

**CS162**  
**Operating Systems and**  
**Systems Programming**  
**Lecture 15**  
**Key-Value Storage, Network Protocols**

March 19, 2014  
Anthony D. Joseph  
<http://inst.eecs.berkeley.edu/~cs162>

### Review: Directory Structure

- How many disk accesses to resolve “/my/book/count”?
  - Read in file header for root (fixed spot on disk)
  - Read in first data block for root
    - » Table of file name/index pairs. Search linearly – ok since directories typically very small
  - Read in file header for “my”
  - Read in first data block for “my”; search for “book”
  - Read in file header for “book”
  - Read in first data block for “book”; search for “count”
  - Read in file header for “count”
- **Current working directory:** Per-address-space pointer to a directory (inode) used for resolving file names
  - Allows user to specify relative filename instead of absolute path (say CWD=“/my/book” can resolve “count”)

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.2

### Where are inodes stored?

- In early UNIX and DOS/Windows' FAT file system, headers stored in special array in outermost cylinders
  - Header not stored anywhere near the data blocks. To read a small file, seek to get header, seek back to data.
  - Fixed size, set when disk is formatted. At formatting time, a fixed number of inodes were created (They were each given a unique number, called an “inumber”)

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.3

### Where are inodes stored?

- Later versions of UNIX moved the header information to be closer to the data blocks
  - Often, inode for file stored in same “cylinder group” as parent directory of the file (makes an `ls` of that directory run fast)
  - Pros:
    - » UNIX BSD 4.2 puts a portion of the file header array on each cylinder. For small directories, can fit all data, file headers, etc. in same cylinder ⇒ no seeks!
    - » File headers much smaller than whole block (a few hundred bytes), so multiple headers fetched from disk at same time
    - » Reliability: whatever happens to the disk, you can find many of the files (even if directories disconnected)
  - Part of the Fast File System (FFS)
    - » General optimization to avoid seeks

3/19/2014

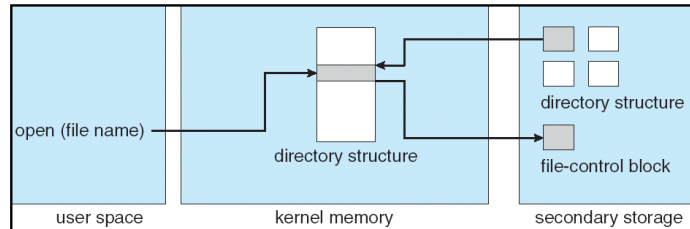
Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.4

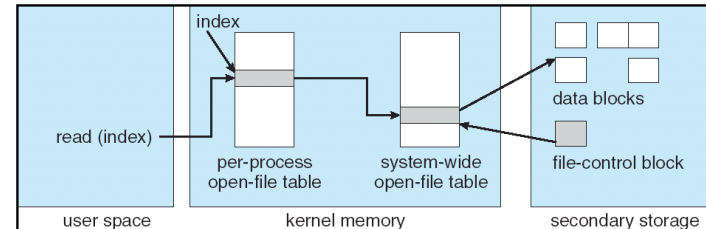
## In-Memory File System Structures



- Open system call:
  - Resolves file name, finds file control block (inode)
  - Makes entries in per-process and system-wide tables
  - Returns index (called “file handle”) in open-file table

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.5

## In-Memory File System Structures



- Read/write system calls:
  - Use file handle to locate inode
  - Perform appropriate reads or writes

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.6

## File System Summary

- File System:
  - Transforms blocks into Files and Directories
  - Optimize for access and usage patterns
  - Maximize sequential access, allow efficient random access
- File (and directory) defined by header, called “inode”
- 4.2 BSD Multilevel index files
  - Inode contains ptrs to actual blocks, indirect blocks, double indirect blocks, etc.
  - Optimizations for sequential access: start new files in open ranges of free blocks, rotational Optimization
- Naming: translate user-visible names to system resources
  - Directories used for naming for local file systems

