

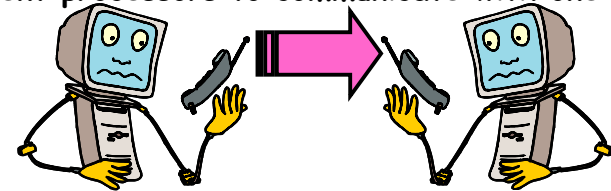
CS162 Operating Systems and Systems Programming Lecture 21

Distributed Decision Making, TCP/IP and Sockets

November 16th, 2015
Prof. John Kubiatowicz
<http://cs162.eecs.Berkeley.edu>

Recall: 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



11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.2

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

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.3

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

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.4

Examples of Protocols in Human Interactions

• Telephone

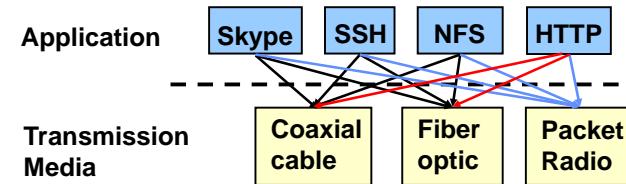
1. (Pick up / open up the phone)
2. Listen for a dial tone / see that you have service
3. Dial
4. Should hear ringing ...
5. Callee: "Hello?"
6. Caller: "Hi, it's John...."
Or: "Hi, it's me" (← what's *that* about?)
7. Caller: "Hey, do you think ... blah blah blah ..." pause
1. Callee: "Yeah, blah blah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

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.5

Global Communication: The Problem



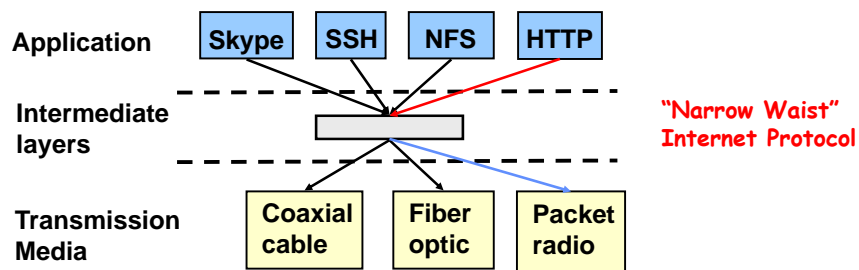
- Many different applications
 - email, web, P2P, etc.
- Many different network styles and technologies
 - Wireless vs. wired vs. optical, etc.
- How do we organize this mess?
 - Re-implement every application for every technology?
- No! But how does the Internet design avoid this?

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.6

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**

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.7

Recall: Use of TCP: Sockets

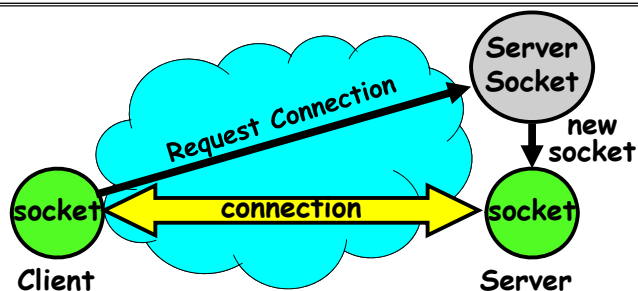
- **Socket**: an abstraction of a network I/O queue
 - Embodies one side of a communication channel
 - » Same interface regardless of location of other end
 - » Could be local machine (called "UNIX socket") or remote machine (called "network socket")
 - First introduced in 4.2 BSD UNIX: big innovation at time
 - » Now most operating systems provide some notion of socket
- Using Sockets for Client-Server (C/C++ interface):
 - On server: set up "server-socket"
 - » Create socket, Bind to protocol (TCP), local address, port
 - » Call listen(): tells server socket to accept incoming requests
 - » Perform multiple accept() calls on socket to accept incoming connection request
 - » Each successful accept() returns a new socket for a new connection; can pass this off to handler thread
 - On client:
 - » Create socket, Bind to protocol (TCP), remote address, port
 - » Perform connect() on socket to make connection
 - » If connect() successful, have socket connected to server

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.8

Recall: Socket Setup over TCP/IP



- **Server Socket:** Listens for new connections
 - Produces new sockets for each unique connection
- **Things to remember:**
 - Connection involves 5 values:
[Client Addr, Client Port, Server Addr, Server Port, Protocol]
 - Often, Client Port "randomly" assigned
 - » Done by OS during client socket setup
 - Server Port often "well known"
 - » 80 (web), 443 (secure web), 25 (sendmail), etc
 - » Well-known ports from 0–1023

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.9

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

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.10

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!

