

CSI 62
Operating Systems and
Systems Programming
Lecture 20

Reliability, Transactions
Distributed Systems

April 11th, 2016
Prof. Anthony D. Joseph
<http://cs162.eecs.Berkeley.edu>

Recall: 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

4/11/16

Joseph CSI 62 ©UCB Spring 2016

Lec 20.2

RAID: Redundant Arrays of Inexpensive Disks

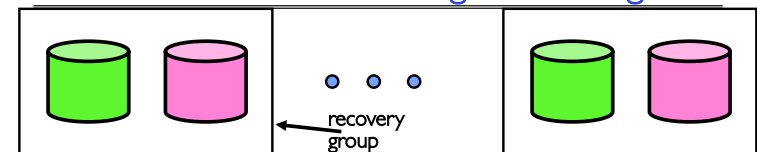
- Invented by David Patterson, Garth A. Gibson, and Randy Katz here at UCB in 1987
- Data stored on multiple disks (redundancy)
- Either in software or hardware
 - In hardware case, done by disk controller; file system may not even know that there is more than one disk in use
- Initially, five levels of RAID (more now)

4/11/16

Joseph CSI 62 ©UCB Spring 2016

Lec 20.3

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

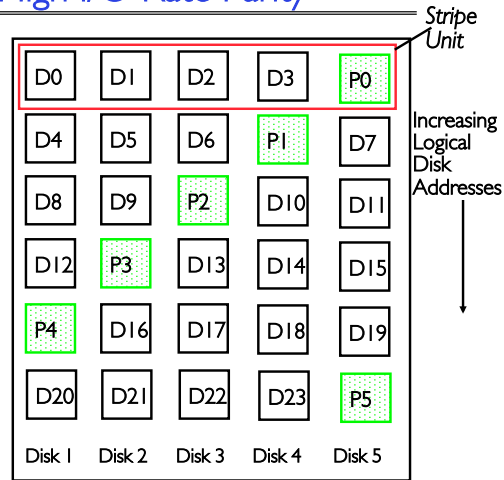
4/11/16

Joseph CSI 62 ©UCB Spring 2016

Lec 20.4

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: $D_3 = D_0 \oplus D_1 \oplus D_2 \oplus P_0$



- Can spread information widely across internet for durability
 - Overview now, more later in semester

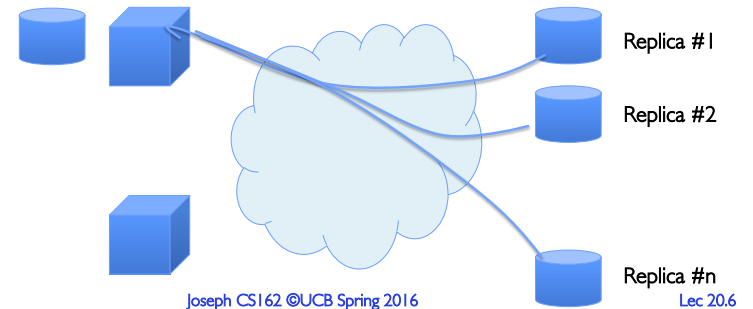
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.5

Higher Durability/Reliability through Geographic Replication

- Highly durable – hard to destroy all copies
- Highly available for reads – read any copy
- 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



4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.6

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
 - 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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.7

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?

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.8

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

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)

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 → inode number
- Update modify time for directory

Recovery:

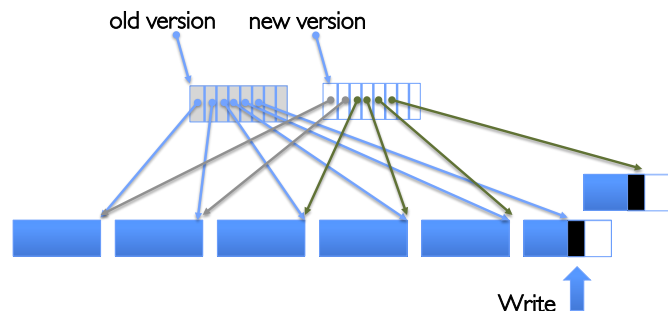
- Scan inode table
- If any unlinked files (not in any directory), delete
- Compare free block bitmap against inode trees
- Scan directories for missing update/access times

