

# CS162

## Operating Systems and Systems Programming

### Lecture 23

## Distributed Storage, Key-Value Stores, Security

April 27<sup>th</sup>, 2015

Prof. John Kubiatowicz

<http://cs162.eecs.Berkeley.edu>

## Recall: Two Phase (2PC) Commit

- Distributed transaction: Two or more machines agree to do something, or not do it, **atomically**
- Two Phase Commit:
  - 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

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.2

## Brief aside: 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:
 

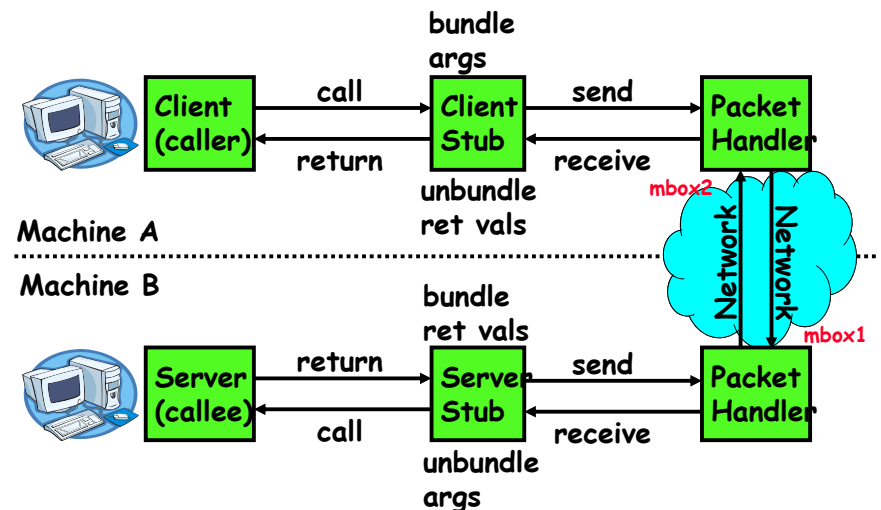
```
fileSys→Read("rutabaga");
```
- Implementation:
  - Request-response message passing (under covers!)
  - "Stub" provides glue on client/server
    - » Client stub is responsible for "marshalling" arguments and "unmarshalling" the return values
    - » Server-side stub is responsible for "unmarshalling" 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.

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.3

## RPC Information Flow



4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.4

## 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).

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.5

## 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

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.6

## 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

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.7

## 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)

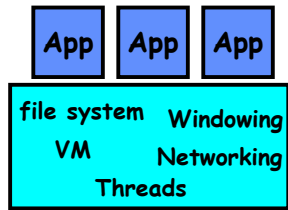
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

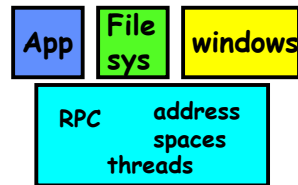
Lec 23.8

## Microkernel operating systems

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



Monolithic Structure



Microkernel Structure

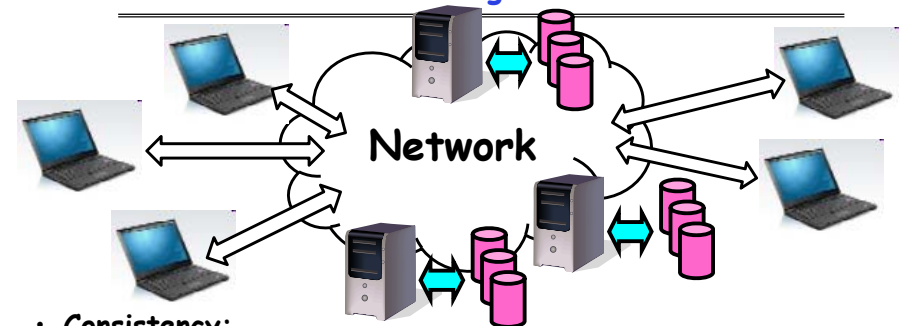
- Why split the OS into separate domains?
  - Fault isolation: bugs are more isolated (build a firewall)
  - Enforces modularity: allows incremental upgrades of pieces of software (client or server)
  - Location transparent: service can be local or remote
    - » For example in the X windowing system: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

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.9

## Network-Attached Storage and the CAP Theorem



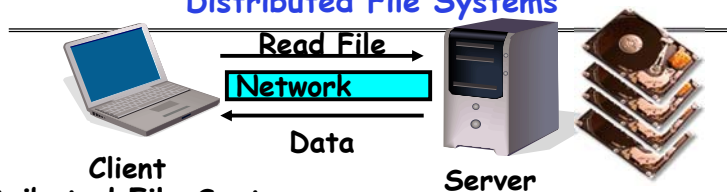
- Consistency:
  - Changes appear to everyone in the same serial order
- Availability:
  - Can get a result at any time
- Partition-Tolerance
  - System continues to work even when network becomes partitioned
- Consistency, Availability, Partition-Tolerance (CAP) Theorem: **Cannot have all three at same time**
  - Otherwise known as "Brewer's Theorem"

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.10

## Distributed File Systems



- Distributed File System:
  - Transparent access to files stored on a remote disk

- Naming choices (always an issue):

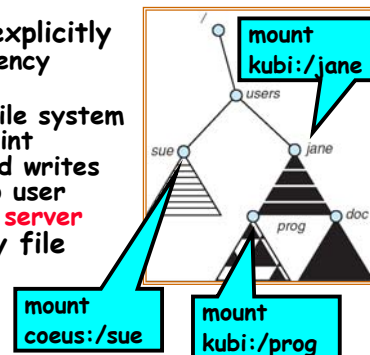
- *Hostname:localname*: Name files explicitly
  - » No location or migration transparency

- *Mounting of remote file systems*

- » System manager mounts remote file system by giving name and local mount point
- » Transparent to user: all reads and writes look like local reads and writes to user e.g. `/users/sue/foo` → `/sue/foo` on server

- *A single, global name space*: every file in the world has unique name

- » Location Transparency: servers can change and files can move without involving user

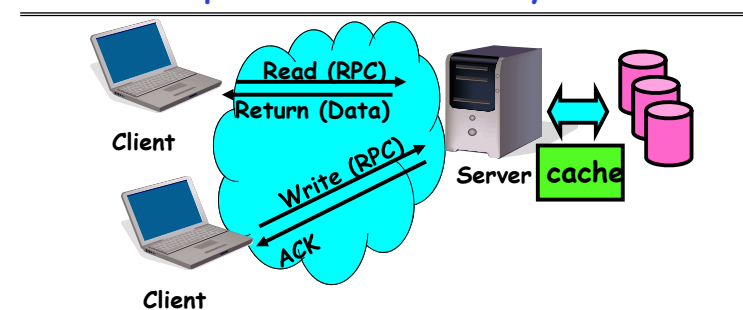


4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.11

## Simple Distributed File System



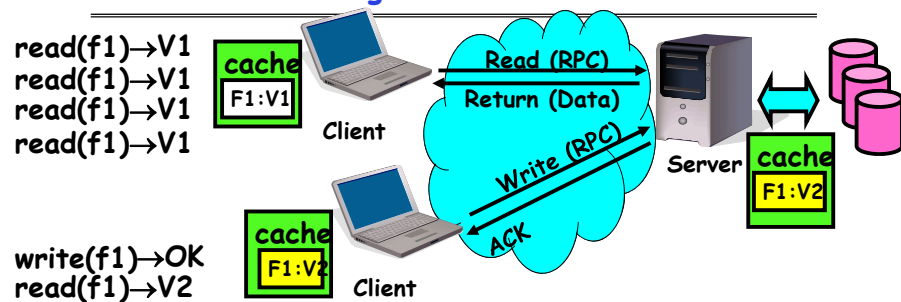
- Remote Disk: Reads and writes forwarded to server
  - Use Remote Procedure Calls (RPC) to translate file system calls into remote requests
  - No local caching/can be caching at server-side
- Advantage: Server provides completely consistent view of file system to multiple clients
- Problems? Performance!
  - Going over network is slower than going to local memory
  - Lots of network traffic/not well pipelined
  - Server can be a bottleneck