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.11

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 the window in TCP: window is size of buffer on far end
 - Restricts sender to forward only what will fit in buffer

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.12

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 \equiv requester, Server \equiv 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);
```

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

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

Request File

Get Response

Receive Request

Send Response

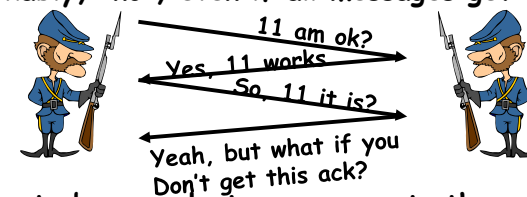
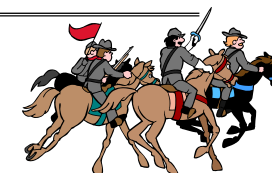
11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.13

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!

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.14

Administrivia

- Midterm II: Next Monday (11/23)
 - Time: 7:00PM-10:00PM
 - Location: VLSB: 2040/2050/2060
 - Division of people between rooms will be posted
 - All topics from Midterm I, up to this Wednesday, including:
 - » Address Translation/TLBs/Paging
 - » I/O subsystems, Storage Layers, Disks/SSD
 - » Performance and Queueing Theory
 - » File systems
 - » Distributed systems, TCP/IP, RPC
- Closed book, one page of notes - both sides
- Review session:
 - 306 Soda Hall
 - Sunday TBA

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.15

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:
 - **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