Time proportional to size of disk

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 and OpenZFS

COW Integrated with File System



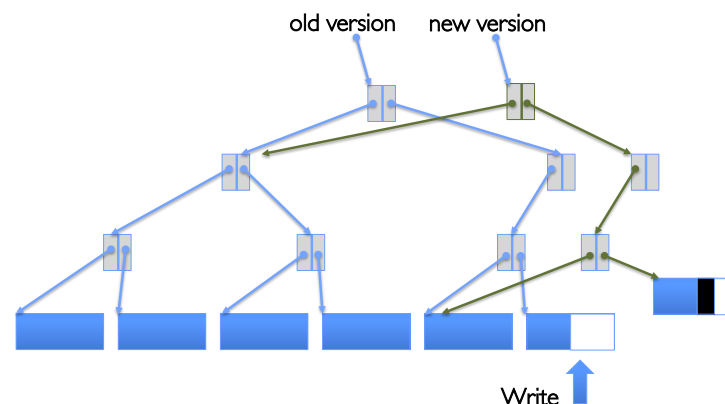
- If file represented as a tree of blocks, just need to update the leading fringe

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.13

COW with Smaller-Radix Blocks



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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.14

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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.15

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)

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.16

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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.17

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



4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.18

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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.19

“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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.20

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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.21

Administrivia

- Prof Joseph's office hours this week: Wed 4/13 10-noon 465 Soda
- **Midterm II: Next week! (4/20)**
 - 6-7:30PM (aa-eh 10 Evans, ej-oa 155 Dwinelle)
 - Covers lectures #13 to 21 (assumes knowledge of #1 – 12)
 - 1 page of hand-written notes, both sides
 - **Review session Mon 4/18, 6:30-8:00 PM in 245 Li Ka Shing**
- Project 2 code due today
- HW4 – Releases today (Due 4/25)
- Project 3 releases tomorrow (Design doc due 4/22)

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.22

BREAK

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.23

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 (can be accelerated with NVRAM)
 - File system may not be updated immediately, data preserved in the log
- Difference between “Log Structured” and “Journaling”
 - In a Log Structured filesystem, data stays in log form
 - In a Journaling 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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.24

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 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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.25

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

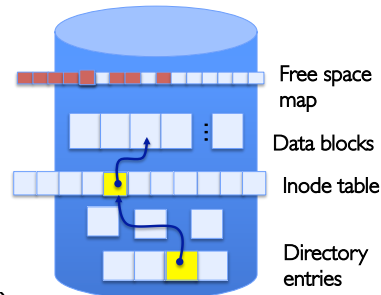
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.26

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



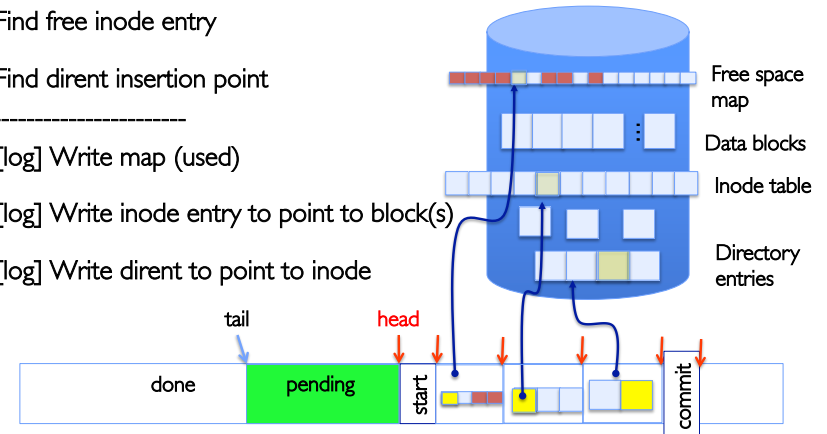
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.27

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)

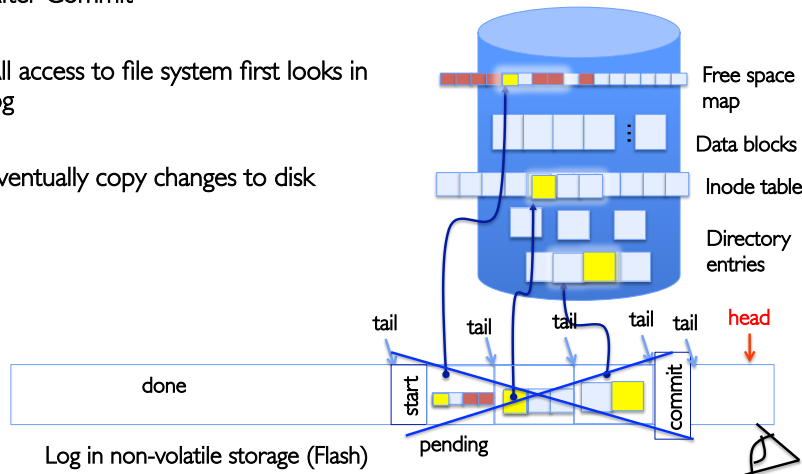
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.28