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.7

## Goals for Today

- Key-Value Storage
  - Interface and Examples
  - Distributed Hash Tables
  - Challenges and Solutions
- Networking
  - What is a protocol?
  - Layering

**Many slides generated from Ion Stoica's lecture notes by Vern Paxson, and Scott Shenker.**

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.8

## Key-Value Storage

- Interface
  - **put**(key, value); // insert/write “value” associated with “key”
  - value = **get**(key); // get/read data associated with “key”
- Abstraction used to implement
  - A simpler and more scalable “database”
  - Content-addressable network storage (CANs)
- Can handle large volumes of data, e.g., Petabytes!
  - Need to distribute data over hundreds, even thousands of machines
  - Designed to be faster with lower overhead (additional storage) than conventional a DBMS

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.9

## Database (DBMS) Attributes

Databases require 4 properties:

- **Atomicity**: When an update happens, it is “all or nothing”
- **Consistency**: The state of various tables must be consistent (relations, constraints) at all times
- **Isolation**: Concurrent execution of transactions produces the same result as if they occurred sequentially
- **Durability**: Once committed, the results of a transaction persist against various problems like power failure etc.

These ACID properties ensure that data is protected even with complex updates and system failures

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.10

## CAP Theorem (Brewer, Gilbert, Lynch)

But we also have the CAP theorem for distributed systems:

**Consistency**: All nodes have the same view of the data

**Availability**: Every request receives a response of success or failure

**Partition Tolerance**: System continues even with loss of messages or part of the data nodes

The theorem states that **you cannot achieve all three at once**

Many systems therefore strive to implement two of the three properties – Key-Value stores often do this

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.11

## KV-stores and Relational Tables

KV-stores seem very simple indeed. They can be viewed as two-column (key, value) tables with a single key column

But they can be used to implement more complicated relational tables:

State	ID	Population	Area	Senator_1
Alabama	1	4,822,023	52,419	Sessions
Alaska	2	731,449	663,267	Begich
Arizona	3	6,553,255	113,998	Boozman
Arkansas	4	2,949,131	53,178	Flake
California	5	38,041,430	163,695	Boxer
Colorado	6	5,187,582	104,094	Bennet
...	...			

↑  
Index

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.12

## KV-stores and Relational Tables

The KV-version of the previous table includes one table indexed by the actual key, and others by an ID

State	ID	ID	Population	ID	Area	ID	Senator_1
Alabama	1	1	4,822,023	1	52,419	1	Sessions
Alaska	2	2	731,449	2	663,267	2	Begich
Arizona	3	3	6,553,255	3	113,998	3	Boozman
Arkansas	4	4	2,949,131	4	53,178	4	Flake
California	5	5	38,041,430	5	163,695	5	Boxer
Colorado	6	6	5,187,582	6	104,094	6	Bennet
...	...	...	...	...	...	...	...

Index

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.13

## KV-stores and Relational Tables

You can add indices with new KV-tables:

Thus KV-tables are used for **column-based storage**, as opposed to row-based storage typical in older DBMS

State	ID	ID	Population	...	Senator_1	ID
Alabama	1	1	4,822,023		Sessions	1
Alaska	2	2	731,449		Begich	2
Arizona	3	3	6,553,255		Boozman	3
Arkansas	4	4	2,949,131		Flake	4
California	5	5	38,041,430		Boxer	5
Colorado	6	6	5,187,582		Bennet	6
...	...	...	...		...	...