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.16

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

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.17

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

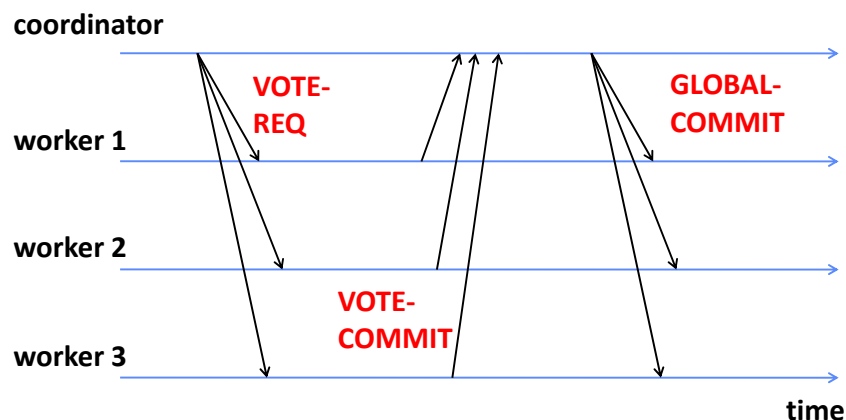
- If receive GLOBAL-COMMIT then commit
- If receive GLOBAL-ABORT then abort

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.18

Failure Free Example Execution



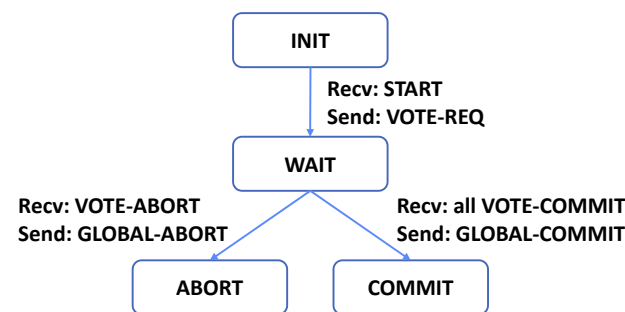
11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.19

State Machine of Coordinator

- Coordinator implements simple state machine:

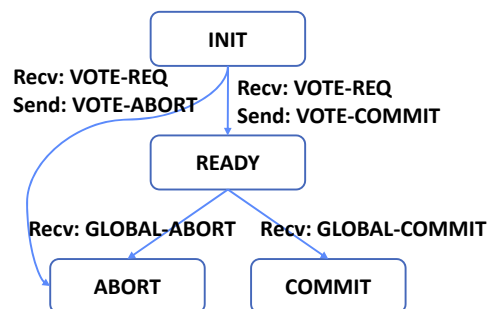


11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.20

State Machine of Workers



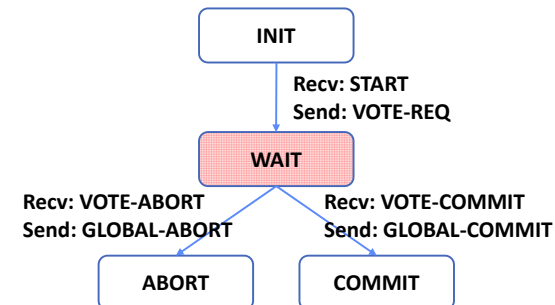
11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.21

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

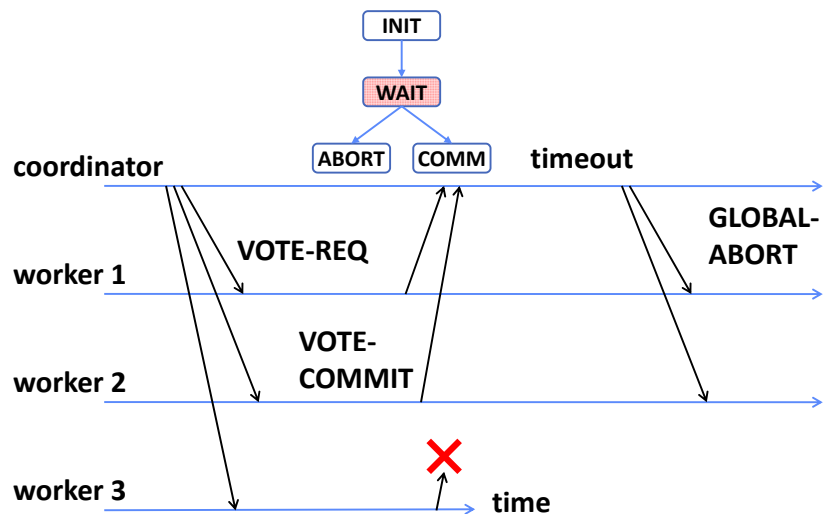


11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.22

Example of Worker Failure



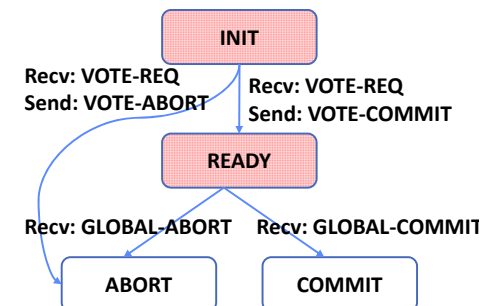
11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.23

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

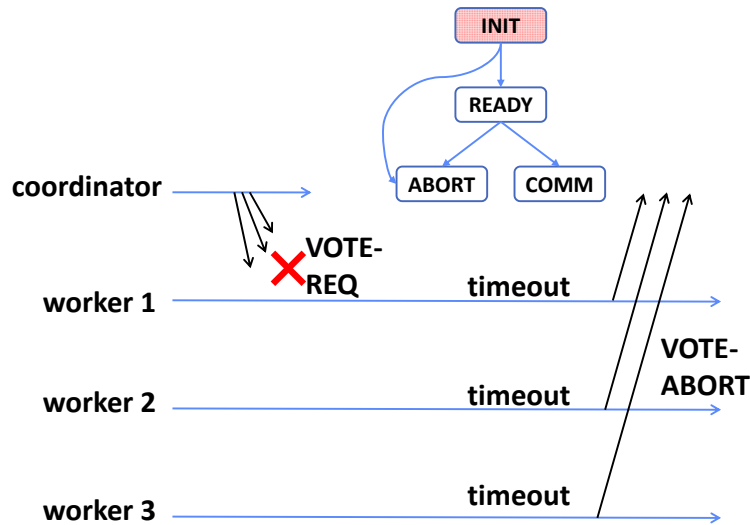


11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.24

Example of Coordinator Failure #1

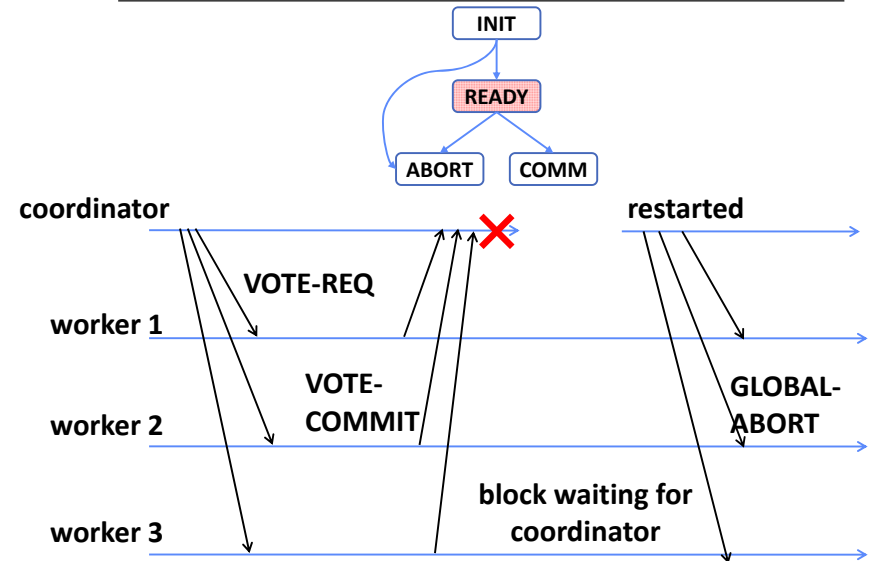


11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.25

Example of Coordinator Failure #2



