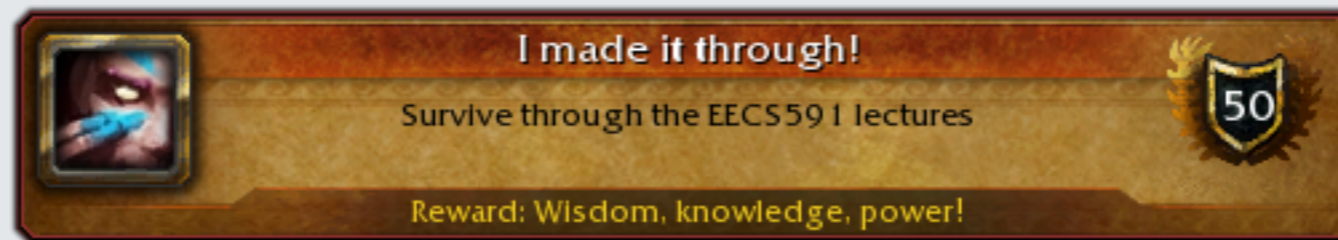# EECS591 CLASS REVIEW

# What a long, strange trip it's been…

# PART ONE: FUNDAMENTALS

# Two generals' problem
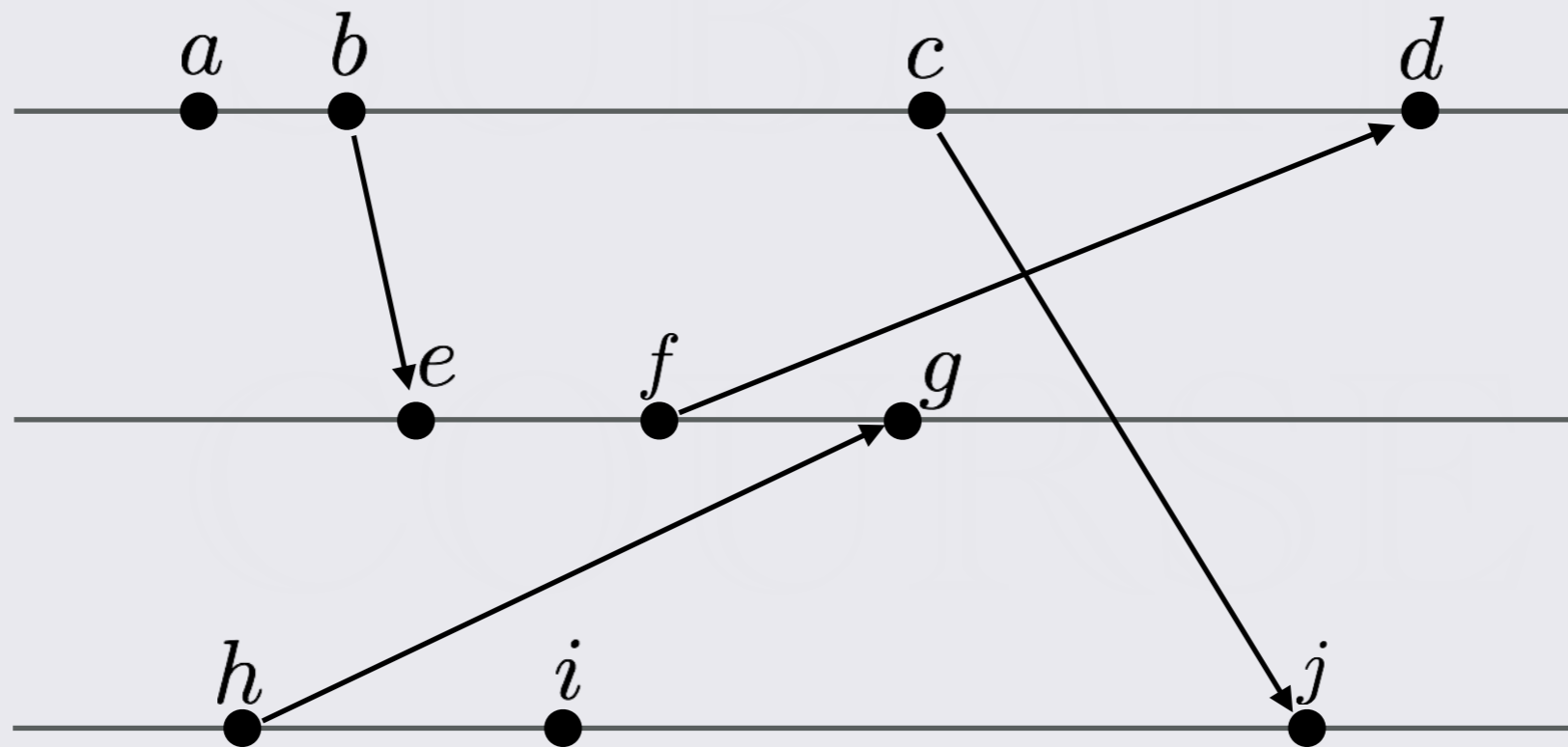
Both generals must attack together or face defeat