Index

Index\_2

OR: the value field can contain complex data (next page):

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.14

## 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, ...)
- iTunes Match:
  - Key: Movie/song name
  - Value: Movie, Song
- Distributed file systems
  - Key: Block ID
  - Value: Block

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.15

## System Examples

- Google File System, Hadoop Dist. File Systems (HDFS)
- Amazon
  - Dynamo: 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” (Facebook)
- Memcached: in-memory key-value store for small chunks of arbitrary data (strings, objects)
- eDonkey/eMule: peer-to-peer sharing system

3/19/2014

Anthony D. Joseph

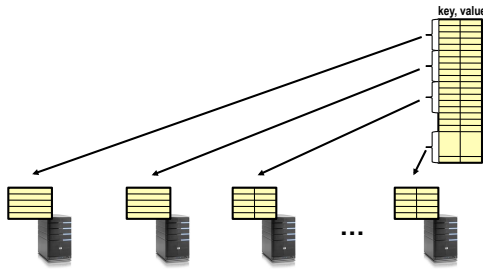
CS162

©UCB Spring 2014

Lec 15.16

## Key-Value Store

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



3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.17

## 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 several Gb/s

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.18

## 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

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.19

## Administrivia

- Midterm 2 is Monday April 28<sup>th</sup> 4-5:30
- Strategy tip: start studying now!
- Each week:
  - After class, review slides/slidescast
  - Pick a few related problems from a prior year midterm (over 35 exams from Sp'96 to Fa'13 are available)
  - See if you can finish the problems within the appropriate amount of time (look at time for exam and points per problem)
  - See me or your TA if you have any questions
  - More finals will be posted pending solutions!

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

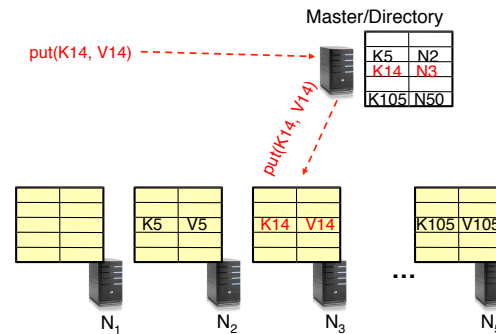
Lec 15.20

3 min Break

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.21

## Directory-Based Architecture

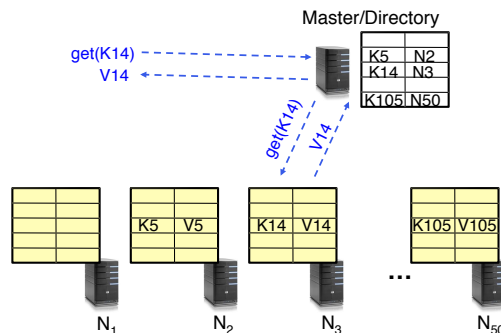
- Have a node maintain the mapping between **keys** and the **machines (nodes)** that store the **values** associated with the **keys**



3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.22

## Directory-Based Architecture

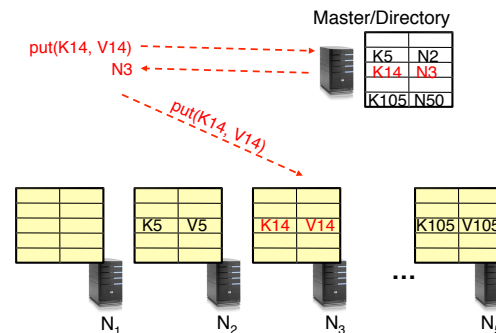
- Have a node maintain the mapping between **keys** and the **machines (nodes)** that store the **values** associated with the **keys**



3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.23

## Directory-Based Architecture

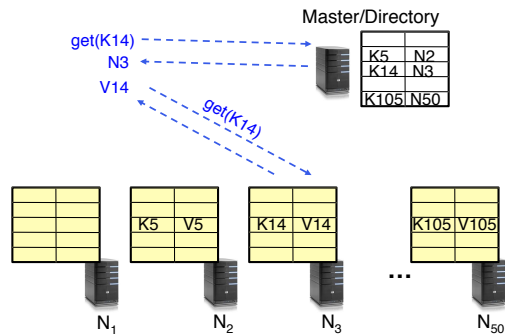
- Having the master relay the requests → **recursive query**
- Another method: **iterative query** (this slide)
  - Return node to requester and let requester contact node