11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.26

Durability

- All nodes use **stable storage** to store current state
 - stable storage is non-volatile storage (e.g. backed by disk) that guarantees atomic writes.
- 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

11/16/15

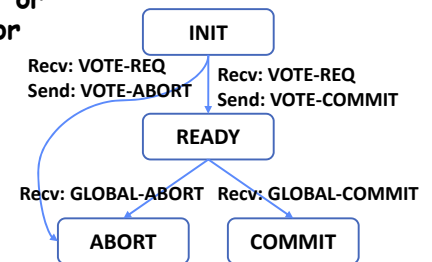
Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.27

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)

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.28

Distributed Decision Making Discussion

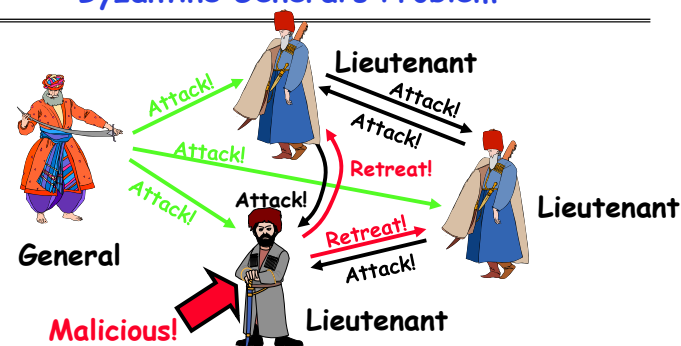
- Why is distributed decision making desirable?
 - Fault Tolerance!
 - A group of machines can come to a decision even if one or more of them fail during the process
 - » Simple failure mode called "failstop" (different modes later)
 - After decision made, result recorded in multiple places
- Undesirable 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

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.29

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

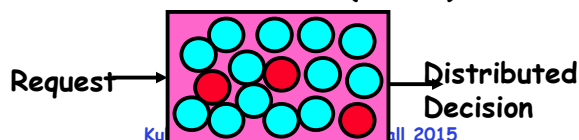
11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.30

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



11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.31

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
- Another option: Remote Procedure Call (RPC)
 - Calls a procedure on a remote machine
 - Client calls:


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

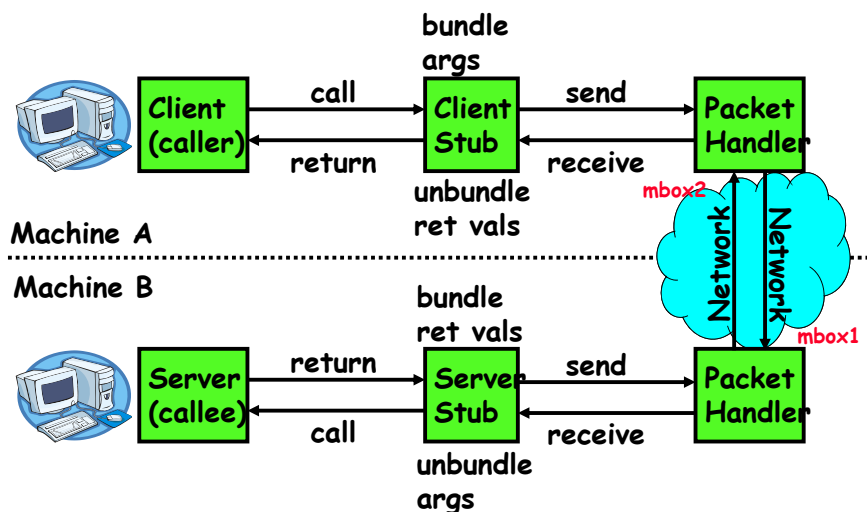