Communication is only by messengers sneaking through the valley
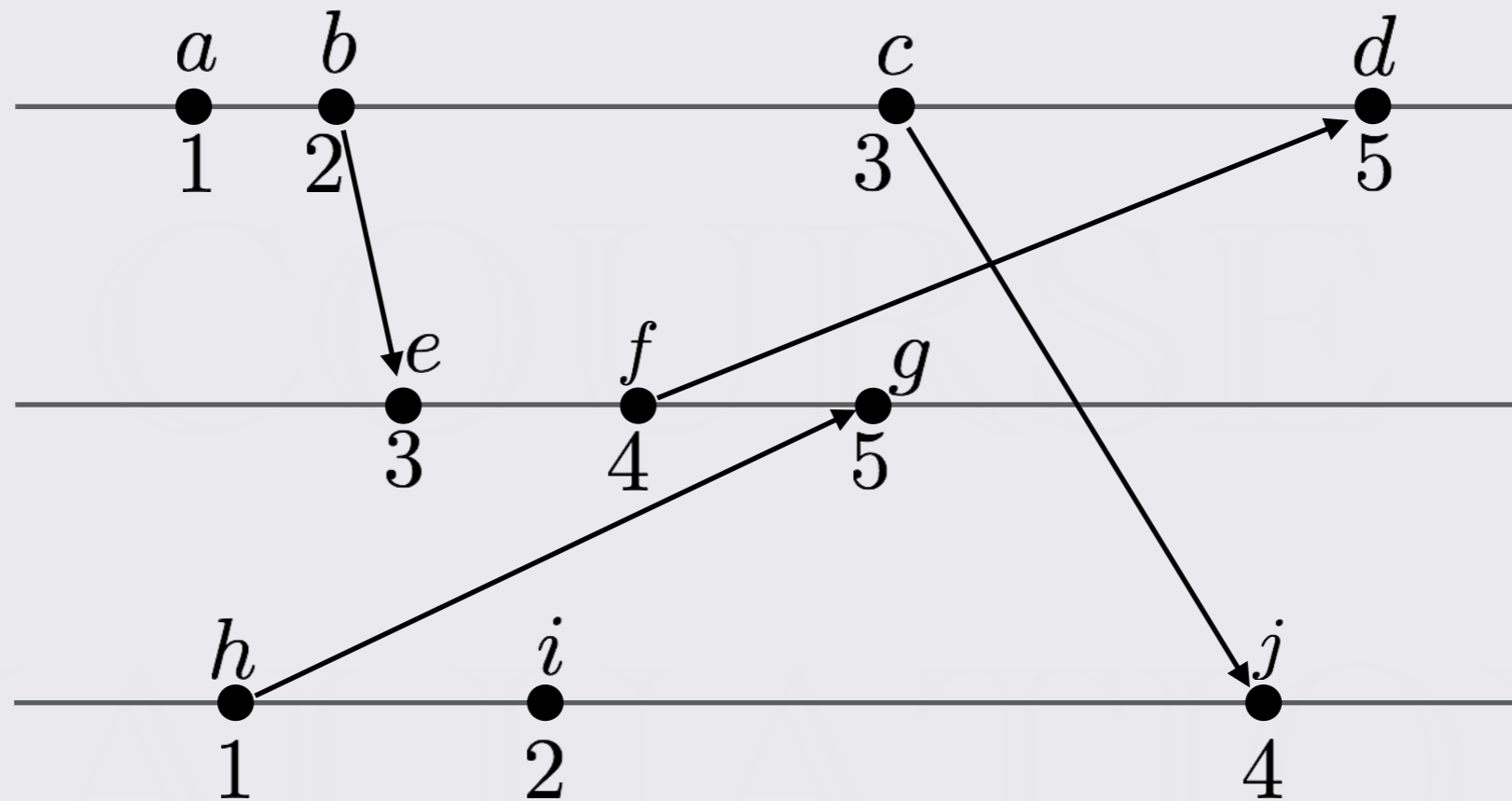
Messengers may not make it through…

# ORDERING EVENTS WITHOUT PHYSICAL CLOCKS



Question 1 (true or false)

