

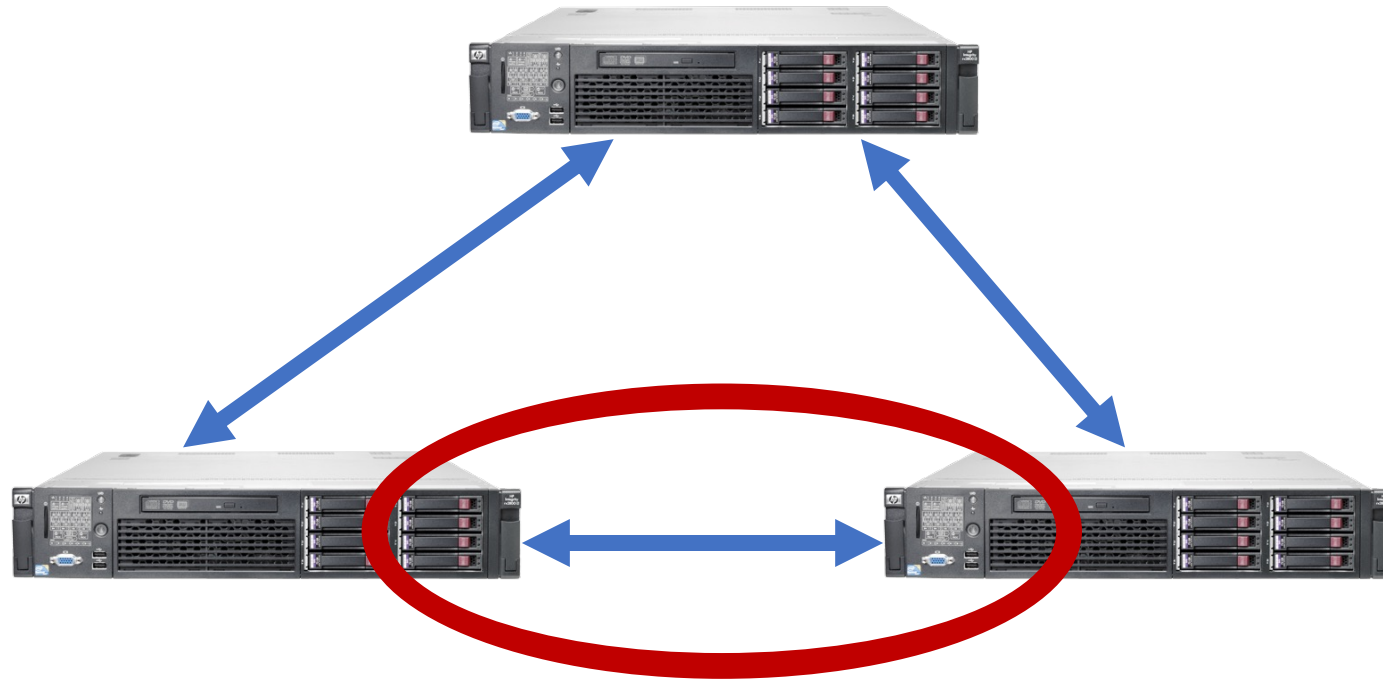
Introduction + Course Overview, MapReduce case study



CS 240: *Computing Systems and Concurrency*
Lecture 1

Marco Canini

Distributed Systems, What?



- 1) Multiple computers
- 2) Connected by a network
- 3) Doing something together

Distributed Systems, Why?