```
fileSys→Read("rutabaga");
```
- Implementation:
 - Request-response message passing (under covers!)
 - "Stub" provides glue on client/server
 - » Client stub is responsible for "marshalling" arguments and "unmarshalling" the return values
 - » Server-side stub is responsible for "unmarshalling" arguments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

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.32

RPC Information Flow



11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.33

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).

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.34

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 \rightarrow 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

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.35

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 \ll same-machine RPC \ll network RPC
 - Means programmers must be aware that RPC is not free
 - » Caching can help, but may make failure handling complex

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.36

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)

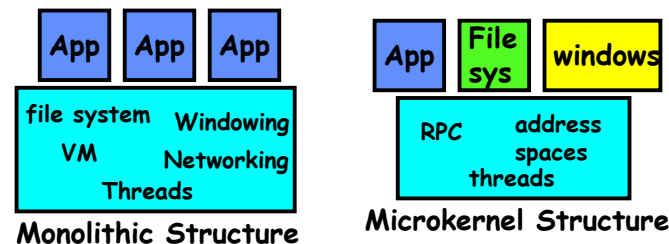
11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.37

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.

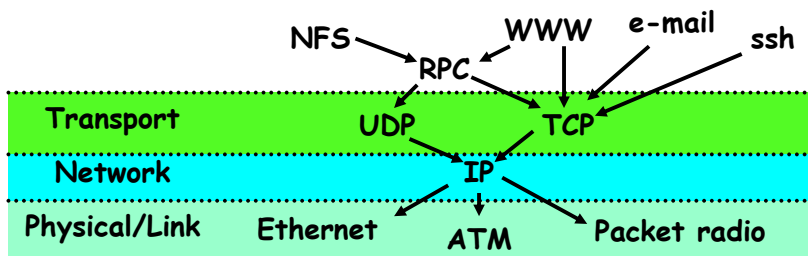
11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.38

Network Protocols

- Networking protocols: many levels
 - Physical level: mechanical and electrical network (e.g. how are 0 and 1 represented)
 - Link level: packet formats/error control (for instance, the CSMA/CD protocol)
 - Network level: network routing, addressing
 - Transport Level: reliable message delivery
- Protocols on today's Internet:



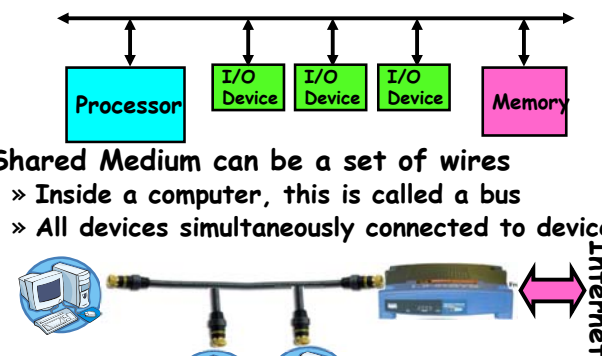
11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.39

Broadcast Networks

- **Broadcast Network: Shared Communication Medium**



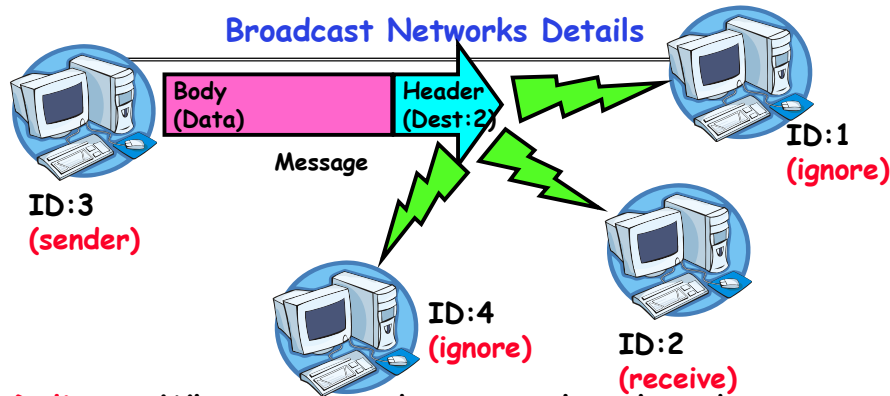
- Shared Medium can be a set of wires
 - » Inside a computer, this is called a bus
 - » All devices simultaneously connected to devices
- Originally, Ethernet was a broadcast network
 - » All computers on local subnet connected to one another
- More examples (wireless: medium is air): cellular phones, GSM GPRS, EDGE, CDMA 1xRTT, and 1EvDO

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.40

Broadcast Networks Details



- **Delivery:** When you broadcast a packet, how does a receiver know who it is for? (packet goes to everyone!)
 - Put header on front of packet: [Destination | Packet]
 - Everyone gets packet, discards if not the target
 - In Ethernet, this check is done in hardware
 - » No OS interrupt if not for particular destination
 - This is layering: we're going to build complex network protocols by layering on top of the packet

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.41

Carrier Sense, Multiple Access/Collision Detection

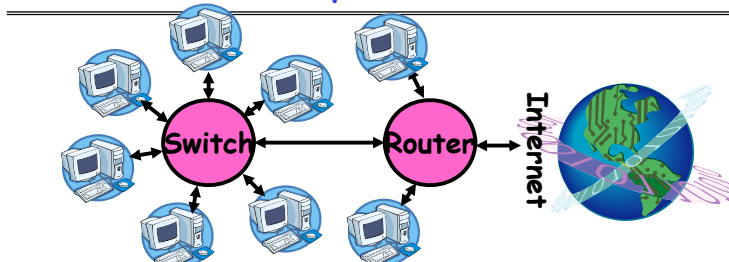
- Ethernet (early 80's): first practical local area network
 - It is the most common LAN for UNIX, PC, and Mac
 - Use wire instead of radio, but still broadcast medium
- Key advance was in arbitration called **CSMA/CD**: Carrier sense, multiple access/collision detection
 - **Carrier Sense:** don't send unless idle
 - » Don't mess up communications already in process
 - **Collision Detect:** sender checks if packet trampled.
 - » If so, abort, wait, and retry.
 - **Backoff Scheme:** Choose wait time before trying again
- How long to wait after trying to send and failing?
