinator failures?
 - worker waits for VOTE-REQ in INIT
 - » Worker can time out and abort (coordinator handles it)
 - worker waits for GLOBAL-* message in READY
 - » If coordinator fails, workers must **BLOCK** waiting for coordinator to recover and send GLOBAL_* message

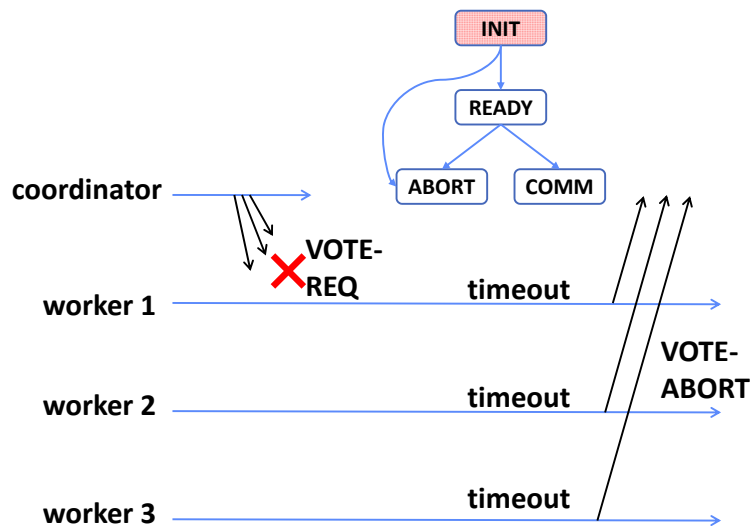


11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.62

Example of Coordinator Failure #1

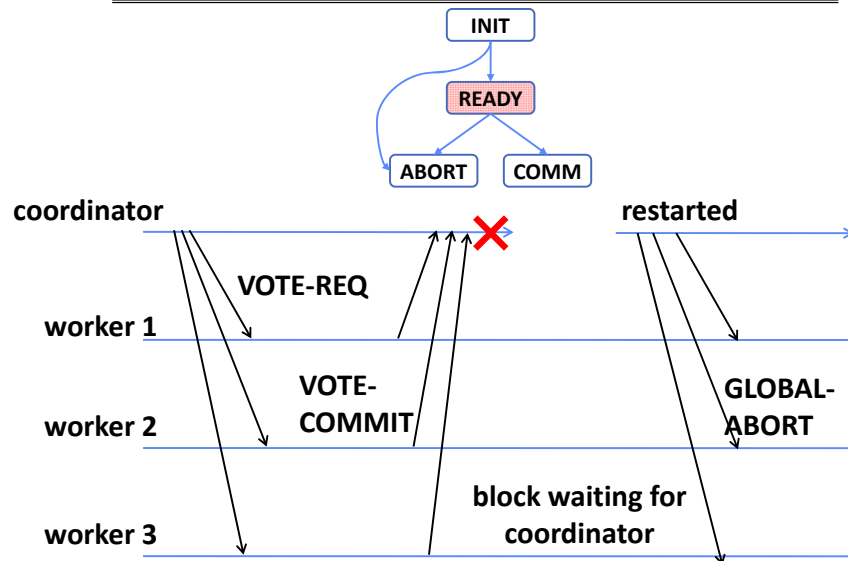


11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.63

Example of Coordinator Failure #2



11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.64

Durability

- All nodes use stable storage* to store which state they are in
- Upon recovery, it can restore state and resume:
 - Coordinator aborts in INIT, WAIT, or ABORT
 - Coordinator commits in COMMIT
 - Worker aborts in INIT, ABORT
 - Worker commits in COMMIT
 - Worker asks Coordinator in READY
- * - stable storage is non-volatile storage (e.g. backed by disk) that guarantees atomic writes.

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.65

Blocking for Coordinator to Recover

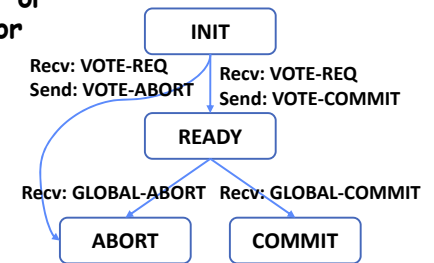
- A worker waiting for global decision can ask fellow workers about their state

- If another worker is in ABORT or COMMIT state then coordinator must have sent GLOBAL-*

» Thus, worker can safely abort or commit, respectively

- If another worker is still in INIT state then both workers can decide to abort

- If all workers are in ready, need to **BLOCK** (don't know if coordinator wanted to abort or commit)