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.12

## Use of caching to reduce network load



- **Idea:** Use caching to reduce network load
  - In practice: use buffer cache at source and destination
- **Advantage:** if open/read/write/close can be done locally, don't need to do any network traffic...fast!
- **Problems:**
  - Failure:
    - » Client caches have data not committed at server
  - Cache consistency!
    - » Client caches not consistent with server/each other

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.13

## Failures



- What if server crashes? Can client wait until server comes back up and continue as before?
  - Any data in server memory but not on disk can be lost
  - Shared state across RPC: What if server crashes after seek? Then, when client does "read", it will fail
  - Message retries: suppose server crashes after it does UNIX "rm foo", but before acknowledgment?
    - » Message system will retry: send it again
    - » How does it know not to delete it again? (could solve with two-phase commit protocol, but NFS takes a more ad hoc approach)
- **Stateless protocol:** A protocol in which all information required to process a request is passed with request
  - Server keeps no state about client, except as hints to help improve performance (e.g. a cache)
  - Thus, if server crashes and restarted, requests can continue where left off (in many cases)
- What if client crashes?
  - Might lose modified data in client cache

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.14

## Administrivia

- **Midterm 2 grading**
  - In progress. Hopefully done by end of week.
  - Solutions have been posted
- **Final Exam**
  - Friday, May 15<sup>th</sup>, 2015.
  - 3-6P, Wheeler Auditorium
  - All material from the course
    - » With slightly more focus on second half, but you are still responsible for all the material
  - Two sheets of notes, both sides
  - Will need dumb calculator
- **Should be working on Project 3!**
  - Checkpoint 1 this Wednesday

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.15

## Administrivia (con't)

- **Final Lecture topics submitted to me:**
  - Real Time Operating systems
  - Peer to peer systems and/or Distributed Systems
  - OS trends in the mobile phone industry (Android, etc)
    - » Differences from traditional OSes?
  - GPU and ManyCore programming (and/or OSes?)
  - Virtual Machines and/or Trusted Hardware for security
  - Systems programming for non-standard computer systems
    - » i.e. Quantum Computers, Biological Computers, ...
  - Net Neutrality and/or making the Internet Faster
  - Mesh networks
  - Device drivers
  - A couple of votes for Dragons...
- This is a lot of topics...

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.16



## Network File System (NFS)

- Three Layers for NFS system
  - **UNIX file-system interface**: open, read, write, close calls + file descriptors
  - **VFS layer**: distinguishes local from remote files
    - » Calls the NFS protocol procedures for remote requests
  - **NFS service layer**: bottom layer of the architecture
    - » Implements the NFS protocol
- NFS Protocol: RPC for file operations on server
  - Reading/searching a directory
  - manipulating links and directories
  - accessing file attributes/reading and writing files
- **Write-through caching**: Modified data committed to server's disk before results are returned to the client
  - lose some of the advantages of caching
  - time to perform write() can be long
  - Need some mechanism for readers to eventually notice changes! (more on this later)

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.17

## NFS Continued

- NFS servers are **stateless**; each request provides all arguments require for execution
  - E.g. reads include information for entire operation, such as ReadAt(inumber, position), not Read(openfile)
  - No need to perform network open() or close() on file - each operation stands on its own
- **Idempotent**: Performing requests multiple times has same effect as performing it exactly once
  - Example: Server crashes between disk I/O and message send, client resend read, server does operation again
  - Example: Read and write file blocks: just re-read or re-write file block - no side effects
  - Example: What about "remove"? NFS does operation twice and second time returns an advisory error
- Failure Model: Transparent to client system
  - Is this a good idea? What if you are in the middle of reading a file and server crashes?
  - Options (NFS Provides both):
    - » Hang until server comes back up (next week?)
    - » Return an error. (Of course, most applications don't know they are talking over network)

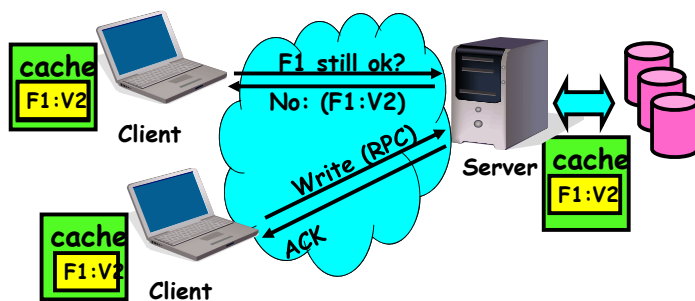
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.18

## NFS Cache consistency

- NFS protocol: weak consistency
  - Client polls server periodically to check for changes
    - » Polls server if data hasn't been checked in last 3-30 seconds (exact timeout is tunable parameter).
    - » Thus, when file is changed on one client, server is notified, but other clients use old version of file until timeout.



- What if multiple clients write to same file?
  - » In NFS, can get either version (or parts of both)
  - » Completely arbitrary!

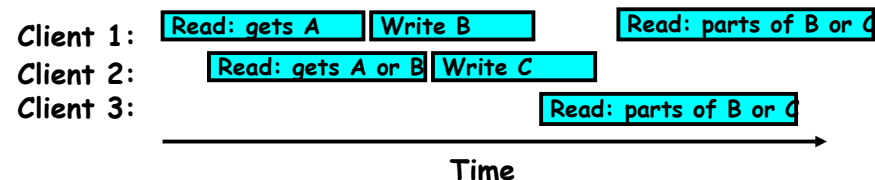
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.19

## Sequential Ordering Constraints

- What sort of cache coherence might we expect?
  - i.e. what if one CPU changes file, and before it's done, another CPU reads file?
- Example: Start with file contents = "A"



- What would we actually want?
  - Assume we want distributed system to behave exactly the same as if all processes are running on single system
    - » If read finishes before write starts, get old copy
    - » If read starts after write finishes, get new copy
    - » Otherwise, get either new or old copy
  - For NFS:
    - » If read starts more than 30 seconds after write, get new copy; otherwise, could get partial update

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.20

## Andrew File System

- Andrew File System (AFS, late 80's) → DCE DFS (commercial product)
- **Callbacks:** Server records who has copy of file
  - On changes, server immediately tells all with old copy
  - No polling bandwidth (continuous checking) needed
- Write through on close
  - Changes not propagated to server until close()
  - Session semantics: updates visible to other clients only after the file is closed
    - » As a result, do not get partial writes: all or nothing!
    - » Although, for processes on local machine, updates visible immediately to other programs who have file open
- In AFS, everyone who has file open sees old version
  - Don't get newer versions until reopen file