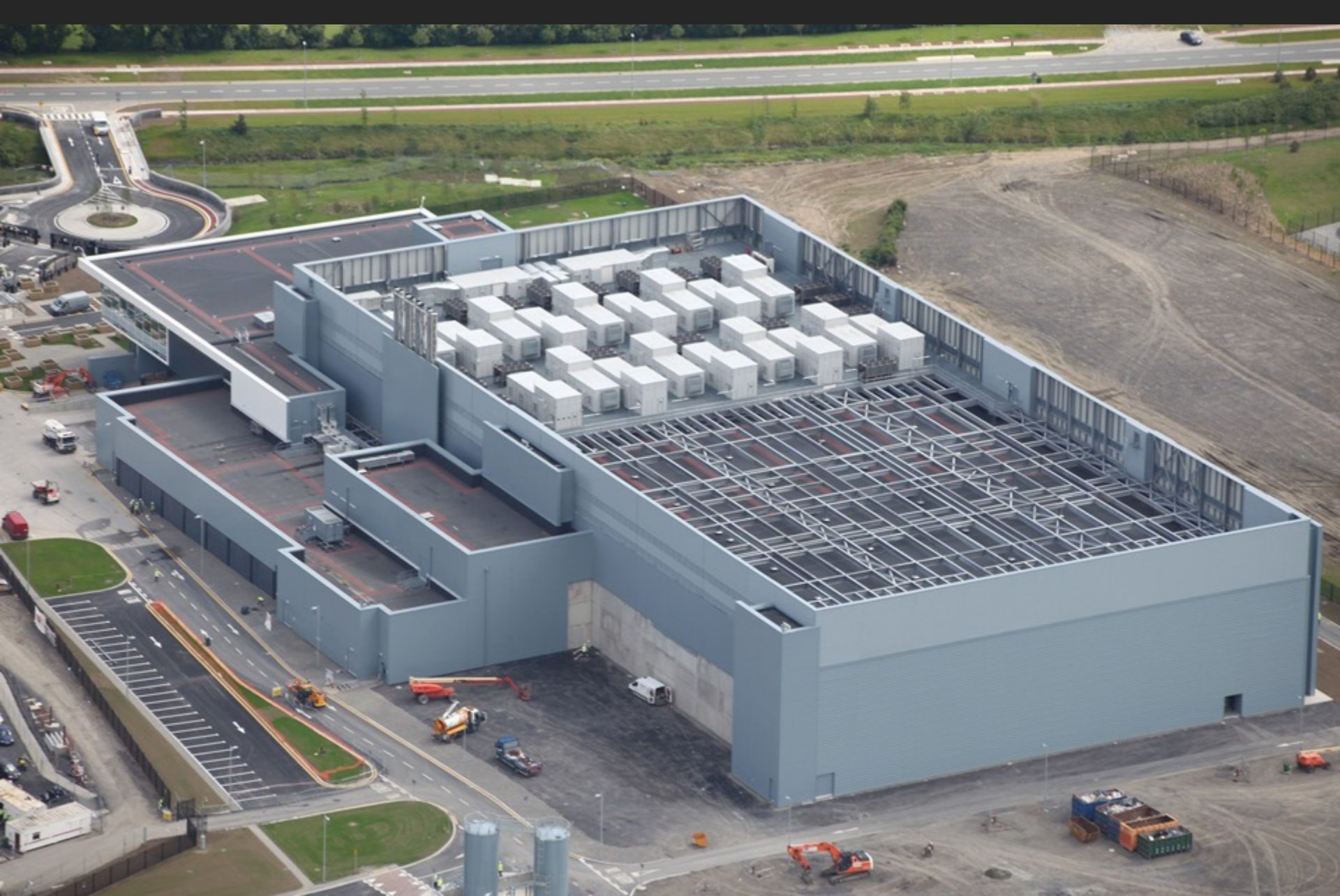
- Or, why not 1 computer to rule them all?
- Failure
- Limited computation/storage/...
- Physical location



Backrub (Google) 1997

Google 2012

“The Cloud” is not amorphous



Microsoft

Google





Facebook







**100,000s of physical servers
10s MW energy consumption**

**Facebook Prineville:
\$250M physical infra, \$1B IT infra**



GDH DC @ KAUST
~10,000 servers
14.4 MW IT load
8,000 m² of DC space



The goal of “distributed systems”

- Service with higher-level abstractions/interface
 - e.g., file system, database, key-value store, programming model, RESTful web service, ...
- Hide complexity
 - Scalable (scale-out)
 - Reliable (fault-tolerant)
 - Well-defined semantics (consistent)
- Do “heavy lifting” so app developer doesn’t need to

What is a distributed system?

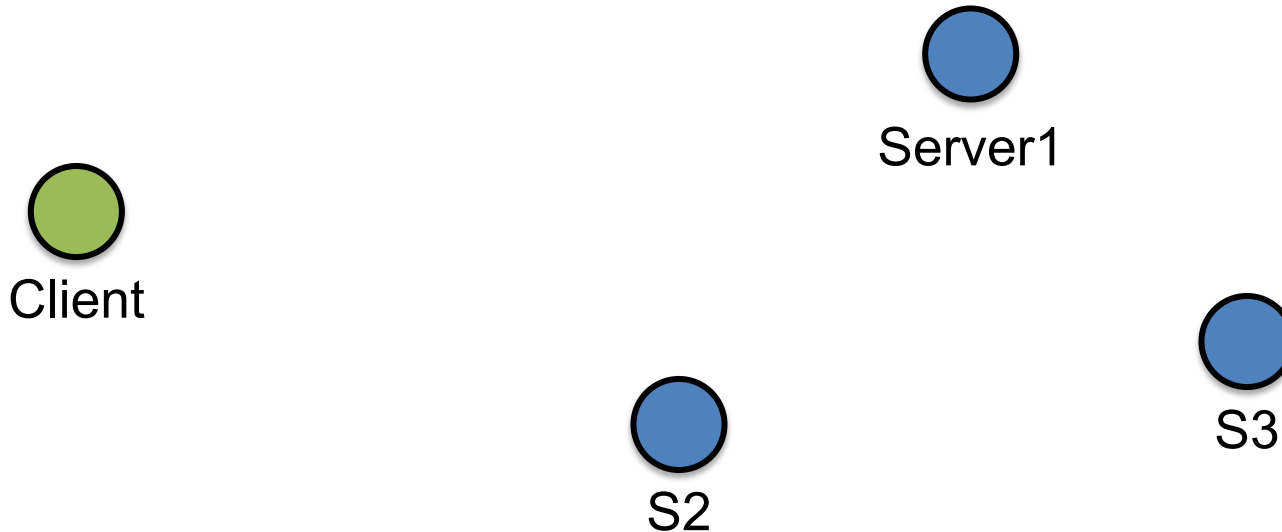
- “A collection of independent computers that appears to its users as a single coherent system”
- Features:
 - No shared memory
 - Message-based communication
 - Each runs its own local OS
 - Heterogeneity
- Ideal: to present a single-system image:
 - The distributed system “looks like” a single computer rather than a collection of separate computers