a. $e \rightarrow d$

b. $a \rightarrow j$

c. $g \rightarrow b$

# LAMPORT CLOCKS
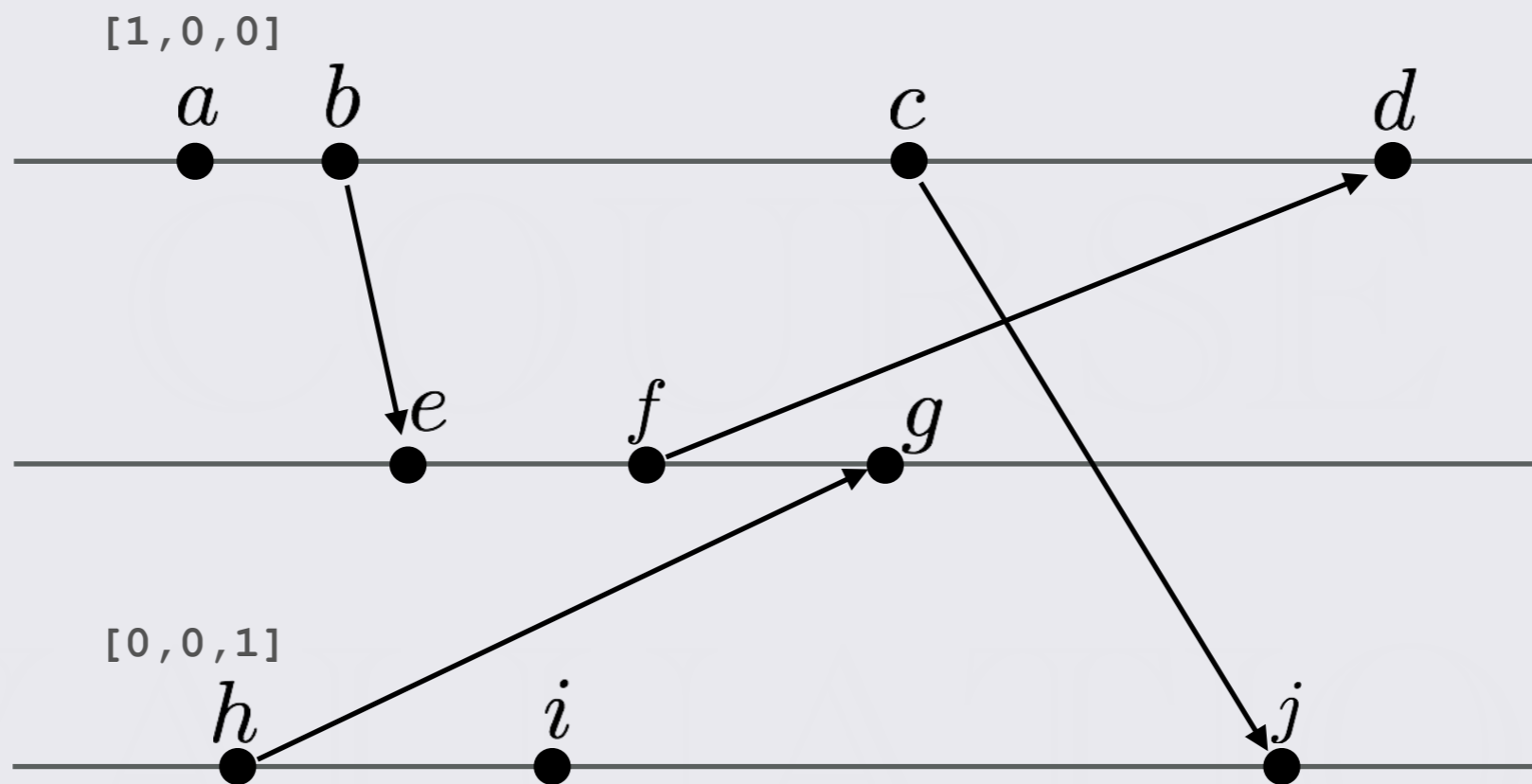


$$p \to q \Rightarrow LC(p) < LC(q)$$

the Clock condition

# VECTOR CLOCKS

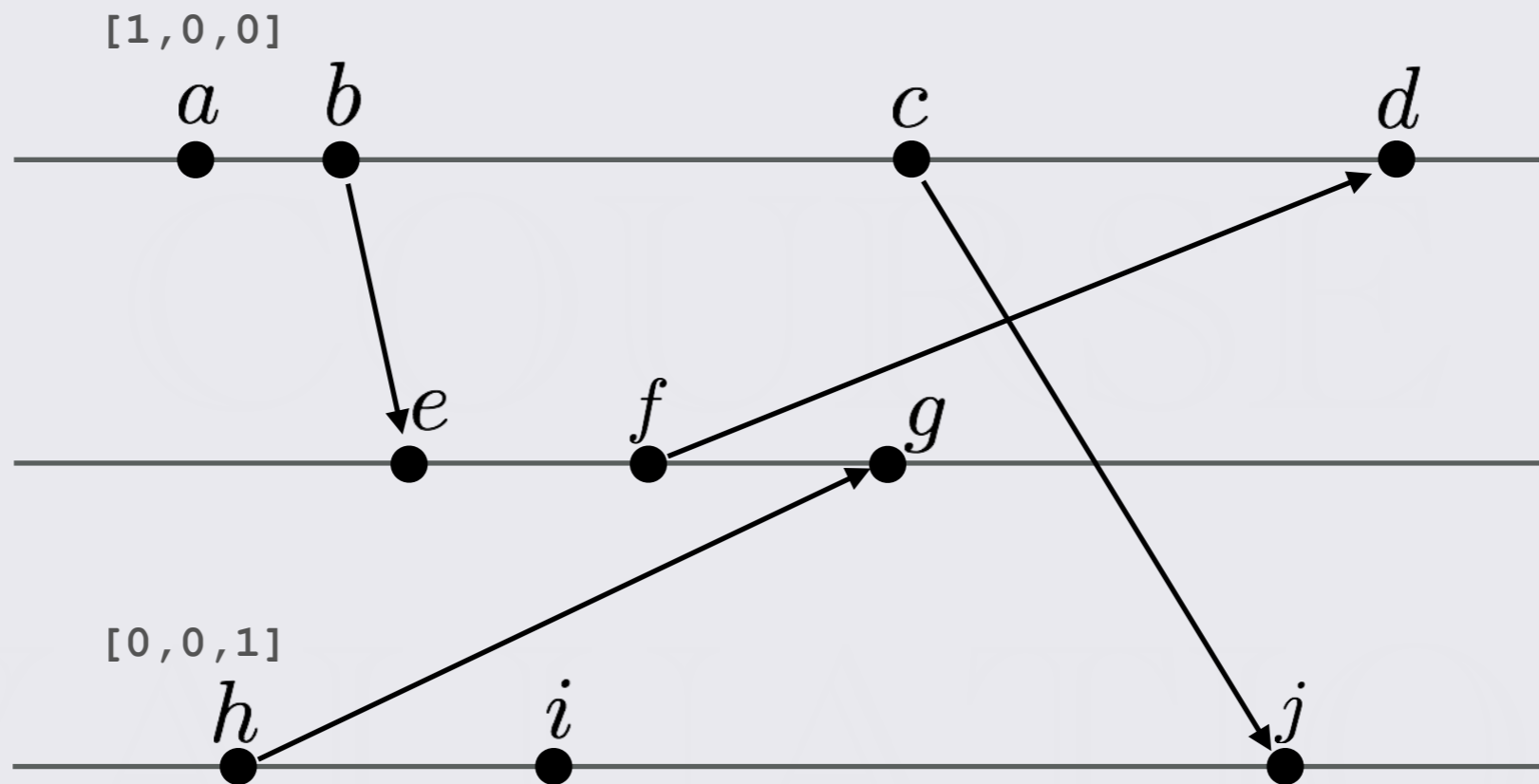$VC(e_i)[j]$ = number of events executed by process j that causally precede $e_i$