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.21

## Andrew File System (con't)

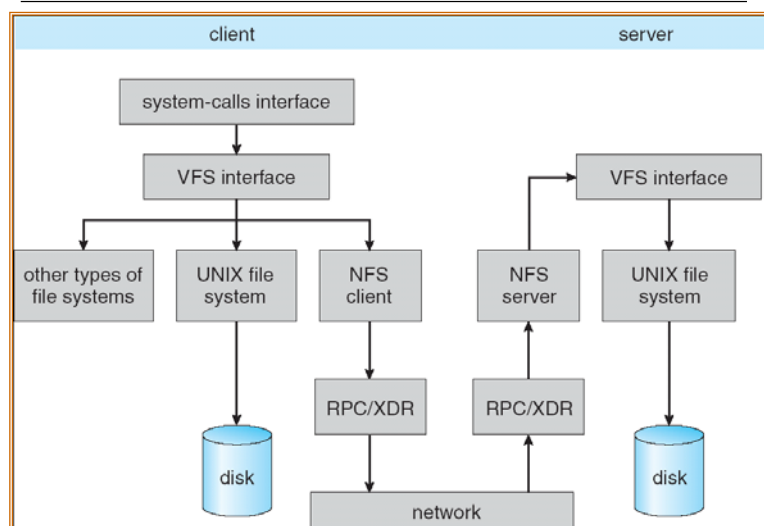
- Data cached on local disk of client as well as memory
  - On open with a cache miss (file not on local disk):
    - » Get file from server, set up callback with server
  - On write followed by close:
    - » Send copy to server; tells all clients with copies to fetch new version from server on next open (using callbacks)
- What if server crashes? Lose all callback state!
  - Reconstruct callback information from client: go ask everyone "who has which files cached?"
- AFS Pro: Relative to NFS, less server load:
  - Disk as cache ⇒ more files can be cached locally
  - Callbacks ⇒ server not involved if file is read-only
- For both AFS and NFS: central server is bottleneck!
  - Performance: all writes→server, cache misses→server
  - Availability: Server is single point of failure
  - Cost: server machine's high cost relative to workstation

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.22

## Implementation of NFS

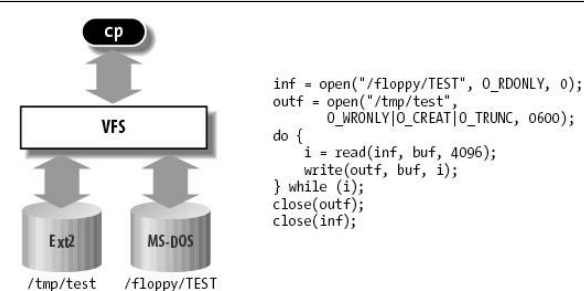


4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.23

## Enabling Factor: Virtual Filesystem (VFS)



```
inf = open("/floppy/TEST", O_RDONLY, 0);
outf = open("/tmp/test",
           O_WRONLY|O_CREAT|O_TRUNC, 0600);
do {
    i = read(inf, buf, 4096);
    write(outf, buf, i);
} while (i);
close(outf);
close(inf);
```

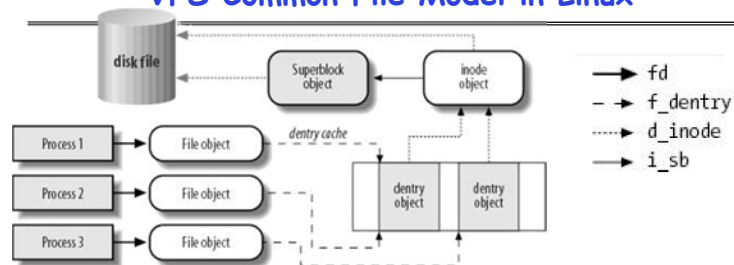
- **VFS:** Virtual abstraction similar to local file system
  - Provides virtual superblocks, inodes, files, etc
  - Compatible with a variety of local and remote file systems
    - » provides object-oriented way of implementing file systems
- VFS allows the same system call interface (the API) to be used for different types of file systems
  - The API is to the VFS interface, rather than any specific type of file system
- In linux, "VFS" stands for "Virtual Filesystem Switch"

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.24

## VFS Common File Model in Linux



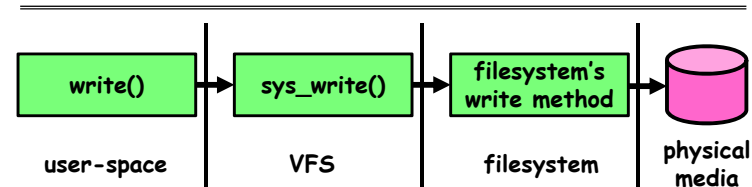
- Four primary object types for VFS:
  - superblock object: represents a specific mounted filesystem
  - inode object: represents a specific file
  - dentry object: represents a directory entry
  - file object: represents open file associated with process
- There is no specific directory object (VFS treats directories as files)
- May need to fit the model by faking it
  - Example: make it look like directories are files
  - Example: make it look like have inodes, superblocks, etc.

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.25

## Linux VFS



- An operations object is contained within each primary object type to set operations of specific filesystems
  - "super\_operations": methods that kernel can invoke on a specific filesystem, i.e. `write_inode()` and `sync_fs()`.
  - "inode\_operations": methods that kernel can invoke on a specific file, such as `create()` and `link()`
  - "dentry\_operations": methods that kernel can invoke on a specific directory entry, such as `d_compare()` or `d_delete()`
  - "file\_operations": methods that process can invoke on an open file, such as `read()` and `write()`
- There are a lot of operations

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.26

## Key Value Storage

- Handle huge volumes of data, e.g., PBs
  - Store (key, value) tuples
- Simple interface
  - `put(key, value);` // insert/write "value" associated with "key"
  - `value = get(key);` // get/read data associated with "key"
- Used sometimes as a simpler but more scalable "database"

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.27

## Key Values: Examples

- Amazon:
  - Key: customerID
  - Value: customer profile (e.g., buying history, credit card, ..)
- Facebook, Twitter:
  - Key: UserID
  - Value: user profile (e.g., posting history, photos, friends, ...)
- iCloud/iTunes:
  - Key: Movie/song name
  - Value: Movie, Song

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.28

## Examples

- **Amazon**
  - DynamoDB: internal key value store used to power Amazon.com (shopping cart)
  - Simple Storage System (S3)
- **BigTable/HBase/Hypertable**: distributed, scalable data storage
- **Cassandra**: “distributed data management system” (developed by Facebook)
- **Memcached**: in-memory key-value store for small chunks of arbitrary data (strings, objects)
- **eDonkey/eMule**: peer-to-peer sharing system
- ...