3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.24

## Directory-Based Architecture

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



3/19/2014

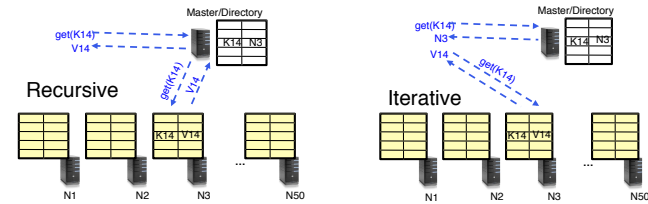
Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.25

## Discussion: Iterative vs. Recursive Query



### Recursive Query:

- Advantages:
  - » Faster (latency), 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 (latency), harder to enforce data consistency

3/19/2014

Anthony D. Joseph

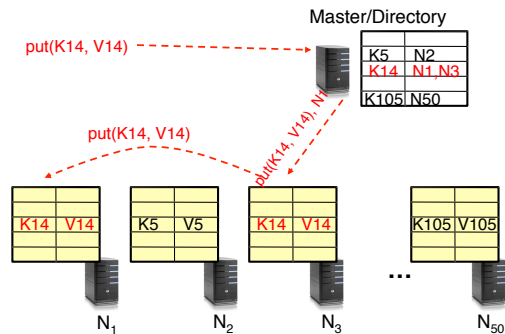
CS162

©UCB Spring 2014

Lec 15.26

## Fault Tolerance

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



3/19/2014

Anthony D. Joseph

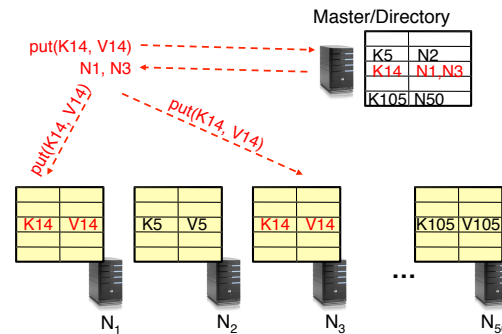
CS162

©UCB Spring 2014

Lec 15.27

## Fault Tolerance

- Again, we can have
  - **Recursive** replication (previous slide)
  - **Iterative** replication (this slide)



3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.28

## Scalability

- Storage: use more nodes
- Request Throughput:
  - Can serve requests from all nodes on which a value is stored in parallel
  - Large “values” can be broken into blocks (HDFS files are broken up this way)
  - 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? (use a structured p2p DHT)

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.29

## 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 failed node to other nodes

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.30

## Replication Challenges

- 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

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.31

## Consistency

- How close does a distributed system emulate a single machine in terms of read and write semantics?
- **Q:** Assume **put(K14, V14')** and **put(K14, V14'')** are concurrent, what value ends up being stored?
- **A:** assuming **put()** is atomic, then either **V14'** or **V14''**, right?
- **Q:** Assume a client calls **put(K14, V14)** and then **get(K14)**, what is the result returned by **get()**?
- **A:** It should be V14, right?
- Above semantics, not trivial to achieve in distributed systems