$$p \rightarrow q \Leftrightarrow LC(p) < LC(q)$$
Strong clock condition

# Vector clocks

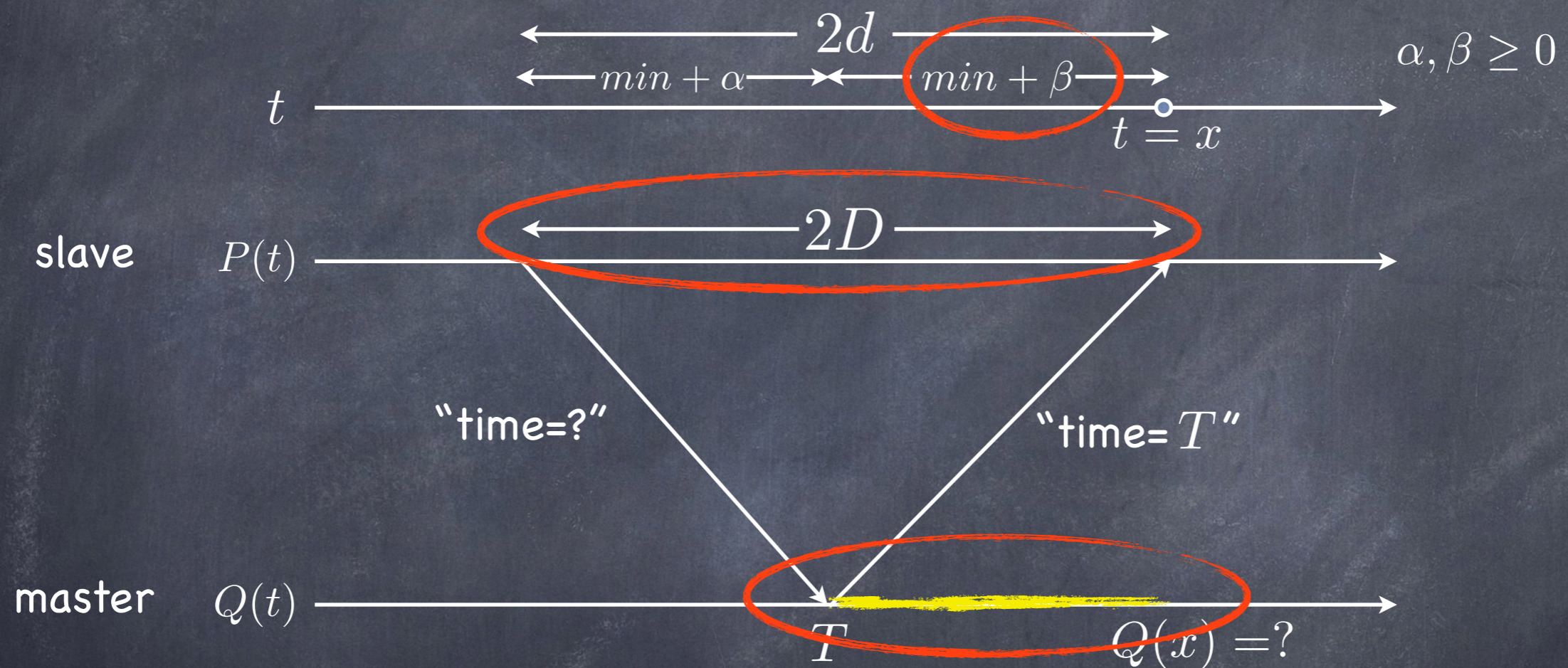$VC(e_i)[j]$ = number of events executed by process j that causally precede $e_i$