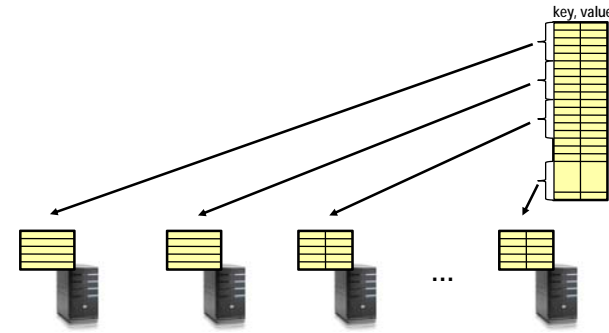
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.29

## Key Value Store

- Also called Distributed Hash Tables (DHT)
- Main idea: partition set of key-values across many machines



4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.30

## Challenges



- **Fault Tolerance**: handle machine failures without losing data and without degradation in performance
- **Scalability**:
  - Need to scale to thousands of machines
  - Need to allow easy addition of new machines
- **Consistency**: maintain data consistency in face of node failures and message losses
- **Heterogeneity** (if deployed as peer-to-peer systems):
  - Latency: 1ms to 1000ms
  - Bandwidth: 32Kb/s to 100Mb/s

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.31

## Key Questions

- **put(key, value)**: where do you store a new (key, value) tuple?
- **get(key)**: where is the value associated with a given “key” stored?
- And, do the above while providing
  - Fault Tolerance
  - Scalability
  - Consistency

4/27/15

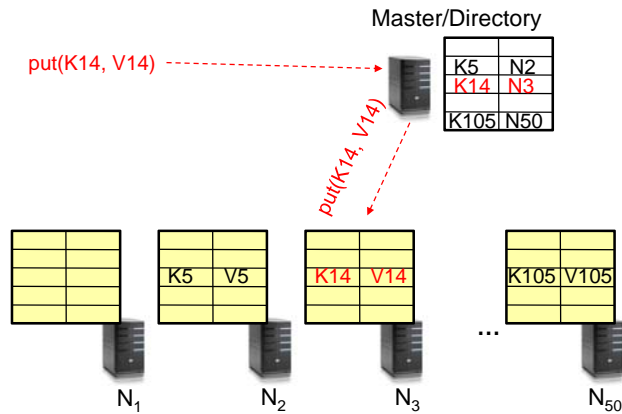
Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.32



## Directory-Based Architecture

- Have a node maintain the mapping between keys and the machines (nodes) that store the values associated with the keys



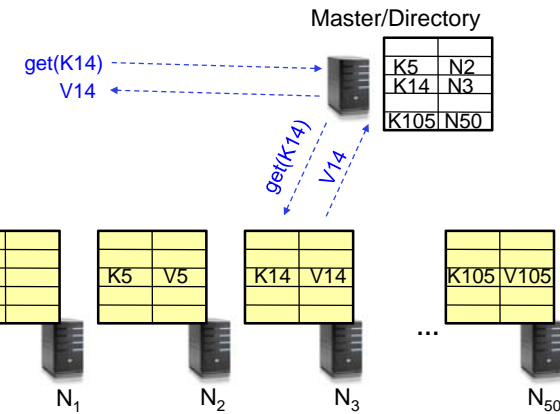
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.33

## Directory-Based Architecture

- Have a node maintain the mapping between keys and the machines (nodes) that store the values associated with the keys



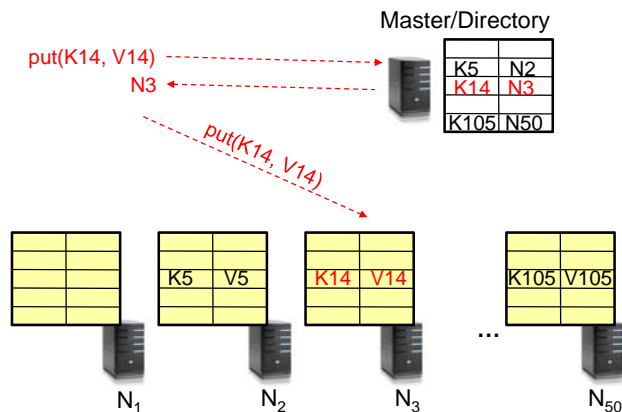
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.34

## Directory-Based Architecture

- Having the master relay the requests → recursive query
- Another method: iterative query (this slide)
  - Return node to requester and let requester contact node



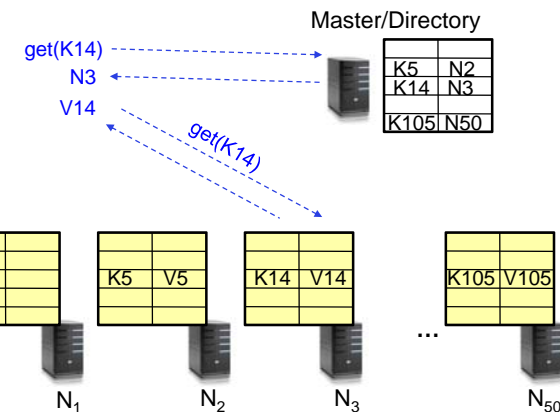
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.35

## Directory-Based Architecture

- Having the master relay the requests → recursive query
- Another method: iterative query
  - Return node to requester and let requester contact node

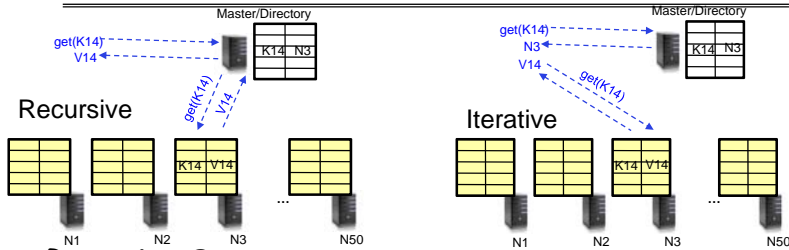


4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.36

## Discussion: Iterative vs. Recursive Query



### Recursive Query:

#### - Advantages:

- » Faster, as typically master/directory closer to nodes
- » Easier to maintain consistency, as master/directory can serialize puts()/gets()

- Disadvantages: scalability bottleneck, as all "Values" go through master/directory

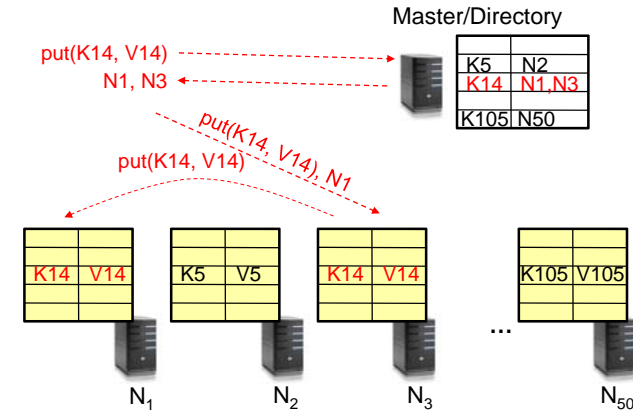
### Iterative Query

- Advantages: more scalable

- Disadvantages: slower, harder to enforce data consistency

## Fault Tolerance

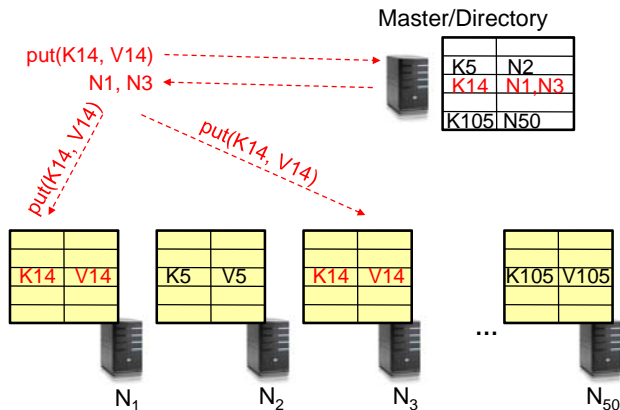
- Replicate value on several nodes
- Usually, place replicas on different racks in a datacenter to guard against rack failures