3/19/2014

Anthony D. Joseph

CS162

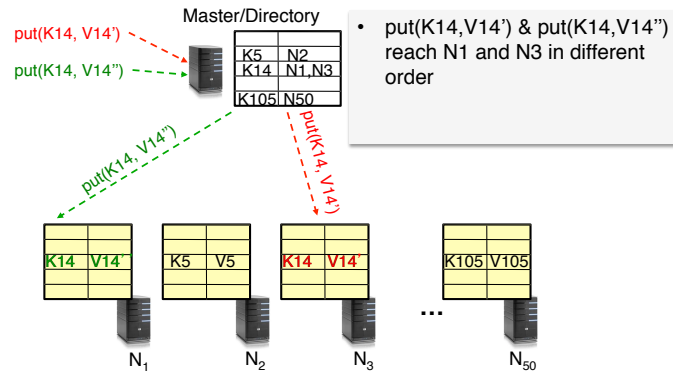
©UCB Spring 2014

Lec 15.32



## Concurrent Writes (Updates)

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



3/19/2014

Anthony D. Joseph

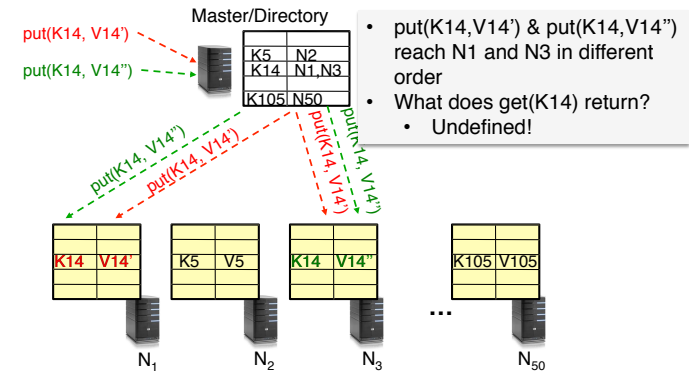
CS162

©UCB Spring 2014

Lec 15.33

## Concurrent Writes (Updates)

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



3/19/2014

Anthony D. Joseph

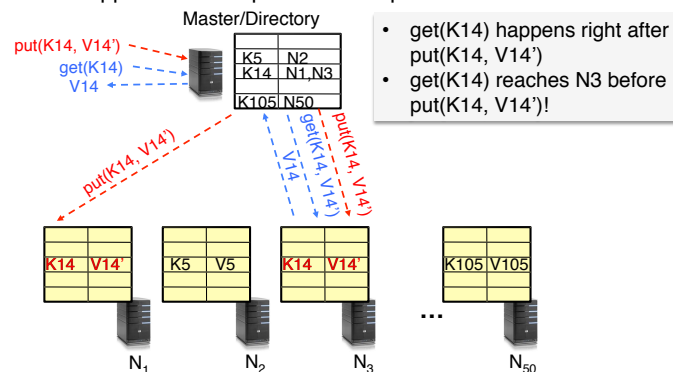
CS162

©UCB Spring 2014

Lec 15.34

## Read after Write

- Read not guaranteed to return value of latest write
  - Can happen if Master processes requests in different threads



3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.35

## 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 (later in the semester)
  - Eventual consistency: given enough time all updates will propagate through the system
    - One of the weakest forms of consistency; used by many systems in practice
- And many others: causal consistency, sequential consistency, strong consistency, ...

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.36

## Strong Consistency

- Assume Master serializes all operations
- Challenge: master becomes a bottleneck
  - Not addressed here
- Still want to improve performance of reads/writes → quorum consensus

3/19/2014

Anthony D. Joseph

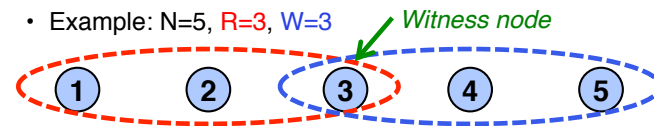
CS162

©UCB Spring 2014

Lec 15.37

## Quorum Consensus

- Improve **put()** and **get()** operation performance
- Define a replica set of size  $N$
- **put()** waits for acks 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
- Example:  $N=5$ ,  $R=3$ ,  $W=3$