[1,0,0]

a  b  c  d

e  f  g

[0,0,1]

h  i  j

Question 2: what is the VC of:
a. event $d$
b. event $g$

# Cristian's algorithm

# 2-PHASE COMMIT

Coordinator $c$

Participant $p_i$

1. sends VOTE-REQ to all participants

2. sends $vote_i$ to Coordinator

    if $vote_i$ = **No** then
        $decision_i$ := **Abort**
    halt

3. if (all votes are **Yes**) then
        $decision_c$ := **Commit**
      send **Commit** to all
  else
        $decision_c$ := **Abort**
      send **Abort** to all who voted **Yes**
halt

4. if received **Commit** then
        $decision_i$ := **Commit**
    else
        $decision_i$ := **Abort**
    halt

# 3-PHASE COMMIT

Coordinator $c$

Participant $p_i$

1. sends VOTE-REQ to all participants

2. sends $vote_i$ to Coordinator

if $vote_i$ = **No** then
$\quad$ $decision_i$ := **Abort**
$\quad$ halt

3. if (all votes are **Yes**) then
$\quad$ send **Precommit** to all
else
$\quad$ $decision_c$ := **Abort**
$\quad$ send **Abort** to all who voted **Yes**
$\quad$ halt

4. if received **Precommit** then
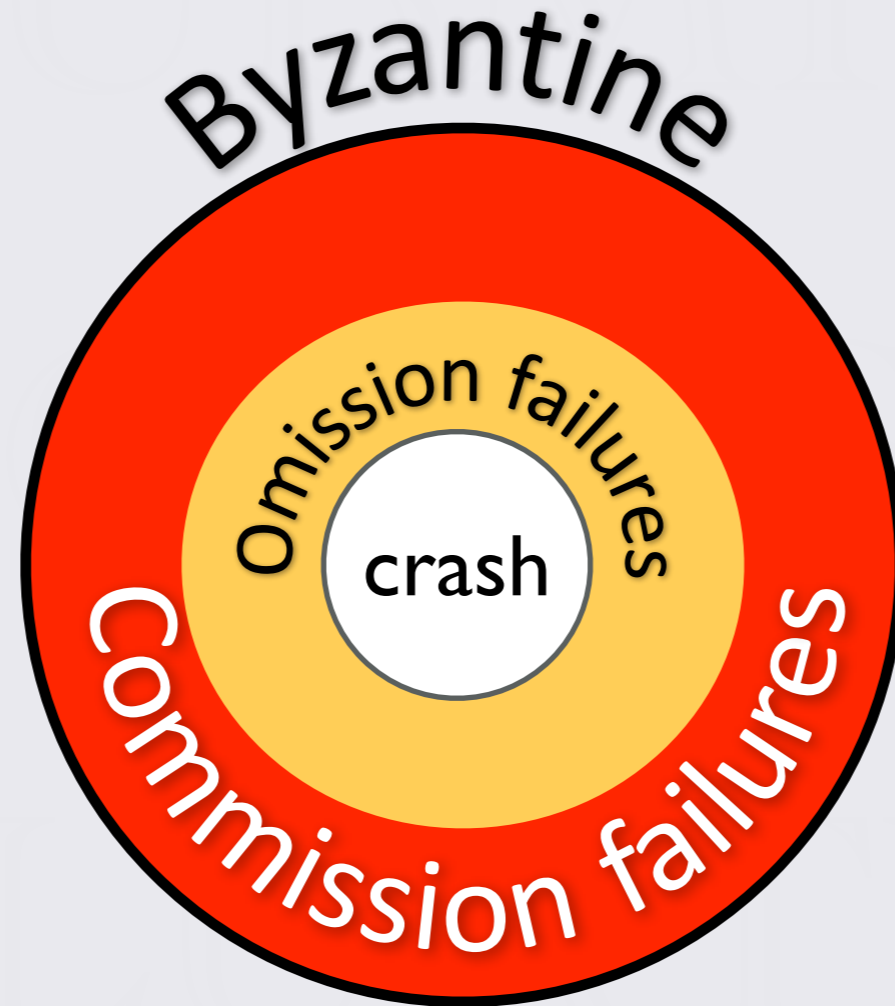$\quad$ send **Ack**

5. collect **Ack** from all participants
$\quad$ When all **Ack**'s have been received:
$\quad$ $decision_c$ := **Commit**
$\quad$ send **Commit** to all