## Fault Tolerance

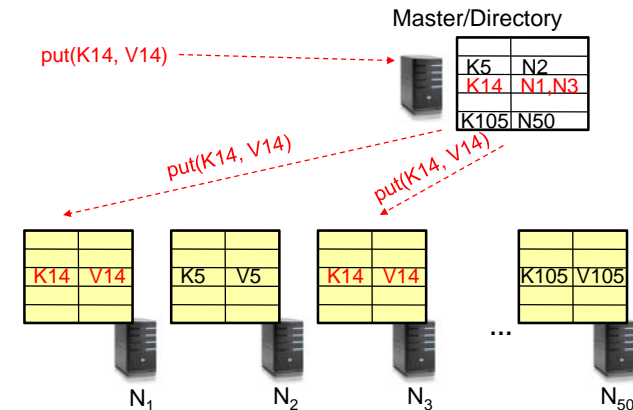
• Again, we can have

- Recursive replication (previous slide)
- Iterative replication (this slide)



## Fault Tolerance

• Or we can use recursive query and iterative replication...



## Scalability

- **Storage: use more nodes**
- **Number of requests:**
  - Can serve requests from all nodes on which a value is stored in parallel
  - Master can replicate a popular value on more nodes
- **Master/directory scalability:**
  - Replicate it
  - Partition it, so different keys are served by different masters/directories
    - » How do you partition?

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.41

## Scalability: Load Balancing

- **Directory keeps track of the storage availability at each node**
  - Preferentially insert new values on nodes with more storage available
- **What happens when a new node is added?**
  - Cannot insert only new values on new node. Why?
  - Move values from the heavy loaded nodes to the new node
- **What happens when a node fails?**
  - Need to replicate values from fail node to other nodes

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.42

## Consistency

- **Need to make sure that a value is replicated correctly**
- **How do you know a value has been replicated on every node?**
  - Wait for acknowledgements from every node
- **What happens if a node fails during replication?**
  - Pick another node and try again
- **What happens if a node is slow?**
  - Slow down the entire put()? Pick another node?
- **In general, with multiple replicas**
  - Slow puts and fast gets

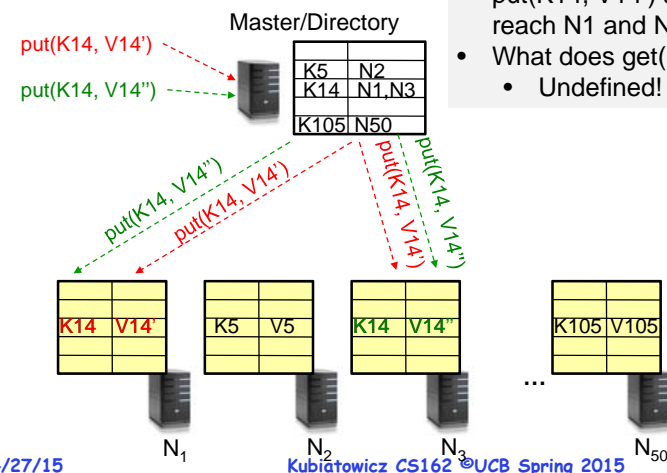
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.43

## Consistency (cont'd)

- **If concurrent updates (i.e., puts to same key) may need to make sure that updates happen in the same order**



4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.44

## Consistency (cont'd)

- Large variety of consistency models:
  - Atomic consistency (linearizability): reads/writes (gets/puts) to replicas appear as if there was a single underlying replica (single system image)
    - » Think "one updated at a time"
    - » Transactions
  - Eventual consistency: given enough time all updates will propagate through the system
    - » One of the weakest form of consistency; used by many systems in practice
  - And many others: causal consistency, sequential consistency, strong consistency, ...

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.45

## Quorum Consensus

- Improve put() and get() operation performance
- Define a replica set of size N
  - put() waits for acknowledgements from at least W replicas
  - get() waits for responses from at least R replicas
  - $W+R > N$
- Why does it work?
  - There is at least one node that contains the update
- Why might you use  $W+R > N+1$ ?

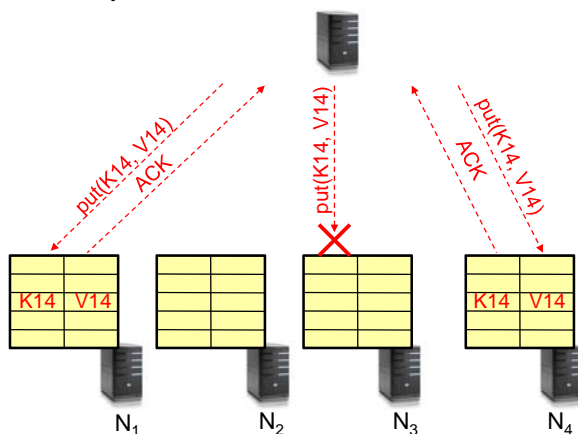
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.46

## Quorum Consensus Example

- $N=3$ ,  $W=2$ ,  $R=2$
- Replica set for K14: {N1, N2, N4}
- Assume put() on N3 fails



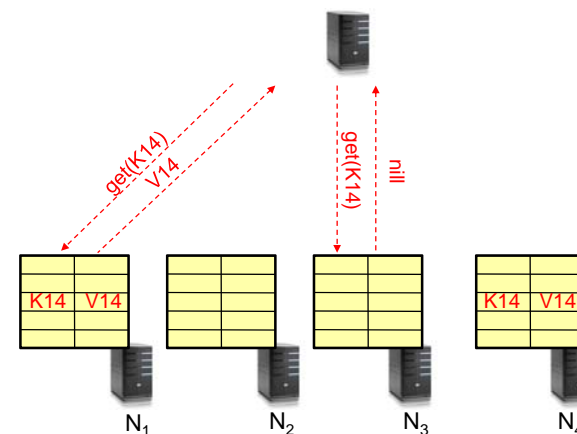
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.47

## Quorum Consensus Example

- Now, issuing get() to any two nodes out of three will return the answer



4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.48

## Scaling Up Directory

- **Challenge:**
  - Directory contains a number of entries equal to number of (key, value) tuples in the system
  - Can be tens or hundreds of billions of entries in the system!
- **Solution: consistent hashing**
- **Associate to each node a unique *id* in an uni-dimensional space  $0..2^m-1$** 
  - Partition this space across  $m$  machines
  - Assume keys are in same uni-dimensional space
  - Each (Key, Value) is stored at the node with the smallest ID larger than Key