3/19/2014

Anthony D. Joseph

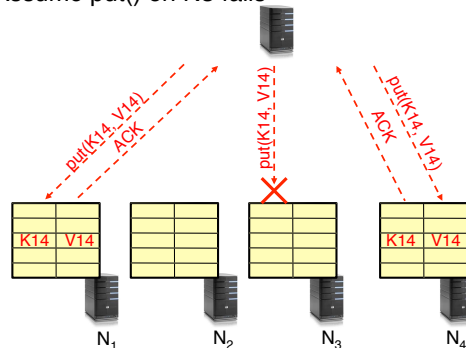
CS162

©UCB Spring 2014

Lec 15.38

## Quorum Consensus Example

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



3/19/2014

Anthony D. Joseph

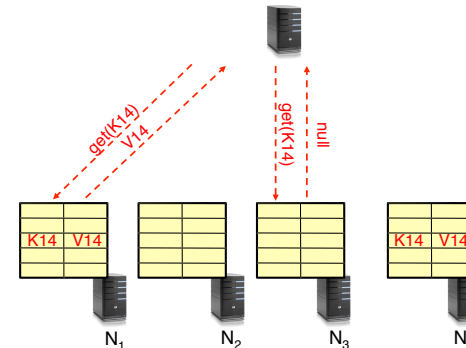
CS162

©UCB Spring 2014

Lec 15.39

## Quorum Consensus Example

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



3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.40

## Summary: Key-Value Store

- Very large scale storage systems
- 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/get performance

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.41

## Administrivia

- Project 2 code due **11:59pm on Thursday 3/20**
- Project 2 group evals due **11:59pm on Friday 3/21**
- Project 3 is being revised!
  - New version posted by end of spring break
  - Design doc not due until Tuesday 4/8

**Watch slip days!** Remember there are only 4 of these, after that there is an automatic (non-negotiable) 10% deduction for each day late. Projects 3 and 4 are challenging!

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.42

## Quiz 15.1: Key-Value Store

- Q1: True \_ False \_ Distributed Key-Value stores should always be Consistent, Available and Partition-Tolerant (CAP)
- Q2: True \_ False \_ On a single node, a key-value store can be implemented by a hash-table
- Q3: True \_ False \_ A Master can be a bottleneck point for a key-value store
- Q4: True \_ False \_ Iterative PUTs achieve lower throughput than recursive PUTs on a loaded system
- Q5: True \_ False \_ With quorum consensus, we can improve read performance at expense of write performance

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.43

**3 min Break**

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.44

### Quiz 15.1: Key-Value Store

- Q1: True \_ False \_ Distributed Key-Value stores should always be Consistent, Available and Partition-Tolerant (CAP)
- Q2: True \_ False \_ On a single node, a key-value store can be implemented by a hash-table
- Q3: True \_ False \_ A Master can be a bottleneck point for a key-value store
- Q4: True \_ False \_ Iterative PUTs achieve lower throughput than recursive PUTs on a loaded system
- Q5: True \_ False \_ With quorum consensus, we can improve read performance at expense of write performance

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.45

### Quiz 15.1: Key-Value Store

- Q1: True \_ False **X** Distributed Key-Value stores should always be Consistent, Available and Partition-Tolerant (CAP)
- Q2: True **X** False \_ On a single node, a key-value store can be implemented by a hash-table
- Q3: True **X** False \_ A Master can be a bottleneck point for a key-value store
- Q4: True \_ False **X** Iterative PUTs achieve lower throughput than recursive PUTs on a loaded system
- Q5: True **X** False \_ With quorum consensus, we can improve read performance at expense of write performance

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.46

### 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

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.47

### 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**
  8. Callee: "Yeah, blah blah blah ..." **pause**
  9. Caller: Bye
  10. Callee: Bye
  11. Hang up

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.48

## Examples of Protocols in Human Interactions

Asking a question

1. Raise your hand
2. Wait to be called on
3. Or: wait for speaker to **pause** and vocalize

3/19/2014

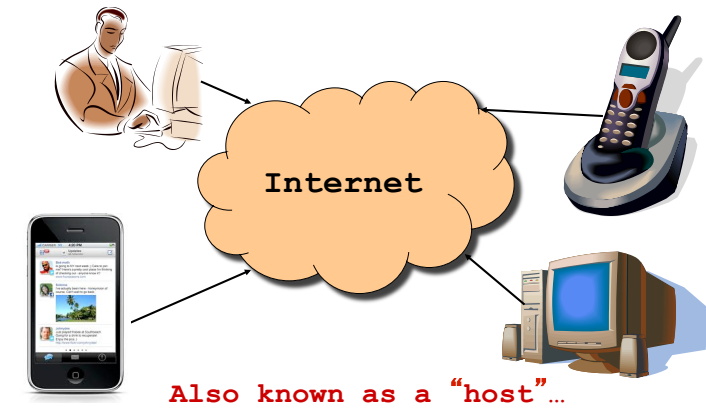
Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.49