Distributed system characteristics

- To present a single-system image:
 - Hide internal organization, communication details
 - Provide uniform interface
- Easily expandable
 - Adding new computers is hidden from users
- Continuous availability
 - Failures in one component can be covered by other components

Example

- Assume a distributed storage
 - Clients can read and write files



System model

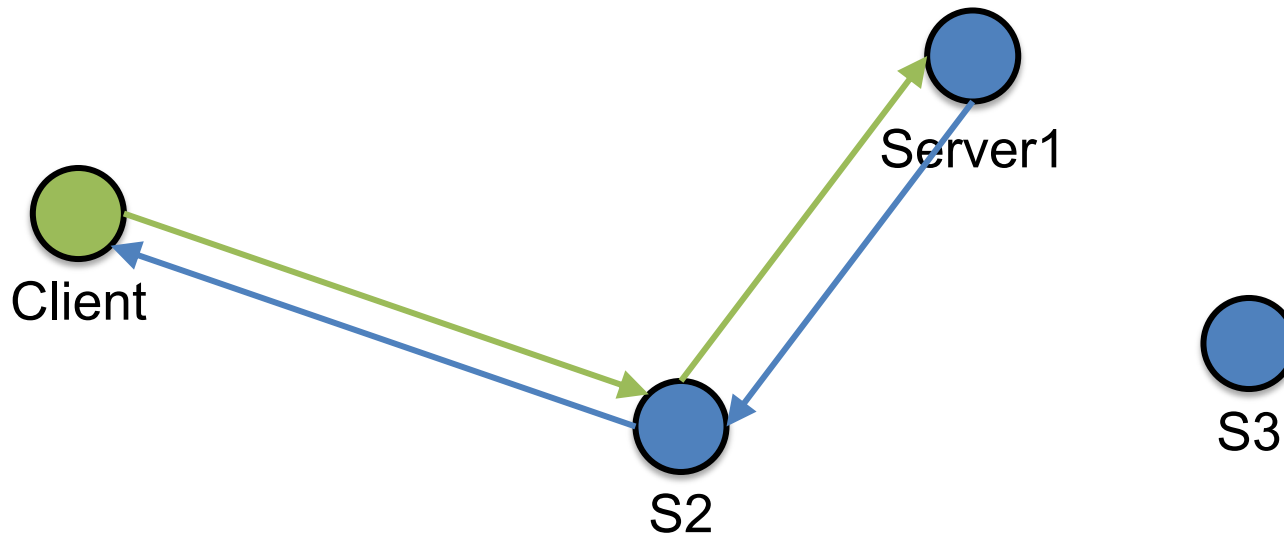
- N **processes** p_1, \dots, p_N in the system (no process failures)
 - Every process executes an algorithm
 - An automation with set of states, set of inputs, set of outputs and a state transition function $S \times I \rightarrow S \times O$
- There are two first-in, first-out, unidirectional **channels** between every process pair p_i and p_j
 - Call them **channel**(p_i, p_j) and **channel**(p_j, p_i)
 - All messages sent on channels arrive intact and in order
 - Channel cannot duplicate, create or modify messages

System model

- Message passing
- No failures (for now)
- Two possible timing assumptions
 1. Synchronous System
 2. Asynchronous System
 - No upper bound on **message delivery**
 - No bound on relative **process speeds**