 - What if everyone waits the same length of time? Then, they all collide again at some time!
 - Must find way to break up shared behavior with nothing more than shared communication channel
- Adaptive randomized waiting strategy:
 - **Adaptive and Random:** First time, pick random wait time with some initial mean. If collide again, pick random value from bigger mean wait time. Etc.
 - Randomness is important to decouple colliding senders
 - Scheme figures out how many people are trying to send!

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.42

Point-to-point networks



- Why have a shared bus at all? Why not simplify and only have point-to-point links + routers/switches?
 - Originally wasn't cost-effective
 - Now, easy to make high-speed switches and routers that can forward packets from a sender to a receiver.
- **Point-to-point network:** a network in which every physical wire is connected to only two computers
- **Switch:** a bridge that transforms a shared-bus (broadcast) configuration into a point-to-point network.
- **Router:** a device that acts as a junction between two networks to transfer data packets among them.

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.43

The Internet Protocol: "IP"

- The Internet is a large network of computers spread across the globe
 - According to the Internet Systems Consortium, there were over 681 million computers as of July 2009
 - In principle, every host can speak with every other one under the right circumstances
- **IP Packet:** a network packet on the internet
- **IP Address:** a 32-bit integer used as the destination of an IP packet
 - Often written as four dot-separated integers, with each integer from 0–255 (thus representing $8 \times 4 = 32$ bits)
 - Example CS file server is: 169.229.60.83 \equiv 0xA9E53C53
- **Internet Host:** a computer connected to the Internet
 - Host has one or more IP addresses used for routing
 - » Some of these may be private and unavailable for routing
 - Not every computer has a unique IP address
 - » Groups of machines may share a single IP address
 - » In this case, machines have private addresses behind a "Network Address Translation" (NAT) gateway

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.44

Address Subnets

- **Subnet:** A network connecting a set of hosts with related destination addresses
- With IP, all the addresses in subnet are related by a prefix of bits
 - **Mask:** The number of matching prefix bits
 - » Expressed as a single value (e.g., 24) or a set of ones in a 32-bit value (e.g., 255.255.255.0)
- A subnet is identified by 32-bit value, with the bits which differ set to zero, followed by a slash and a mask
 - Example: 128.32.131.0/24 designates a subnet in which all the addresses look like 128.32.131.XX
 - Same subnet: 128.32.131.0/255.255.255.0
- Difference between subnet and complete network range
 - Subnet is always a subset of address range
 - Once, subnet meant single physical broadcast wire; now, less clear exactly what it means (virtualized by switches)

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.45

Address Ranges in IP

- IP address space divided into prefix-delimited ranges:
 - Class A: NN.0.0.0/8
 - » NN is 1-126 (126 of these networks)
 - » 16,777,214 IP addresses per network
 - » 10.xx.yy.zz is private
 - » 127.xx.yy.zz is loopback
 - Class B: NN.MM.0.0/16
 - » NN is 128-191, MM is 0-255 (16,384 of these networks)
 - » 65,534 IP addresses per network
 - » 172.[16-31].xx.yy are private
 - Class C: NN.MM.LL.0/24