6. When $p_i$ receives **Commit**,
sets $decision_i$ := **Commit** and halts

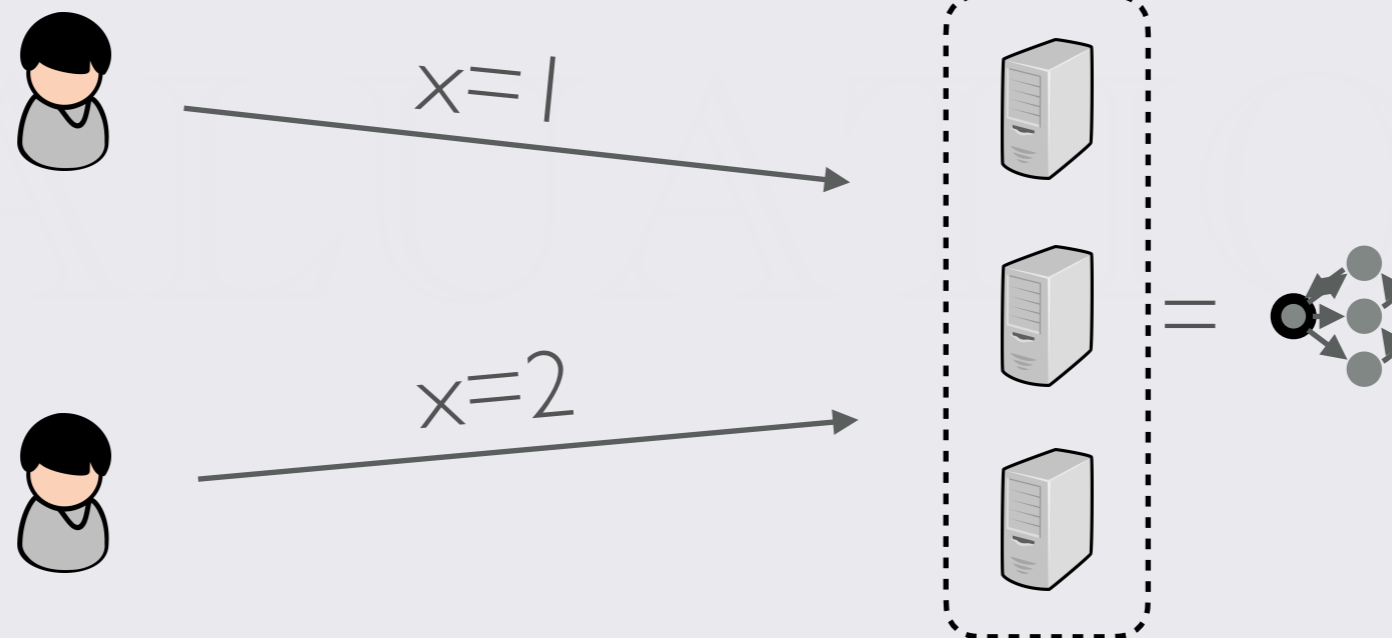# A hierarchy of failure models

# State Machine Replication



Ingredients: a server
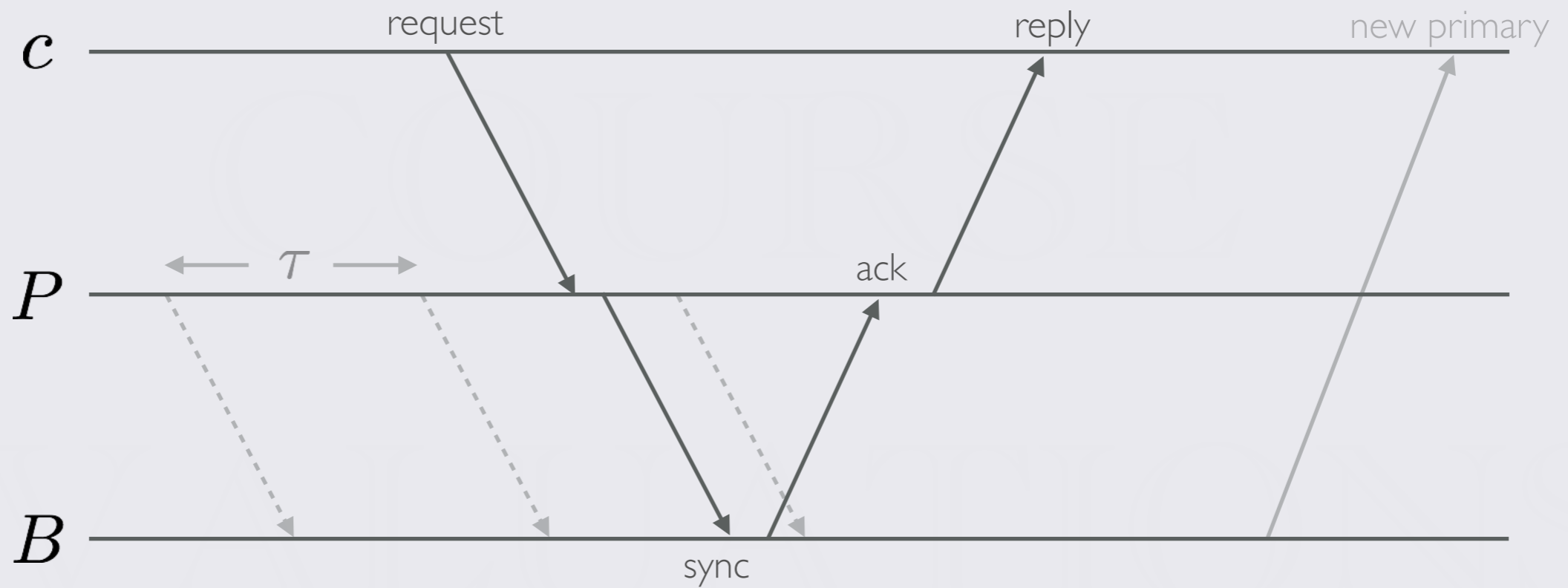
1. Make server deterministic (state machine)