Example execution

- Assume a distributed storage
 - Clients can read and write files

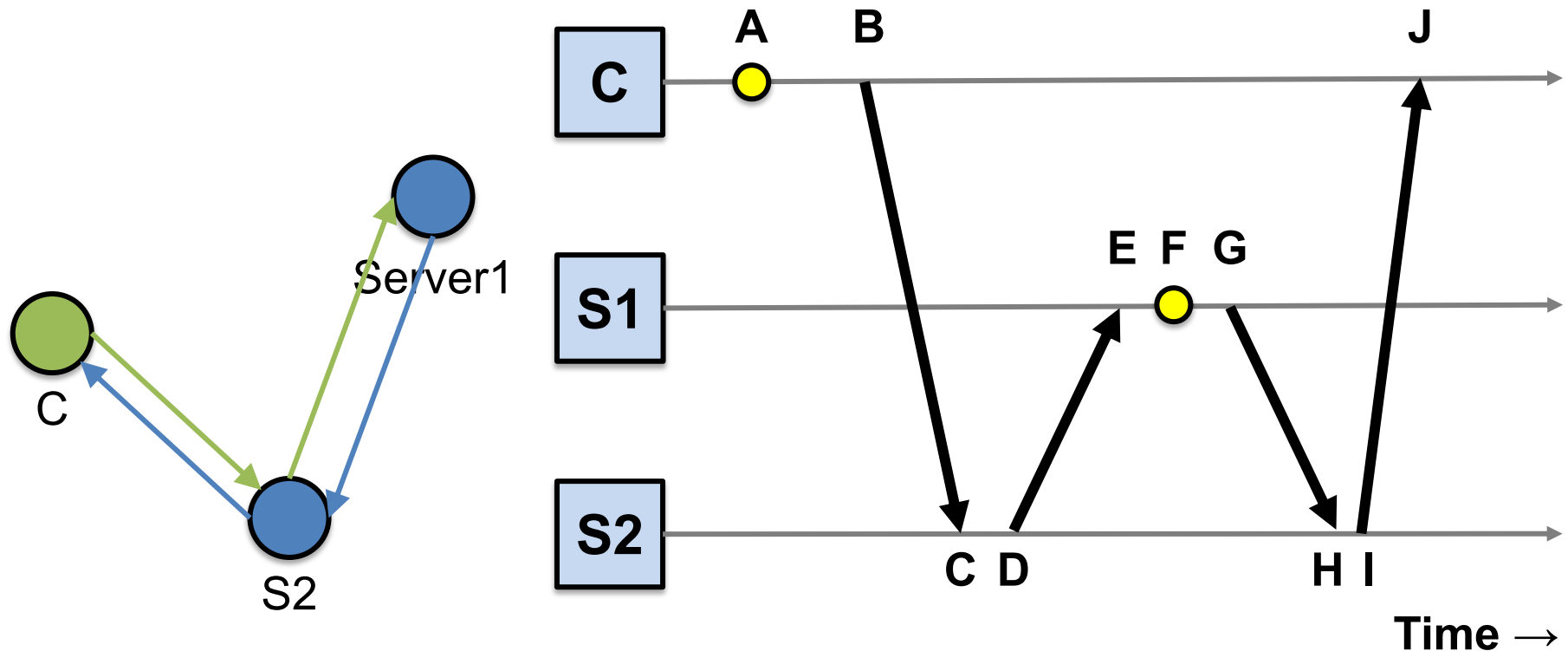


Execution of the system

- Processes execute sequences of events
 - events can be of 3 types: local, send or receive
- An execution (or run) is a sequence of events that respect the system-wide distributed algorithm
 - each process is consistent with the local sequences
 - a message is sent by a process only if its (local) algorithm prescribes it to do it given the preceding sequence of its inputs
 - every received message was previously sent, and no message is received twice

Space-Time diagrams

- A graphic representation of distributed execution



Common failure assumption

- Generally, a failure occurs when a process deviates from the algorithm assigned to it
- A process is *correct* if it never fails
- *crash* failure: the faulty process prematurely stops taking steps of its algorithm
- A typical assumption is that, in every possible execution out of N processes, at most $f < N$ can be faulty
- We call such a system *f-resilient*

Scalable systems in this class

- Scale computation across many machines
 - MapReduce
- Scale storage across many machines
 - Chord, Dynamo, COPS, Spanner

Fault tolerant systems in this class

- Retry on another machine
 - MapReduce
- Maintain replicas on multiple machines
 - Primary-backup replication
 - Paxos
 - RAFT
 - Bayou
 - Dynamo, COPS, Spanner

Range of abstractions and guarantees

- Eventual Consistency
 - Dynamo
- Causal Consistency
 - Bayou, COPS
- Linearizability
 - Paxos, RAFT, Primary-backup replication
- Strict Serializability
 - 2PL, Spanner

Summary

- Distributed Systems
 - Multiple machines doing something together
 - Pretty much everywhere and everything computing now
- “Systems”
 - Hide complexity and do the heavy lifting (i.e., **interesting!**)
 - Scalability, fault tolerance, guarantees

Course Overview

Philosophy and Recurring Themes

- Keep it real! This is the real world:
 - Things break. Components fail.
 - Latency matters. Can't beat speed of light.
 - Certain things are impossible. Need work arounds.
- How do we build systems that work at **very large scale** and **tolerate failures**?
- Given systems span many nodes, how do we enable different nodes to **agree** on “things” (e.g., time, order of operations, state of the system)?