11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.66

Distributed Decision Making Discussion

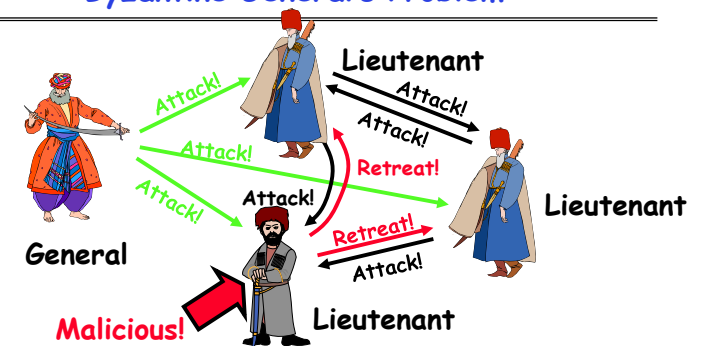
- Why is distributed decision making desirable?
 - Fault Tolerance!
 - A group of machines can come to a decision even if one or more of them fail during the process
 - » Simple failure mode called "failstop" (different modes later)
 - After decision made, result recorded in multiple places
- Undesirable feature of Two-Phase Commit: Blocking
 - One machine can be stalled until another site recovers:
 - » Site B writes "prepared to commit" record to its log, sends a "yes" vote to the coordinator (site A) and crashes
 - » Site A crashes
 - » Site B wakes up, check its log, and realizes that it has voted "yes" on the update. It sends a message to site A asking what happened. At this point, B cannot decide to abort, because update may have committed
 - » B is blocked until A comes back
 - A blocked site holds resources (locks on updated items, pages pinned in memory, etc) until learns fate of update
- **PAXOS**: An alternative used by GOOGLE and others that does not have this blocking problem
- What happens if one or more of the nodes is malicious?
 - **Malicious**: attempting to compromise the decision making

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.67

Byzantine General's Problem



- Byzantine General's Problem (n players):
 - One General
 - n-1 Lieutenants
 - Some number of these (f) can be insane or malicious
- The commanding general must send an order to his n-1 lieutenants such that:
 - IC1: All loyal lieutenants obey the same order
 - IC2: If the commanding general is loyal, then all loyal lieutenants obey the order he sends

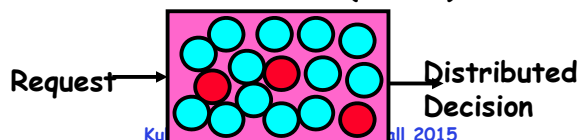
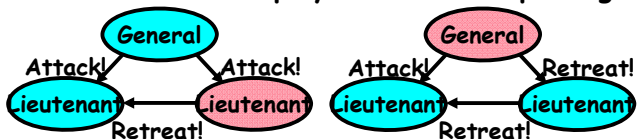
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.68

Byzantine General's Problem (con't)

- **Impossibility Results:**
 - Cannot solve Byzantine General's Problem with $n=3$ because one malicious player can mess up things
- With f faults, need $n > 3f$ to solve problem
- Various algorithms exist to solve problem
 - Original algorithm has #messages exponential in n
 - Newer algorithms have message complexity $O(n^2)$
 - » One from MIT, for instance (Castro and Liskov, 1999)
- Use of BFT (Byzantine Fault Tolerance) algorithm
 - Allow multiple machines to make a coordinated decision even if some subset of them ($< n/3$) are malicious



11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.69

Remote Procedure Call

- Raw messaging is a bit too low-level for programming
 - Must wrap up information into message at source
 - Must decide what to do with message at destination
 - May need to sit and wait for multiple messages to arrive
- Better option: Remote Procedure Call (RPC)
 - Calls a procedure on a remote machine
 - Client calls:


```
remoteFileSystem→Read("rutabaga");
```
 - Translated automatically into call on server:

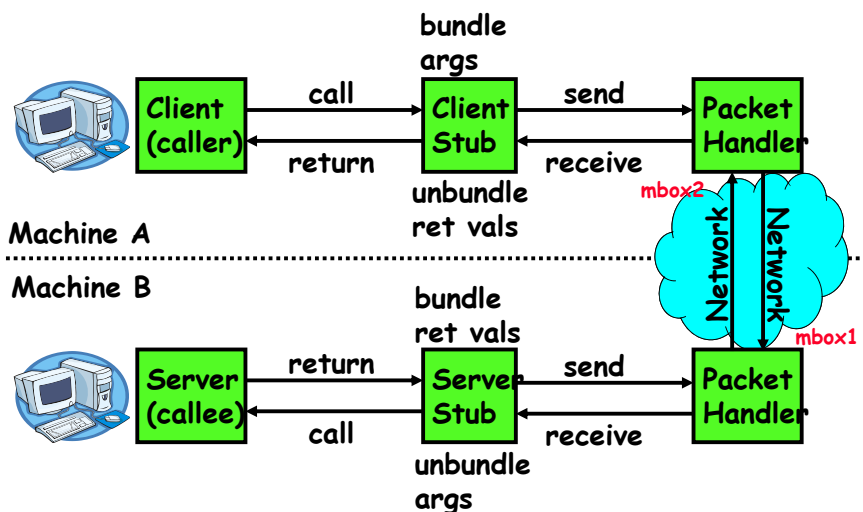

```
fileSys→Read("rutabaga");
```
- Implementation:
 - Request-response message passing (under covers!)
 - "Stub" provides glue on client/server
 - » Client stub is responsible for "marshalling" arguments and "unmarshalling" the return values
 - » Server-side stub is responsible for "unmarshalling" arguments and "marshalling" the return values.
- **Marshalling** involves (depending on system)
 - Converting values to a canonical form, serializing objects, copying arguments passed by reference, etc.