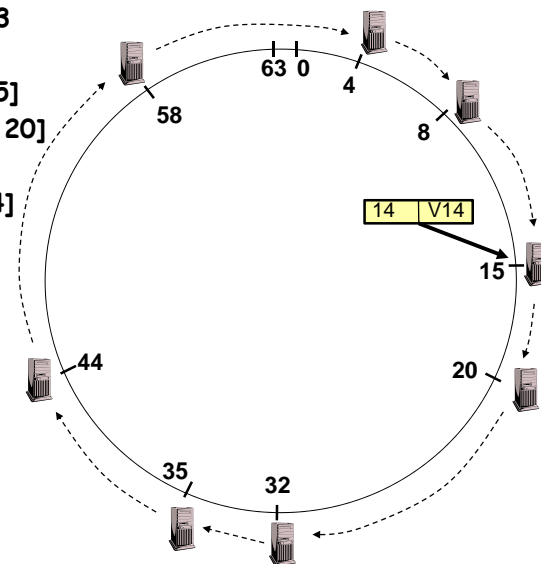
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.49

## Key to Node Mapping Example

- $m = 6 \rightarrow$  ID space:  $0..63$
- Node 8 maps keys  $[5, 8]$
- Node 15 maps keys  $[9, 15]$
- Node 20 maps keys  $[16, 20]$
- ...
- Node 4 maps keys  $[59, 4]$



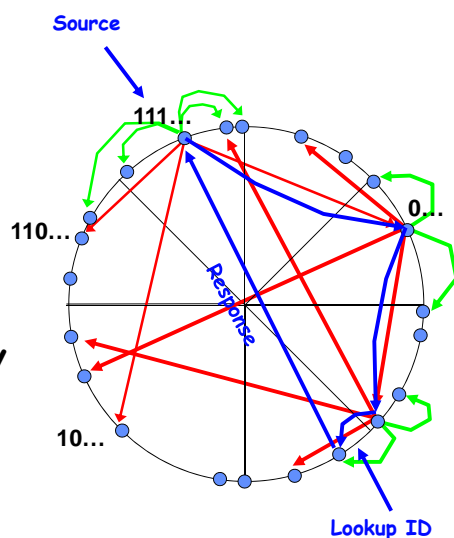
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.50

## Lookup in Chord-like system (with Leaf Set)

- **Assign IDs to nodes**
  - Map hash values to node with closest ID
- **Leaf set is successors and predecessors**
  - All that's needed for correctness
- **Routing table matches successively longer prefixes**
  - Allows efficient lookups
- **Data Replication:**
  - On leaf set



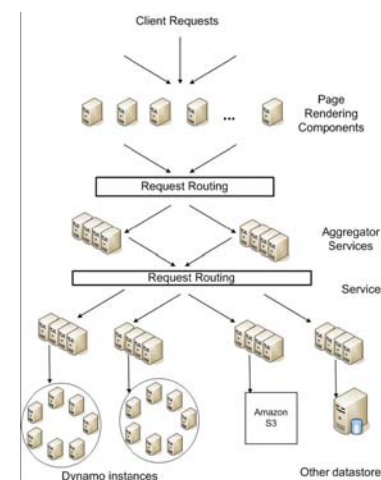
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.51

## DynamoDB Example: Service Level Agreements (SLA)

- Application can deliver its functionality in a bounded time:
  - Every dependency in the platform needs to deliver its functionality with even tighter bounds.
- Example: service guaranteeing that it will provide a response within 300ms for 99.9% of its requests for a peak client load of 500 requests per second
- Contrast to services which focus on mean response time



Service-oriented architecture of Amazon's platform

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.52



## What is Computer Security Today?

- Computing in the presence of an adversary!
  - Adversary is the security field's defining characteristic
- Reliability, robustness, and fault tolerance
  - Dealing with Mother Nature (random failures)
- Security
  - Dealing with actions of a knowledgeable attacker dedicated to causing harm
  - Surviving malice, and not just mischance
- Wherever there is an adversary, there is a computer security problem!

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.53

## Protection vs. Security

- **Protection**: mechanisms for controlling access of programs, processes, or users to resources
  - Page table mechanism
  - Round-robin schedule
  - Data encryption
- **Security**: use of protection mech. to prevent misuse of resources
  - Misuse defined with respect to policy
    - » E.g.: prevent exposure of certain sensitive information
    - » E.g.: prevent unauthorized modification/deletion of data
  - Need to consider external environment the system operates in
    - » Most well-constructed system cannot protect information if user accidentally reveals password – social engineering challenge

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.54

## Security Requirements

- Authentication
  - Ensures that a user is who is claiming to be
- Data integrity
  - Ensure that data is not changed from source to destination or after being written on a storage device
- Confidentiality
  - Ensures that data is read only by authorized users
- Non-repudiation
  - Sender/client can't later claim didn't send/write data
  - Receiver/server can't claim didn't receive/write data

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.55

## Authentication: Identifying Users

- How to identify users to the system?
  - Passwords
    - » Shared secret between two parties
    - » Since only user knows password, someone types correct password ⇒ must be user typing it
    - » Very common technique
  - Smart Cards
    - » Electronics embedded in card capable of providing long passwords or satisfying challenge → response queries
    - » May have display to allow reading of password
    - » Or can be plugged in directly; several credit cards now in this category
  - Biometrics
    - » Use of one or more intrinsic physical or behavioral traits to identify someone
    - » Examples: fingerprint reader, palm reader, retinal scan
    - » Becoming quite a bit more common



4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.56

## Passwords: Secrecy



- System must keep copy of secret to check against passwords
  - What if malicious user gains access to list of passwords?
    - » Need to obscure information somehow
  - Mechanism: utilize a transformation that is difficult to reverse without the right key (e.g. encryption)
- Example: UNIX /etc/passwd file
  - passwd → one way transform(hash) → encrypted passwd
  - System stores only encrypted version, so OK even if someone reads the file!
  - When you type in your password, system compares encrypted version
- Problem: Can you trust encryption algorithm?
  - Example: one algorithm thought safe had back door
    - » Governments want back door so they can snoop
  - Also, security through obscurity doesn't work
    - » GSM encryption algorithm was secret; accidentally released; Berkeley grad students cracked in a few hours

4/27/15

Kubiatowicz CS162 @UCB Spring 2015

Lec 23.57

## Passwords: How easy to guess?

- Ways of Compromising Passwords
  - Password Guessing:
    - » Often people use obvious information like birthday, favorite color, girlfriend's name, etc...
  - Dictionary Attack:
    - » Work way through dictionary and compare encrypted version of dictionary words with entries in /etc/passwd
  - Dumpster Diving:
    - » Find pieces of paper with passwords written on them
    - » (Also used to get social-security numbers, etc)
- Paradox:
  - Short passwords are easy to crack
  - Long ones, people write down!
- Technology means we have to use longer passwords
  - UNIX initially required lowercase, 5-letter passwords: total of  $26^5 = 10$  million passwords
    - » In 1975, 10ms to check a password → 1 day to crack
    - » In 2005, .01μs to check a password → 0.1 seconds to crack
    - » Even faster today (use multiple processors)
  - Takes less time to check for all words in the dictionary!