Learning Objectives

- Reasoning about concurrency
- Reasoning about failure
- Reasoning about performance

- Building systems that correctly handle concurrency and failure

- Knowing specific system designs and design components

Course Goals

- Gain an understanding of the principles and techniques behind the design of modern, reliable, and high-performance systems
- In particular learn about distributed systems
 - Learn general systems principles (modularity, layering, naming, security, ...)
 - Practice implementing real, larger systems that must run in nasty environment
- One consequence: Must pass exams and projects independently as well as in total
 - Note, if you fail either you will not pass the class

Keep the Big Picture in Mind

- Course: many topics, grouped around key areas
- Might feel like lectures are disconnected...
- ... and first need to cover some background
- **Big Picture:**
 - real systems have complex requirements that span the concepts of multiple topics
 - E.g., we want fault tolerance, consistency and scalability

Course Organization

<http://sands.kaust.edu.sa/classes/CS240/F22/>

Learning the material: People

- Lecture
 - Professor Marco Canini
 - Slides available on course website
 - Office hours: by appointment
- TAs
 - Jihao Xin: W 16:30-18:00, 1-4409-WS26
- Main Q&A forum: www.campuswire.com
 - No anonymous (to instructors) posts or questions
 - Can send private messages to instructors



Learning the Material: Books

- Lecture notes!
- No required textbooks
- References on website available in the Library:
 - Programming reference:
 - *The Go Programming Language*. Alan Donovan and Brian Kernighan
 - Topic reference:
 - *Distributed Systems: Principles and Paradigms*. Andrew S. Tanenbaum and Maaten Van Steen
 - *Guide to Reliable Distributed Systems*. Kenneth Birman

Grading

- Four programming assignments (50% total)
 - 10% each for 1 & 2
 - 15% each for 3 & 4
- Two exams (50% total)
 - Midterm exam on October 24 (15%)
 - Final exam during exam period (35%)

Exams

- Test learning objectives mostly using designs covered in lectures
- And test knowledge of specific design patterns and designs
- Open book (but if you don't study it will create time pressure)
- Recipe for success:
 - Attend lecture and actively think through problems
 - Ask questions during lecture and afterwards in my office hours
 - Actively work through problems
 - Complete programming assignments
 - Study lecture materials for specific design patterns and designs
 - Run the system designs in your mind and see what happens

About Assignments

- Systems programming somewhat different from what you might have done before
 - Low-level (C / Go)
 - Often designed to run indefinitely (error handling must be rock solid)
 - Must be secure - horrible environment
 - Concurrency
 - Interfaces specified by documented protocols
- TAs' Office Hours
- Read: Dave Andersen's "[Software Engineering for System Hackers](#)"
 - Practical techniques designed to save you time & pain

Why use Go?

- Easy concurrency w/ goroutines (lightweight threads)
- Garbage collection and memory safety
- Libraries provide easy RPC
- Channels for communication between goroutines

Where is Go used?

- Google, of course!
- Docker (container management)
- CloudFlare (Content delivery Network)
- Digital Ocean (Virtual Machine hosting)
- Dropbox (Cloud storage/file sharing)
- ... and many more!

About Assignments

- Reinforce / demonstrate all learning objectives!
- 1: Sequential Map/Reduce (due September 14)
- 2: Distributed Map/Reduce (due September 21)
- 3-1: Raft Leader Election (due November 16)
- 3-2: Raft Log Consensus (due November 30)
- 4: Key-Value Storage Service (due December 8)

Programming Assignments

- Recipe for disaster
 - Start day assignment is due
 - Write code first, think later
 - Test doesn't pass => randomly flip some bits
 - Assume you know what program is doing

Programming Assignments

- Recipe for success
 - Start early (weeks early)
 - Think through a complete design
 - Progressively build out your design (using tests to help)
 - Checkpoint progress in git (and to gitlab) frequently
 - Debug, debug, debug
 - Verify program state is what you expect (print it out!)
 - Write your own smaller test cases
 - Reconsider your complete design
 - Attend office hours

Policies: Collaboration

- Working together important
 - Discuss course material
 - Work on problem debugging
- Parts **must** be your own work
 - Midterm, final, programming assignments
- What we hate to say: we run cheat checkers...
they work surprisingly well
- Please ***do not*** put code on ***public*** repositories

Policies: Write Your Own Code

Programming is an individual creative process. At first, discussions with friends is fine. When writing code, however, the program must be your own work.

Do not copy another person's programs, comments, README description, or any part of submitted assignment. This includes character-by-character transliteration but also derivative works. Cannot use another's code, etc. even while "citing" them.

Writing code for use by another or using another's code is academic fraud in context of coursework.

Do not publish your code e.g., on Github, during/after course!

Policies: Late Work

- 72 late hours to use throughout the semester
 - (but not beyond December 8)
- After that, each additional day late will incur a 10% lateness penalty
 - (1 min late counts as 1 day late)
- Submissions late by 3 days or more will no longer be accepted
 - (Fri and Sat count as days)
- In case of illness or extraordinary circumstance (e.g., emergency), talk to us early!

Summary

- Attend lecture, attend labs, think actively!
- Start programming assignments early, use the right strategy!