 - » NN is 192-223, MM and LL 0-255 (2,097,151 of these networks)
 - » 254 IP addresses per networks
 - » 192.168.xx.yy are private
- Address ranges are often owned by organizations
 - Can be further divided into subnets

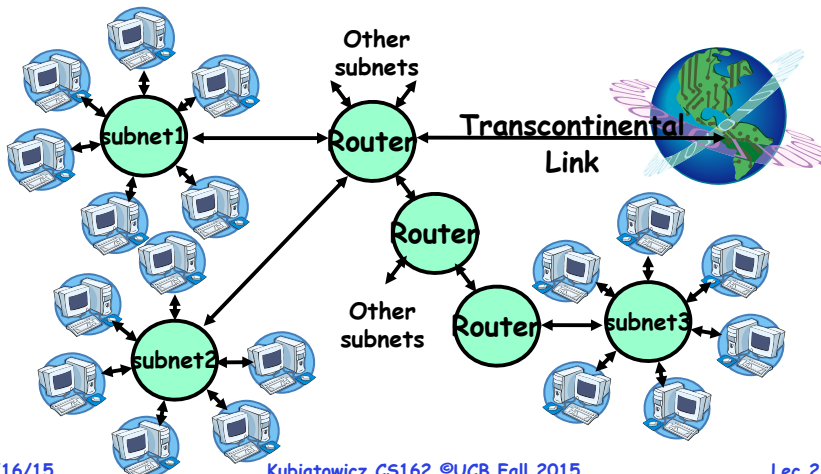
11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.46

Hierarchical Networking: The Internet

- How can we build a network with millions of hosts?
 - Hierarchy! Not every host connected to every other one
 - Use a network of Routers to connect subnets together
 - » Routing is often by prefix: e.g. first router matches first 8 bits of address, next router matches more, etc.



11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.47

Simple Network Terminology

- Local-Area Network (LAN) - designed to cover small geographical area
 - Multi-access bus, ring, or star network
 - Speed \approx 10 - 1000 Megabits/second (even 40-100GB/s)
 - Broadcast is fast and cheap
 - In small organization, a LAN could consist of a single subnet. In large organizations (like UC Berkeley), a LAN contains many subnets
- Wide-Area Network (WAN) - links geographically separated sites
 - Point-to-point connections over long-haul lines (often leased from a phone company)
 - Speed \approx 1.544 - 155 Megabits/second
 - Broadcast usually requires multiple messages

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.48

Routing

- **Routing: the process of forwarding packets hop-by-hop through routers to reach their destination**
 - Need more than just a destination address!
 - » Need a path
 - Post Office Analogy:
 - » Destination address on each letter is not sufficient to get it to the destination
 - » To get a letter from here to Florida, must route to local post office, sorted and sent on plane to somewhere in Florida, be routed to post office, sorted and sent with carrier who knows where street and house is...
- **Internet routing mechanism: routing tables**
 - Each router does table lookup to decide which link to use to get packet closer to destination
 - Don't need 4 billion entries in table: routing is by subnet
 - Could packets be sent in a loop? Yes, if tables incorrect
- **Routing table contains:**
 - Destination address range → output link closer to destination
 - Default entry (for subnets without explicit entries)



11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.49

Setting up Routing Tables

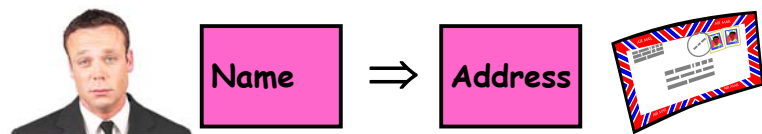
- **How do you set up routing tables?**
 - Internet has no centralized state!
 - » No single machine knows entire topology
 - » Topology constantly changing (faults, reconfiguration, etc)
 - Need dynamic algorithm that acquires routing tables
 - » Ideally, have one entry per subnet or portion of address
 - » Could have "default" routes that send packets for unknown subnets to a different router that has more information
- **Possible algorithm for acquiring routing table**
 - Routing table has "cost" for each entry
 - » Includes number of hops to destination, congestion, etc.
 - » Entries for unknown subnets have infinite cost
 - Neighbors periodically exchange routing tables
 - » If neighbor knows cheaper route to a subnet, replace your entry with neighbors entry (+1 for hop to neighbor)
- **In reality:**
 - Internet has networks of many different scales
 - Different algorithms run at different scales

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.50

Naming in the Internet



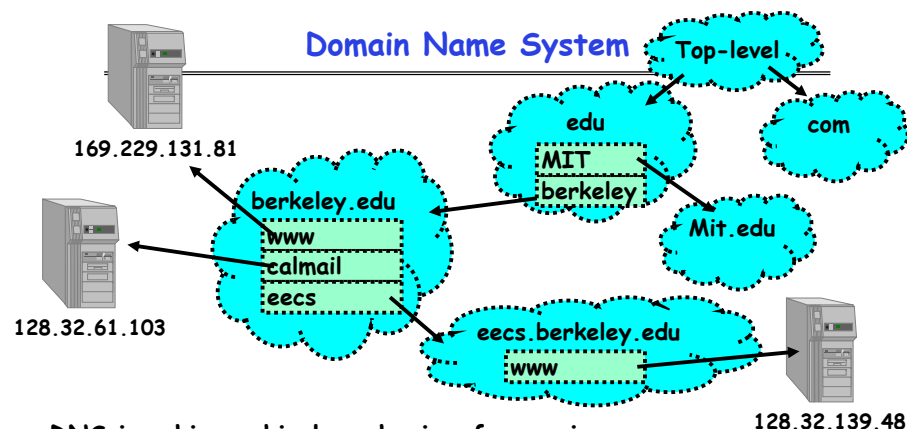
- **How to map human-readable names to IP addresses?**
 - E.g. `www.berkeley.edu` ⇒ `128.32.139.48`