4/27/15

Kubiatowicz CS162 @UCB Spring 2015

Lec 23.58

<p>UNCOMMON (NON-GIBBERISH) BASE WORD</p> <p>ORDER UNKNOWN</p> <p>Tr0ub4dor &amp; 3</p> <p>CAPS? COMMON SUBSTITUTIONS NUMERALS PUNCTUATION</p> <p>(YOU CAN ADD A FEW MORE BITS TO ACCOUNT FOR THE FACT THAT THIS IS ONLY ONE OF A FEW COMMON FORMATS.)</p>	<p>~28 BITS OF ENTROPY</p> <p><math>2^{28} = 3</math> DAYS AT 1000 GUESSES/SEC</p> <p>(PLAUSIBLE ATTACK ON A WEBA REMOTE WEB SERVICE: YES, CRACKING A STOLEN HASH IS FASTER, BUT IT'S NOT WHAT THE AVERAGE USER SHOULD WORRY ABOUT.)</p> <p>DIFFICULTY TO GUESS: EASY</p>	<p>WAS IT TROUBADOR? NO, TROUBADOR. AND ONE OF THE 0s WAS A ZERO?</p> <p>AND THERE WAS SOME SYMBOL...</p> <p>DIFFICULTY TO REMEMBER: HARD</p>
<p>correct horse battery staple</p> <p>FOUR RANDOM COMMON WORDS</p>	<p>~44 BITS OF ENTROPY</p> <p><math>2^{44} = 550</math> YEARS AT 1000 GUESSES/SEC</p> <p>DIFFICULTY TO GUESS: HARD</p>	<p>THAT'S A BATTERY STAPLE.</p> <p>CORRECT!</p> <p>DIFFICULTY TO REMEMBER: YOU'VE ALREADY MEMORIZED IT</p>

THROUGH 20 YEARS OF EFFORT, WE'VE SUCCESSFULLY TRAINED EVERYONE TO USE PASSWORDS THAT ARE HARD FOR HUMANS TO REMEMBER, BUT EASY FOR COMPUTERS TO GUESS.

<https://xkcd.com/936/>

4/27/15

Kubiatowicz CS162 @UCB Spring 2015

Lec 23.59

## Securing Communication: Cryptography

- Cryptography: communication in the presence of adversaries
- Studied for thousands of years
  - See the Simon Singh's **The Code Book** for an excellent, highly readable history
- Central goal: confidentiality
  - How to encode information so that an adversary can't extract it, but a friend can
- General premise: there is a key, possession of which allows decoding, but without which decoding is infeasible
  - Thus, key must be kept secret and not guessable

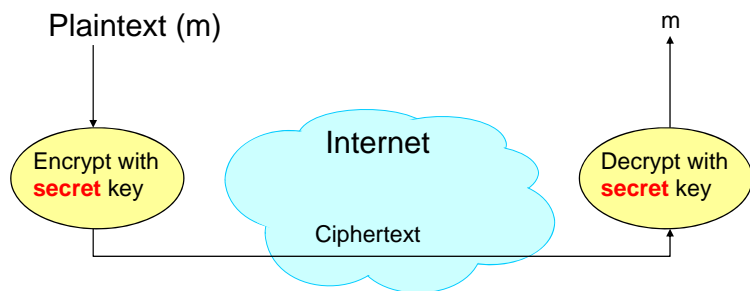
4/27/15

Kubiatowicz CS162 @UCB Spring 2015

Lec 23.60

## Using Symmetric Keys

- Same key for encryption and decryption
- Achieves confidentiality
- Vulnerable to tampering and replay attacks



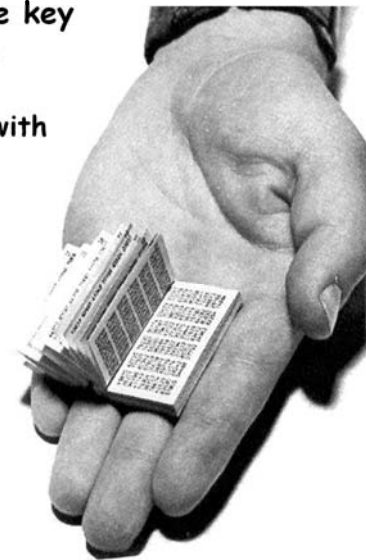
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.61

## Symmetric Keys

- Can just XOR plaintext with the key
  - Easy to implement, but easy to break using frequency analysis
  - Unbreakable alternative: XOR with one-time pad
    - » Use a different key for each message



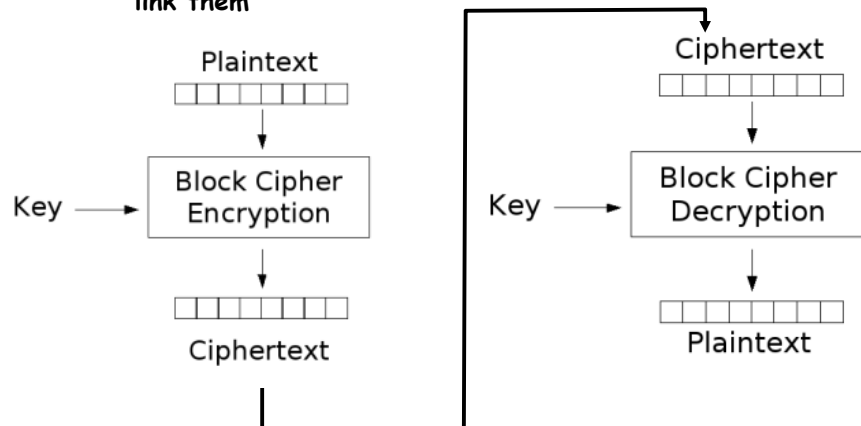
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.62

## Symmetric Keys

- More sophisticated (e.g., block cipher) algorithms
  - Works with a block size (e.g., 64 bits)
    - » To encrypt a stream, can encrypt blocks separately, or link them



4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.63

## Symmetric Key Ciphers - DES & AES

- Data Encryption Standard (DES)
  - Developed by IBM in 1970s, standardized by NBS/NIST
  - 56-bit key (decreased from 64 bits at NSA's request)
  - Still fairly strong other than brute-forcing the key space
    - » But custom hardware can crack a key in < 24 hours
  - Today many financial institutions use Triple DES
    - » DES applied 3 times, with 3 keys totaling 168 bits
- Advanced Encryption Standard (AES)
  - Replacement for DES standardized in 2002
  - Key size: 128, 192 or 256 bits
- How fundamentally strong are they?
  - No one knows (no proofs exist)

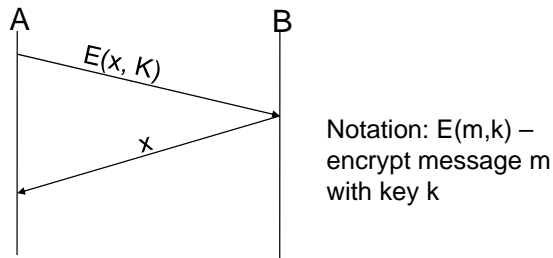
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.64

## Authentication via Secret Key

- **Main idea: entity proves identity by decrypting a secret encrypted with its own key**
  - $K$  - secret key shared only by A and B
- **A can ask B to authenticate itself by decrypting a nonce, i.e., random value,  $x$** 
  - Avoid replay attacks (attacker impersonating client or server)
- **Vulnerable to man-in-the middle attack**