Case Study: MapReduce

(Data-parallel programming at scale)

Application: Word Count

```
SELECT count(word) FROM data  
GROUP BY word
```

```
cat data.txt
```

```
| tr -s '[:punct:][:space:]' '\n'
```

```
| sort | uniq -c
```


Using partial aggregation

1. Compute word counts from individual files
2. Then merge intermediate output
3. Compute word count on merged outputs

Using partial aggregation

1. In parallel, send to worker:
 - Compute word counts from individual files
 - Collect result, wait until all finished
2. Then merge intermediate output
3. Compute word count on merged intermediates

MapReduce: Programming Interface

`map(key, value) -> list(<k', v'>)`

- Apply function to (key, value) pair and produces set of intermediate pairs

`reduce(key, list<value>) -> <k', v'>`

- Applies aggregation function to values
- Outputs result

MapReduce: Programming Interface

```
map(key, value) :  
    for each word w in value:  
        EmitIntermediate(w, "1");  
  
reduce(key, list(values) :  
    int result = 0;  
    for each v in values:  
        result += ParseInt(v);  
    Emit(key, AsString(result));
```

MapReduce: Optimizations

`combine(list<key, value>) -> list<k, v>`

- Perform partial aggregation on mapper node:

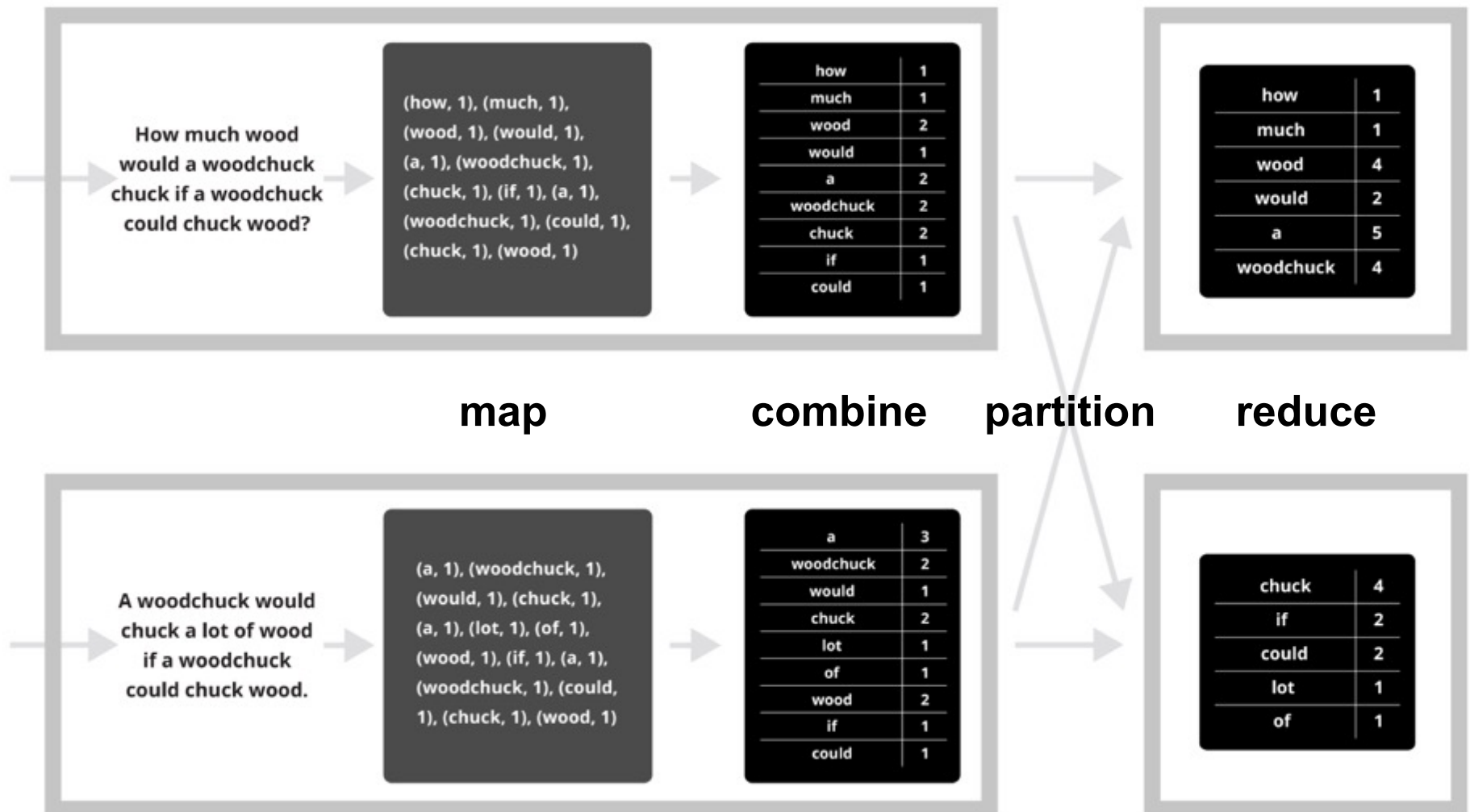
`<the, 1>, <the, 1>, <the, 1> → <the, 3>`