## End System: Computer on the 'Net



3/19/2014

Anthony D. Joseph

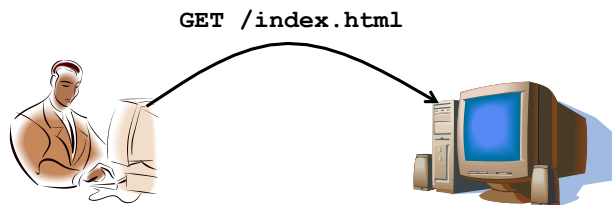
CS162

©UCB Spring 2014

Lec 15.50

## Clients and Servers

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



3/19/2014

Anthony D. Joseph

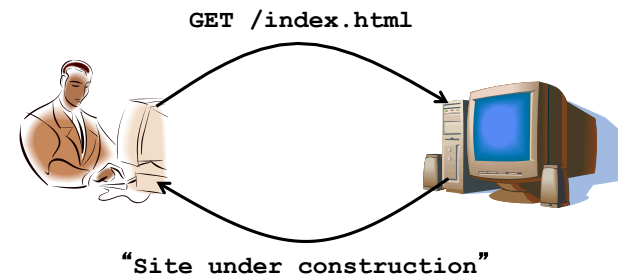
CS162

©UCB Spring 2014

Lec 15.51

## 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



3/19/2014

Anthony D. Joseph

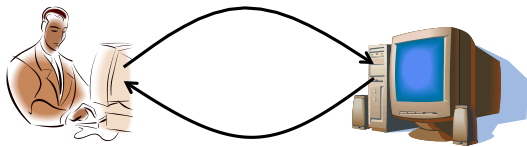
CS162

©UCB Spring 2014

Lec 15.52

## 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



3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.53

## 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**

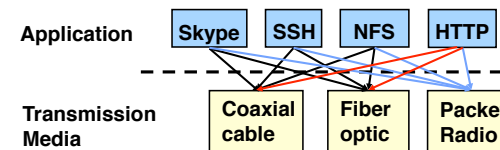
3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.54

## 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?

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.55

## The Problem (cont'd)

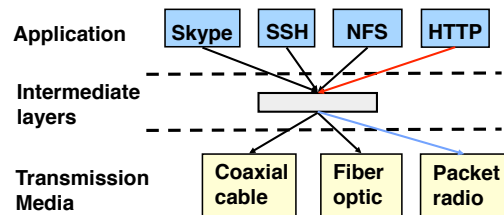


- Re-implement every application for every technology?
- No! But how does the Internet design avoid this?

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.56

## 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”



3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.57

## Software System Modularity

Partition system into modules & abstractions:

- Well-defined interfaces give flexibility
  - **Hides** implementation - thus, it can be freely changed
  - Extend functionality of system by adding new modules
- E.g., libraries encapsulating set of functionality
- E.g., programming language + compiler abstracts away not only how the particular CPU works ...
  - ... but also the **basic computational model**
- Well-defined interfaces hide information
  - Present high-level **abstractions**
  - **But can impair performance**

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.58

## Network System Modularity

Like software modularity, but:

- Implementation distributed across many machines (routers and hosts)
- Must decide:
  - How to break system into modules:
    - » **Layering**
  - What functionality does each module implement:
    - » **End-to-End Principle**: don't put it in the network if you can do it in the endpoints.
- We will address these choices next lecture