4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.65

## Integrity: Cryptographic Hashes

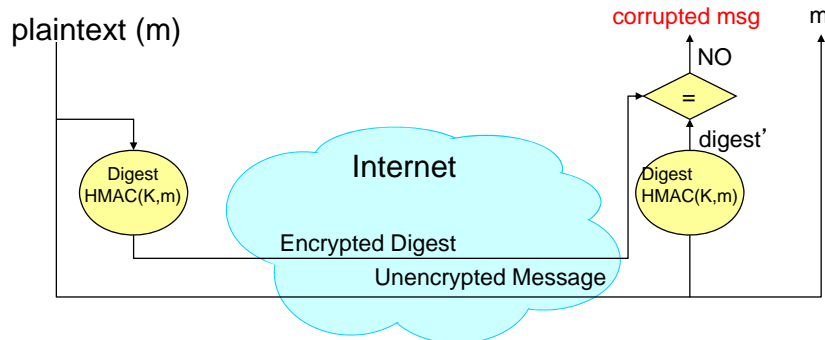
- **Basic building block for integrity: cryptographic hashing**
  - Associate hash with byte-stream, receiver verifies match
    - » Assures data hasn't been modified, either accidentally - or maliciously
- **Approach:**
  - **Sender computes a secure digest of message  $m$  using  $H(x)$** 
    - »  $H(x)$  is a publicly known hash function
    - » Digest  $d = \text{HMAC}(K, m) = H(K \parallel H(K \parallel m))$
    - »  $\text{HMAC}(K, m)$  is a hash-based message authentication function
  - **Send digest  $d$  and message  $m$  to receiver**
  - **Upon receiving  $m$  and  $d$ , receiver uses shared secret key,  $K$ , to recompute  $\text{HMAC}(K, m)$  and see whether result agrees with  $d$**

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.66

## Using Hashing for Integrity



Can encrypt  $m$  for confidentiality

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.67

## Standard Cryptographic Hash Functions

- **MD5 (Message Digest version 5)**
  - Developed in 1991 (Rivest), produces 128 bit hashes
  - Widely used (RFC 1321)
  - Broken (1996-2008): attacks that find collisions
- **SHA-1 (Secure Hash Algorithm)**
  - Developed in 1995 (NSA) as MD5 successor with 160 bit hashes
  - Widely used (SSL/TLS, SSH, PGP, IPSEC)
  - Broken in 2005, government use discontinued in 2010
- **SHA-2 (2001)**
  - Family of SHA-224, SHA-256, SHA-384, SHA-512 functions
- **HMAC's are secure even with older "insecure" hash functions**

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.68

## Asymmetric Encryption (Public Key)

- Idea: use two different keys, one to encrypt ( $e$ ) and one to decrypt ( $d$ )
  - A **key pair**
- Crucial property: knowing  $e$  does not give away  $d$
- Therefore  $e$  can be public: everyone knows it!
- If Alice wants to send to Bob, she fetches Bob's public key (say from Bob's home page) and encrypts with it
  - Alice can't decrypt what she's sending to Bob ...
  - ... but then, neither can anyone else (except Bob)

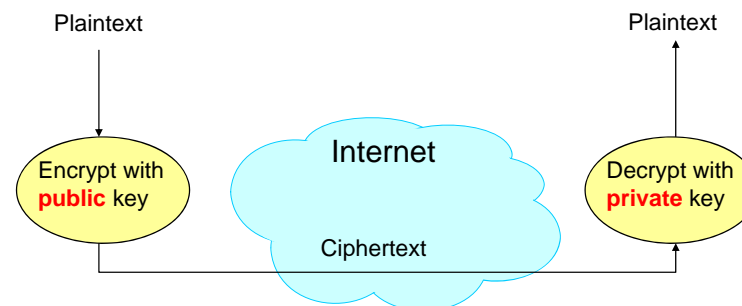
4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.69

## Public Key / Asymmetric Encryption

- Sender uses receiver's **public** key
  - Advertised to everyone
- Receiver uses complementary **private** key
  - Must be kept secret



4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.70

## Public Key Cryptography

- Invented in the 1970s
  - Revolutionized cryptography
  - (Was actually invented earlier by British intelligence)
- How can we construct an encryption/decryption algorithm using a key pair with the public/private properties?
  - Answer: Number Theory
- Most fully developed approach: RSA
  - Rivest / Shamir / Adleman, 1977; RFC 3447
  - Based on modular multiplication of very large integers
  - Very widely used (e.g., ssh, SSL/TLS for https)

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.71

## Properties of RSA

- Requires generating large, random prime numbers
  - Algorithms exist for quickly finding these (probabilistic!)
- Requires exponentiating very large numbers
  - Again, fairly fast algorithms exist
- Overall, much slower than symmetric key crypto
  - One general strategy: use public key crypto to exchange a (short) symmetric session key
    - » Use that key then with AES or such
- How difficult is recovering  $d$ , the private key?
  - Equivalent to finding prime factors of a large number
    - » Many have tried - believed to be very hard (= brute force only)
    - » (Though quantum computers can do so in polynomial time!)

4/27/15

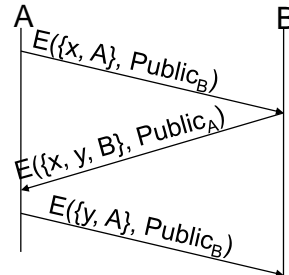
Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.72



## Simple Public Key Authentication

- Each side need only to know the other side's public key
  - No secret key need be shared
- A encrypts a nonce (random num.)  $x$ 
  - Avoid **replay attacks**, e.g., attacker impersonating client or server
- B proves it can recover  $x$
- A can authenticate itself to B in the same way
- *Many more details to make this work securely in practice!*



Notation:  $E(m,k)$  – encrypt message  $m$  with key  $k$

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.73

## Summary (1/2)

- **Distributed File System:**
  - Transparent access to files stored on a remote disk
  - Caching for performance
- **Cache Consistency:** Keeping client caches consistent with one another
  - If multiple clients, some reading and some writing, how do stale cached copies get updated?
  - NFS: check periodically for changes
  - AFS: clients register callbacks to be notified by server of changes
- **Remote Procedure Call (RPC):** Call procedure on remote machine
  - Provides same interface as procedure
  - Automatic packing and unpacking of arguments (in stub)
- **VFS:** Virtual File System layer
  - Provides mechanism which gives same system call interface for different types of file systems

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.74

## Summary

- **Key-Value Store:**
  - Two operations
    - »  $put(key, value)$
    - »  $value = get(key)$
  - Challenges
    - » Fault Tolerance → replication
    - » Scalability → serve  $get()$ 's in parallel; replicate/cache hot tuples
    - » Consistency → quorum consensus to improve  $put()$  performance
- **Distributed identity:** Use cryptography
- **Symmetrical (or Private Key) Encryption**
  - Single Key used to encode and decode
  - Introduces key-distribution problem
- **Public-Key Encryption**
  - Two keys: a public key and a private key
  - Slower than private key, but simplifies key-distribution
- **Secure Hash Function**
  - Used to summarize data
  - Hard to find another block of data with same hash

4/27/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 23.75