- `combine()` should be commutative and associative

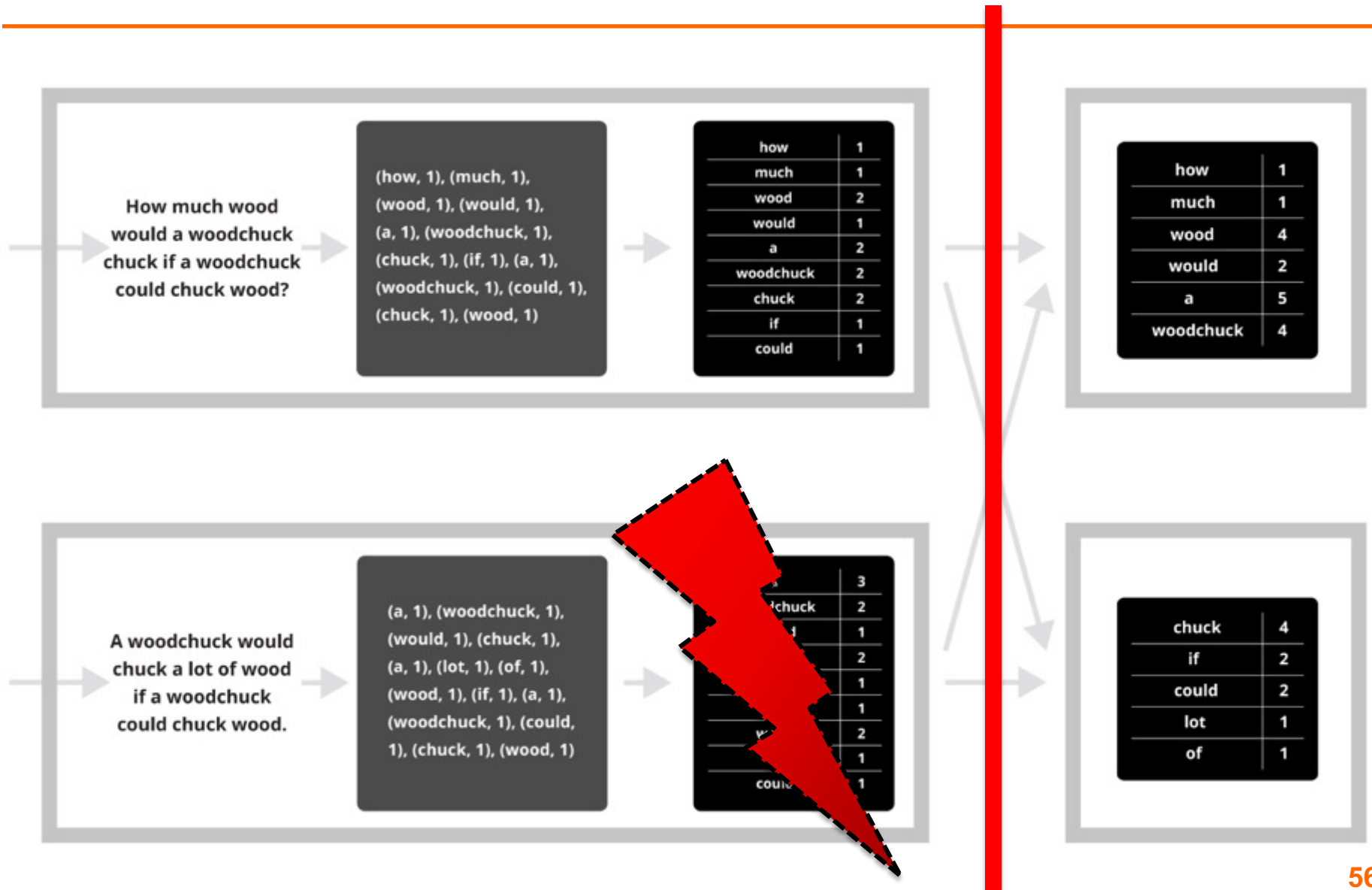
`partition(key, int) -> int`

- Need to aggregate intermediate vals with same key
- Given n partitions, map key to partition $0 \leq i < n$
- Typically via `hash(key) mod n`

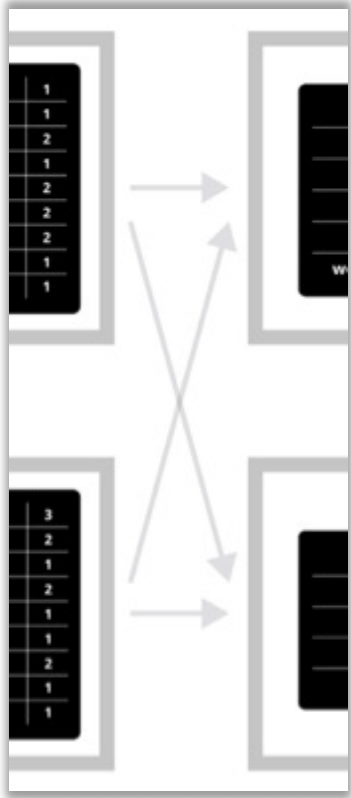
Putting it together...