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.59

## Layering: A Modular Approach

- Partition the system
  - Each layer **solely** relies on services from layer below
  - Each layer **solely** exports services to layer above
- Interface between layers defines interaction
  - Hides implementation details
  - Layers can change without disturbing other layers

3/19/2014 Anthony D. Joseph CS162 ©UCB Spring 2014 Lec 15.60

## Protocol Standardization

- Ensure communicating hosts speak the same protocol
  - Standardization to enable multiple implementations
  - Or, the same folks have to write all the software
- Standardization: Internet Engineering Task Force
  - Based on working groups that focus on specific issues
  - Produces “Request For Comments” (RFCs)
    - » Promoted to standards via rough consensus and running code
  - IETF Web site is <http://www.ietf.org/>
  - RFCs archived at <http://www.rfc-editor.org/>
- De facto standards: same folks writing the code
  - P2P file sharing, Skype, <your protocol here>...

3/19/2014

Anthony D. Joseph

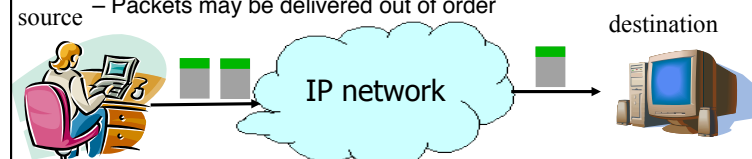
CS162

©UCB Spring 2014

Lec 15.61

## Example: The Internet Protocol (IP): “Best-Effort” Packet Delivery

- Datagram packet switching
  - Send data in packets
  - Header with source & destination address
- Service it provides:
  - Packet arrives quickly (if it does)
  - Packets may be lost
  - Packets may be corrupted
  - Packets may be delivered out of order



3/19/2014

Anthony D. Joseph

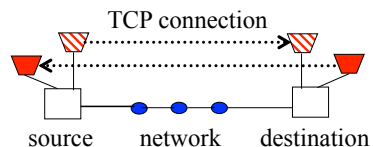
CS162

©UCB Spring 2014

Lec 15.62

## Example: Transmission Control Protocol (TCP)

- Communication service
  - Ordered, reliable byte stream
  - Simultaneous transmission in both directions
- Key mechanisms at end hosts
  - Retransmit lost and corrupted packets
  - Discard duplicate packets and put packets in order
  - **Flow control** to avoid overloading the receiver buffer
  - **Congestion control** to adapt sending rate to network load



3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.63

## Quiz 15.2: Protocols

- Q1: True \_ False \_ Protocols specify the syntax and semantics of communication
- Q2: True \_ False \_ Protocols specify the implementation
- Q3: True \_ False \_ Layering helps to improve application performance
- Q4: True \_ False \_ “Best Effort” packet delivery ensures that packets are delivered in order
- Q5: True \_ False \_ In p2p systems a node is both a client and a server
- Q6: True \_ False \_ TCP ensures that each packet is delivered within a predefined amount of time

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.64



## Quiz 15.2: Protocols

- Q1: True ☒ False ☐ Protocols specify the syntax and semantics of communication
- Q2: True ☐ False ☒ Protocols specify the implementation
- Q3: True ☐ False ☒ Layering helps to improve application performance
- Q4: True ☐ False ☒ “Best Effort” packet delivery ensures that packets are delivered in order
- Q5: True ☒ False ☐ In p2p systems a node is both a client and a server
- Q6: True ☐ False ☒ TCP ensures that each packet is delivered within a predefined amount of time

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.65

## Summary

- Important roles of
  - Protocols, standardization
  - Clients, servers, peer-to-peer
- A layered architecture is a powerful means for organizing complex networks
  - But, layering has its drawbacks too
- Next lecture
  - Layering
  - End-to-End arguments (please read the paper before lecture!)

3/19/2014

Anthony D. Joseph

CS162

©UCB Spring 2014

Lec 15.66