 - E.g. `www.google.com` ⇒ different addresses depending on location, and load
- **Why is this necessary?**
 - IP addresses are hard to remember
 - IP addresses change:
 - » Say, Server 1 crashes gets replaced by Server 2
 - » Or - `google.com` handled by different servers
- **Mechanism: Domain Naming System (DNS)**

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.51

Domain Name System



- **DNS is a hierarchical mechanism for naming**
 - Name divided in domains, right to left: `www.eecs.berkeley.edu`
- **Each domain owned by a particular organization**
 - Top level handled by ICANN (Internet Corporation for Assigned Numbers and Names)
 - Subsequent levels owned by organizations
- **Resolution: series of queries to successive servers**
- **Caching: queries take time, so results cached for period of time**

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.52

How Important is Correct Resolution?

- **If attacker manages to give incorrect mapping:**
 - Can get someone to route to server, thinking that they are routing to a different server
 - » Get them to log into "bank" - give up username and password
- **Is DNS Secure?**
 - Definitely a weak link
 - » What if "response" returned from different server than original query?
 - » Get person to use incorrect IP address!
 - Attempt to avoid substitution attacks:
 - » Query includes random number which must be returned
- **In July 2008, hole in DNS security located!**
 - Dan Kaminsky (security researcher) discovered an attack that broke DNS globally
 - » One person in an ISP convinced to load particular web page, then *all* users of that ISP end up pointing at wrong address
 - High profile, highly advertised need for patching DNS
 - » Big press release, lots of mystery
 - » Security researchers told no speculation until patches applied

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.53

Network Layering

- **Layering:** building complex services from simpler ones
 - Each layer provides services needed by higher layers by utilizing services provided by lower layers
- The physical/link layer is pretty limited
 - Packets are of limited size (called the "Maximum Transfer Unit or MTU: often 200-1500 bytes in size)
 - Routing is limited to within a physical link (wire) or perhaps through a switch
- Our goal in the following is to show how to construct a secure, ordered, message service routed to anywhere:

Physical Reality: Packets	Abstraction: Messages
Limited Size	Arbitrary Size
Unordered (sometimes)	Ordered
Unreliable	Reliable
Machine-to-machine	Process-to-process
Only on local area net	Routed anywhere
Asynchronous	Synchronous
Insecure	Secure

11/16/15

Lec 21.54

Summary (1/2)

- **Protocol: Agreement between two parties as to how information is to be transmitted**
- **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 $n \geq 3f+1$

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.55

Summary (2/2)

- **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)
- **Internet Protocol (IP)**
 - Used to route messages through routes across globe
 - 32-bit addresses, 16-bit ports
- **DNS: System for mapping from names \Rightarrow IP addresses**
 - Hierarchical mapping from authoritative domains
 - Recent flaws discovered
- **Datagram: a self-contained message whose arrival, arrival time, and content are not guaranteed**
- **Ordered messages:**
 - Use sequence numbers and reorder at destination
- **Reliable messages:**
 - Use Acknowledgements

11/16/15

Kubiatowicz CS162 ©UCB Fall 2015

Lec 21.56