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.70

RPC Information Flow



11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.71

RPC Details

- Equivalence with regular procedure call
 - Parameters \leftrightarrow Request Message
 - Result \leftrightarrow Reply message
 - Name of Procedure: Passed in request message
 - Return Address: mbox2 (client return mail box)
- Stub generator: Compiler that generates stubs
 - Input: interface definitions in an "interface definition language (IDL)"
 - » Contains, among other things, types of arguments/return
 - Output: stub code in the appropriate source language
 - » Code for client to pack message, send it off, wait for result, unpack result and return to caller
 - » Code for server to unpack message, call procedure, pack results, send them off
- Cross-platform issues:
 - What if client/server machines are different architectures or in different languages?
 - » Convert everything to/from some canonical form
 - » Tag every item with an indication of how it is encoded (avoids unnecessary conversions).

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.72

RPC Details (continued)

- How does client know which mbox to send to?
 - Need to translate name of remote service into network endpoint (Remote machine, port, possibly other info)
 - **Binding**: the process of converting a user-visible name into a network endpoint
 - » This is another word for "naming" at network level
 - » Static: fixed at compile time
 - » Dynamic: performed at runtime
- Dynamic Binding
 - Most RPC systems use dynamic binding via name service
 - » Name service provides dynamic translation of service→mbox
 - Why dynamic binding?
 - » Access control: check who is permitted to access service
 - » Fail-over: If server fails, use a different one
- What if there are multiple servers?
 - Could give flexibility at binding time
 - » Choose unloaded server for each new client
 - Could provide same mbox (router level redirect)
 - » Choose unloaded server for each new request
 - » Only works if no state carried from one call to next
- What if multiple clients?
 - Pass pointer to client-specific return mbox in request

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.73

Problems with RPC

- Non-Atomic failures
 - Different failure modes in distributed system than on a single machine
 - Consider many different types of failures
 - » User-level bug causes address space to crash
 - » Machine failure, kernel bug causes all processes on same machine to fail
 - » Some machine is compromised by malicious party
 - Before RPC: whole system would crash/die
 - After RPC: One machine crashes/compromised while others keep working
 - Can easily result in inconsistent view of the world
 - » Did my cached data get written back or not?
 - » Did server do what I requested or not?
 - Answer? Distributed transactions/Byzantine Commit
- Performance
 - Cost of Procedure call « same-machine RPC « network RPC
 - Means programmers must be aware that RPC is not free
 - » Caching can help, but may make failure handling complex

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.74

Cross-Domain Communication/Location Transparency

- How do address spaces communicate with one another?
 - Shared Memory with Semaphores, monitors, etc...
 - File System
 - Pipes (1-way communication)
 - "Remote" procedure call (2-way communication)
- RPC's can be used to communicate between address spaces on different machines or the same machine
 - Services can be run wherever it's most appropriate
 - Access to local and remote services looks the same
- Examples of modern RPC systems:
 - CORBA (Common Object Request Broker Architecture)
 - DCOM (Distributed COM)
 - RMI (Java Remote Method Invocation)

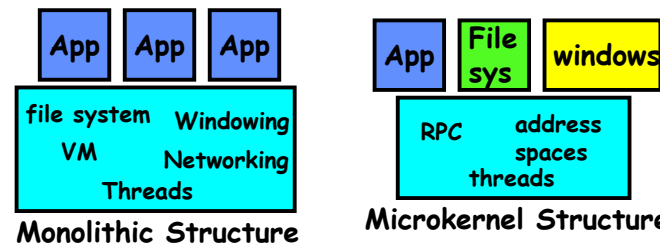
11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.75

Microkernel operating systems

- Example: split kernel into application-level servers.
 - File system looks remote, even though on same machine



- Why split the OS into separate domains?
 - Fault isolation: bugs are more isolated (build a firewall)
 - Enforces modularity: allows incremental upgrades of pieces of software (client or server)
 - Location transparent: service can be local or remote
 - » For example in the X windowing system: Each X client can be on a separate machine from X server; Neither has to run on the machine with the frame buffer.

11/9/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 20.76

Summary (1/2)

- **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

Summary (2/2)

- **Two-phase commit**: distributed decision making
 - First, make sure everyone guarantees that they will commit if asked (prepare)
 - Next, ask everyone to commit
- **Byzantine General's Problem**: distributed decision making with malicious failures
 - One general, $n-1$ lieutenants: some number of them may be malicious (often "f" of them)
 - All non-malicious lieutenants must come to same decision
 - If general not malicious, lieutenants must follow general
 - Only solvable if $n \geq 3f+1$
- **Remote Procedure Call (RPC)**: Call procedure on remote machine
 - Provides same interface as procedure
 - Automatic packing and unpacking of arguments without user programming (in stub)