2. Replicate server

3. Ensure that all replicas go through the same sequence of state transitions
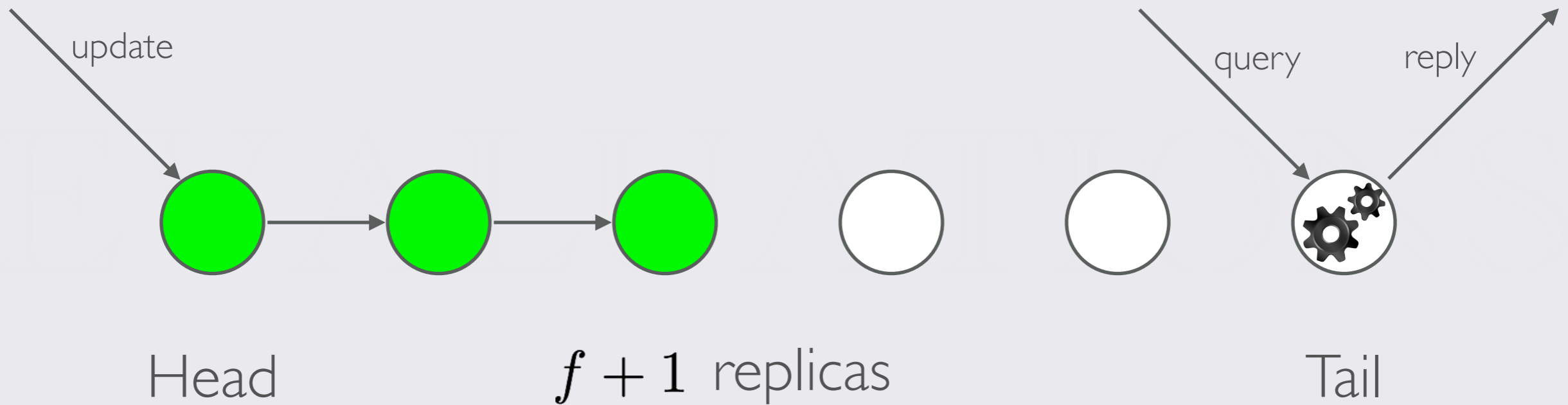
4. Vote on replica outputs

# A PRIMARY-BACKUP PROTOCOL
$$(f = 1)$$

# CHAIN REPLICATION

Tail can respond immediately, without
waiting for the new update



update

query     reply

Head        $f + 1$ replicas        Tail

# Consensus

**Validity** If all processes that propose a value propose $v$, then all correct processes eventually decide $v$

**Agreement** If a correct process decides $v$, then all correct processes eventually decide $v$

**Integrity** Every correct process decides at most one value, and if it decides $v$, then some process must have proposed $v$

**Termination** Every correct process eventually decides some value

# GOOD NEWS

Our algorithm implementing consensus in a synchronous setting is correct! That is, it is both safe and live.

# BAD NEWS

**The FLP result:**

There is no protocol that solves consensus in an asynchronous system where one process may crash
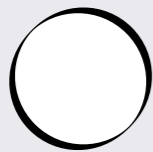
Fischer, Lynch, Paterson 1985

# Paxos

**Abstract**

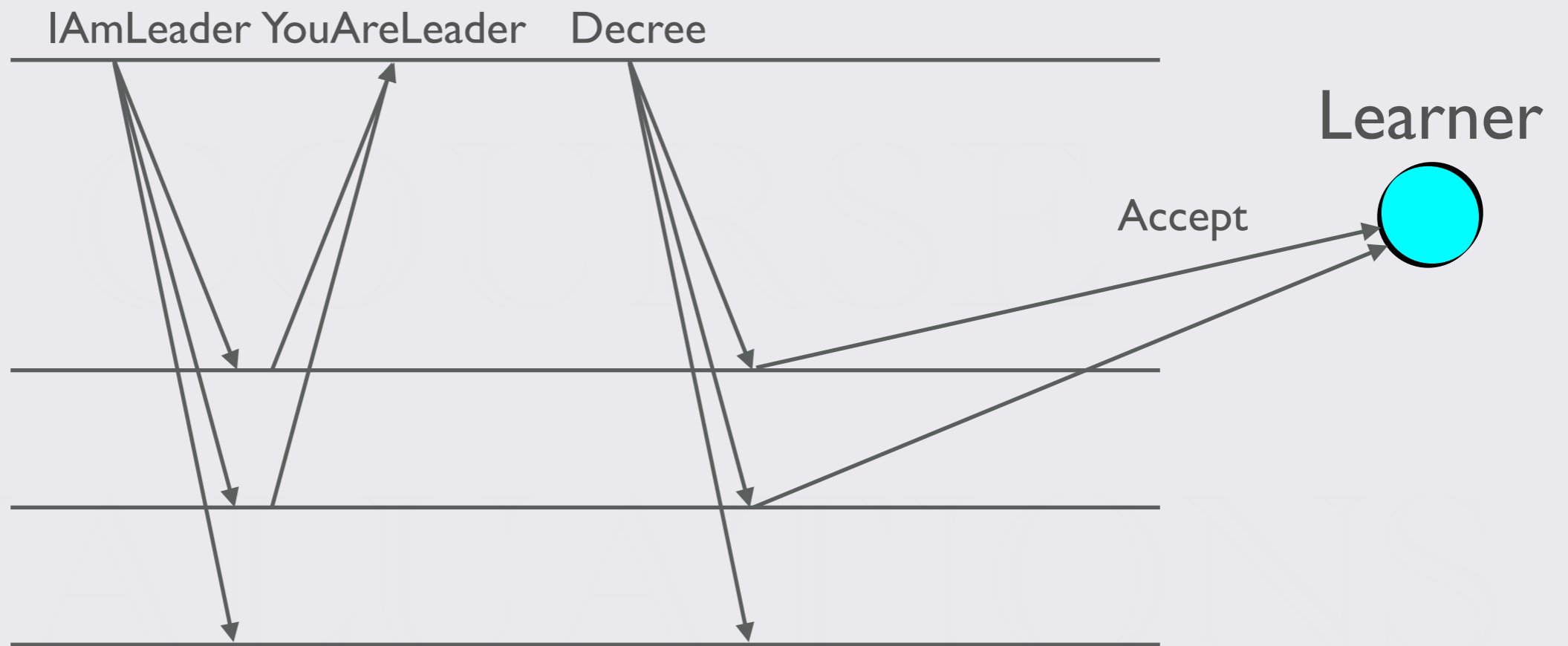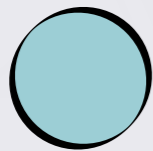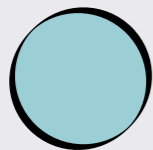The Paxos algorithm, when presented in plain English, is very simple.