Synchronization Barrier



Fault Tolerance in MapReduce

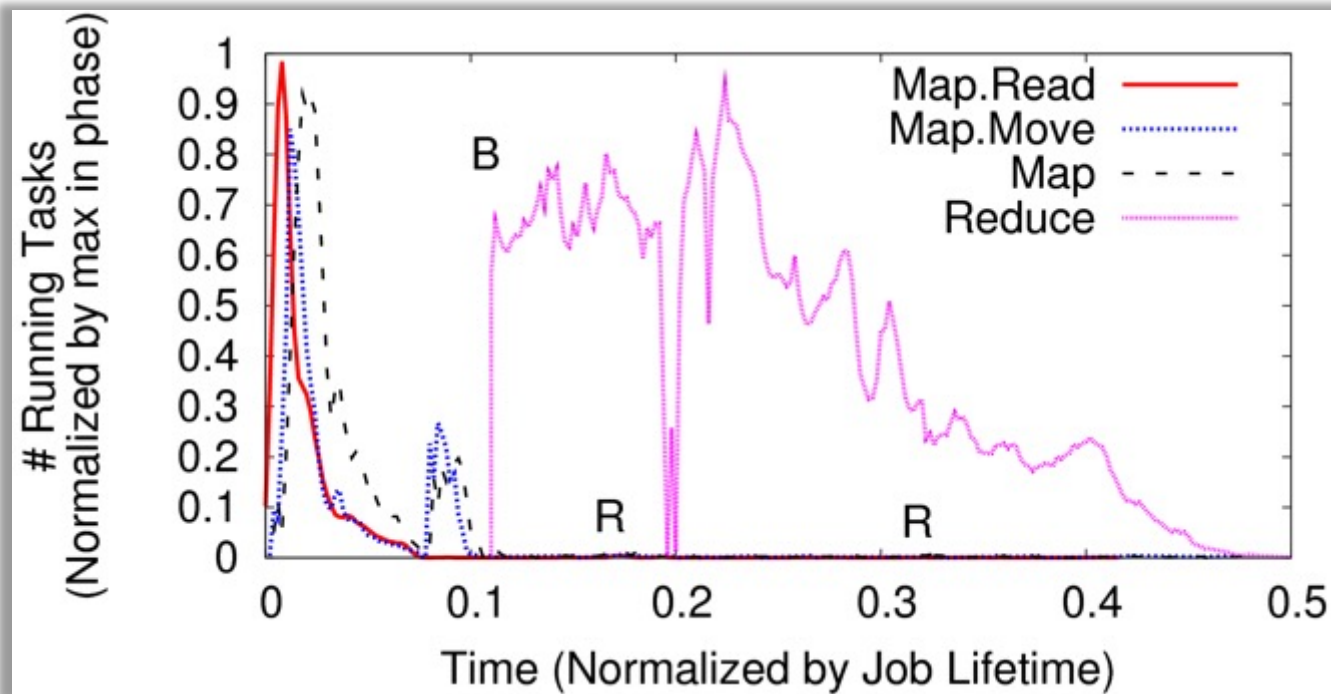


- Map worker writes intermediate output to local disk, separated by partitioning. Once completed, tells master node.
- Reduce worker told of location of map task outputs, pulls their partition's data from each mapper, execute function across data
- Note:
 - “All-to-all” shuffle b/w mappers and reducers
 - Written to disk (“materialized”) b/w *each* stage

Fault Tolerance in MapReduce

- Master node monitors state of system
 - If master failures, job aborts and client notified
- Map worker failure
 - Both in-progress/completed tasks marked as idle
 - Reduce workers notified when map task is re-executed on another map worker
- Reducer worker failure
 - In-progress tasks are reset to idle (and re-executed)
 - Completed tasks had been written to global file system

Straggler Mitigation in MapReduce



- Tail latency means some workers finish late
- For slow map tasks, execute in parallel on second map worker as “backup”, race to complete task

You'll build (simplified) MapReduce!

- Assignment 1: Sequential MapReduce
 - Learn to program in Go!
 - Due September 14
- Assignment 2: Distributed MapReduce
 - Learn Go's concurrency, network I/O, and RPCs
 - Due September 21

Conclusion

- Attend lecture, attend labs, think actively!
- Start programming assignments early, use the right strategy!
- Super cool distributed systems stuff starts Monday!