ReDo Log

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



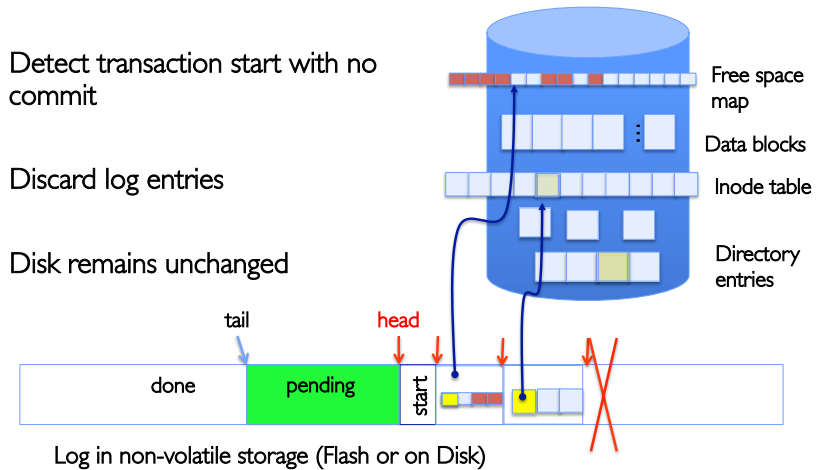
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.29

Crash During Logging – Recover

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



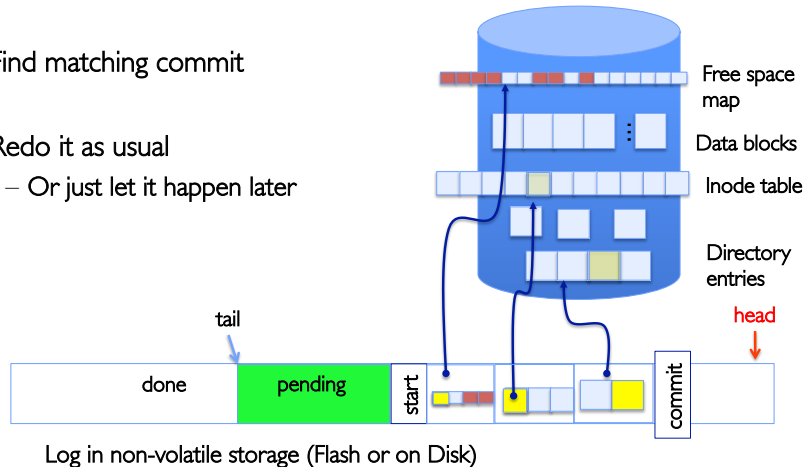
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.30

Recovery After Commit

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

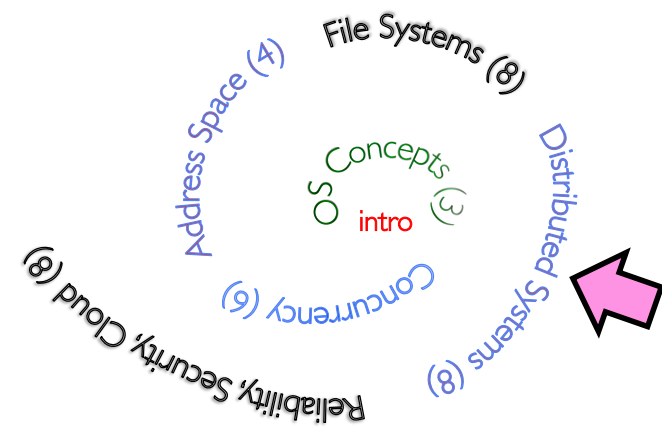


4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.31

Course Structure: Spiral



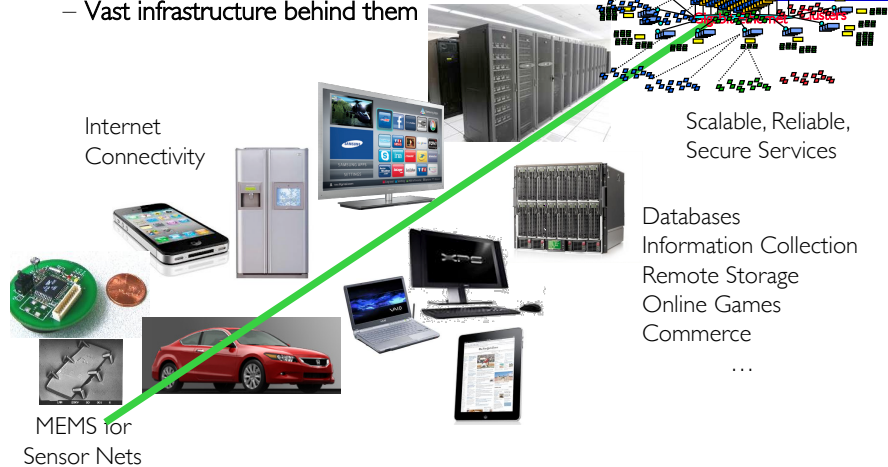
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.32

Societal Scale Information Systems

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

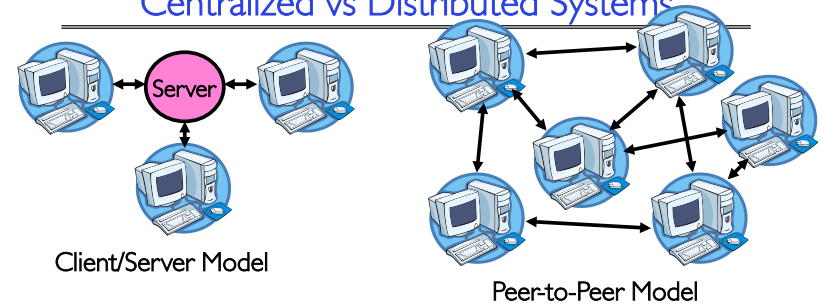


4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.33

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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.34

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

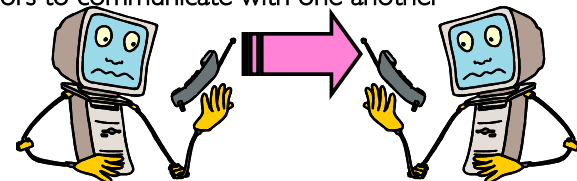
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.35

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



4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.36

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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.37

Examples of Protocols in Human Interactions

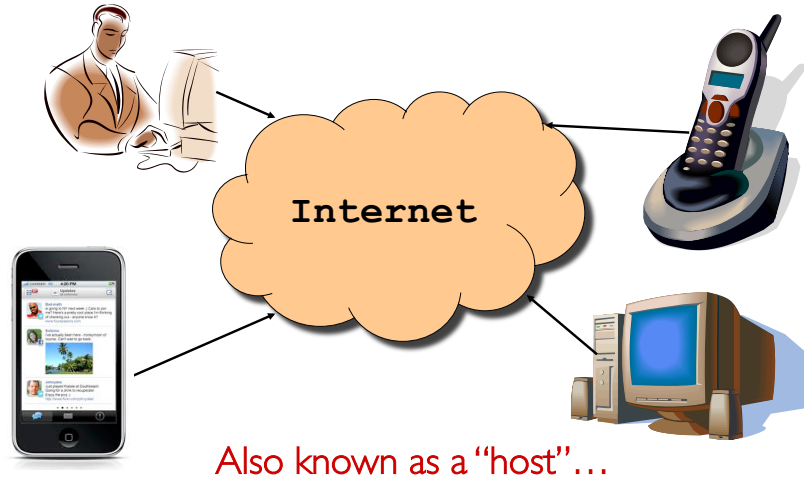
- Telephone
 1. (Pick up / open up the phone)
 2. Listen for a dial tone / see that you have service
 3. Dial
 4. Should hear ringing ...
 - 5.
 6. Caller: "Hi, it's Anthony...."
Or: "Hi, it's me" (← what's *that* about?)
 7. Caller: "Hey, do you think ... blah blah blah ..." pause
 - 8.
 9. Caller: "Bye"
 - 10.
 11. Hang up
- Diagram illustrating the sequence of events in a telephone conversation:
- Caller: "Hi, it's Anthony...."
Or: "Hi, it's me" (← what's *that* about?)
 - Caller: "Hey, do you think ... blah blah blah ..." pause
 - Caller: "Bye"
 - Hang up
 - Callee: "Hello?"
 - Callee: "Yeah, blah blah blah ..." pause
 - Callee: "Bye"

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.38

End System: Computer on the 'Net



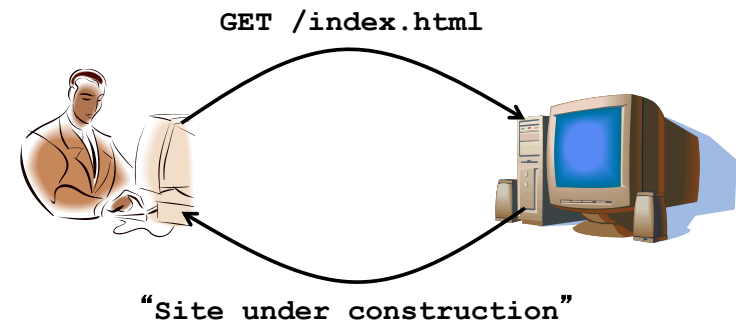
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.39

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



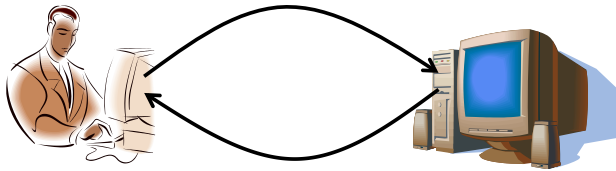
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.40

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



4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.41

Peer-to-Peer Communication

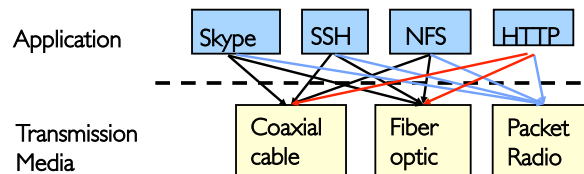
- No always-on server at the center of it all
 - Hosts can come and go, and change addresses
 - Hosts may have a different address each time
- Example: peer-to-peer file sharing (e.g., BitTorrent)
 - Any host can request files, send files, query to find where a file is located, respond to queries, and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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.42

Global Communication: The Problem



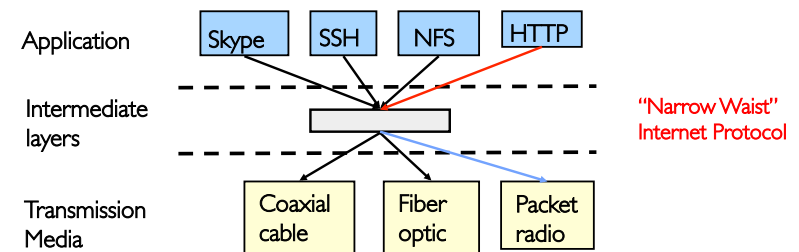
- Many different applications
 - email, web, P2P, etc.
- Many different network styles and technologies
 - Wireless vs. wired vs. optical, etc.
- How do we organize this mess?
 - Re-implement every application for every technology?
- No! But how does the Internet design avoid this?

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.43

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”
- **Goal: Reliable communication channels on which to build distributed applications**

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.44

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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.45

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!

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.46

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
 - Next time: will discuss fact that this is 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

4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.47

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];
```

Request
File

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

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

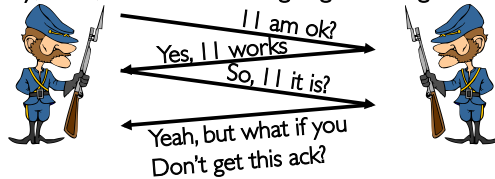
4/11/16

Joseph CS162 ©UCB Spring 2016

Lec 20.48

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!



Summary

- **RAID**: Redundant Arrays of Inexpensive Disks
 - RAID 1: mirroring, RAID 5: Parity block
- Use of Log to improve Reliability
 - Journaling file systems such as ext3, NTFS
- **Transactions**: ACID semantics
 - Atomicity
 - Consistency
 - Isolation
 - Durability
- Protocol: Agreement between two parties as to how information is to be transmitted