# Paxos at work

Proposer

Acceptors

Learner

IAmLeader    YouAreLeader    Decree
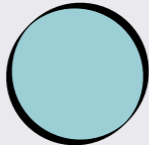
Accept

# ACCEPTOR STATES
## (as leader #50 comes to power)

| Acceptors | Value | By leader |
|:---:|:---:|:---:|
| ⬤ | x | 37 |
| ⬤ | - | - |
| ⬤ | - | - |
| ⬤ | - | - |
| ⬤ | y | 41 |

Question 4:
What is the set of possible values that leader #50 can propose?

# Examples of acceptor states
## (as leader #50 comes to power)

| Acceptors | Value | By leader |
|:---:|:---:|:---:|
| ⬤ | × | 37 |
| ⬤ | y | 42 |
| ⬤ | - | - |
| ⬤ | × | 37 |
| ⬤ | × | 41 |

Question 5:
What is the set of possible values that leader #50 can propose?

# THE THREAT TO LIVENESS: DUELING PROPOSERS

# Paxos/SMR in real life

Proposers, acceptors and learners
are all collocated on $2f + 1$ replicas

# PBFT

| Pre-prepare phase | Prepare phase | Commit phase | Reply phase |

# EXECUTE-VERIFY

Execute

Verify

First execute...

...then verify

(multithreaded and without
agreeing on the order)

(that replicas agree
on the outcome)

# THINGS I HOPE YOU WILL REMEMBER

1. Need causality? Don't reinvent vector clocks!

2. No perfect clock sync; but we can get very close.

3. Fewer than 2f+1 replicas → you **don't** tolerate asynchrony

4. Fewer than 3f+1 replicas → you **don't** tolerate non-benign faults

# 4b. Be able to tell Lorenzo apart from his evil twin



1. Evil Lorenzo Speaks French
2. And was born in Corsica
3. Went to Dartmouth instead of Cornell
4. Rides a Ducati instead of a Moto Guzzi
5. Still listens opera, but doesn't care for Puccini
5. Evil Lorenzo thinks that 2f+1 is good enough

# THINGS I HOPE YOU WILL REMEMBER (CONT.)

5. Always write to disk *before* sending a message

6. 2PC is blocking; and so is 3PC (just less frequently)

7. Be careful when messing with Paxos, or you'll get it wrong :-)

# ADMINISTRIVIA

Research project report due December 14th
(note new date)

Course evaluations due today

# Administrivia

Research project report due December 14th
(note new date)

Course evaluations due today

# HOW TO STRUCTURE A RESEARCH PAPER

- Introduction
  - **Most important** part of the paper
- Related work
- Design
- Implementation
- Evaluation
- Conclusion

# PART TWO: RESEARCH

# SYSTEMS ON REPLICATION AND FAULT TOLERANCE

**Paxos optimizations**
FastPaxos
Flexible Paxos

**Replication in the real world**
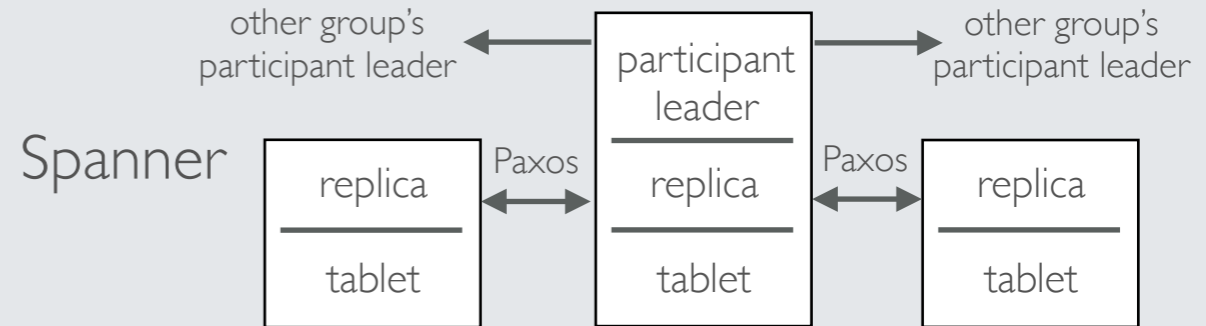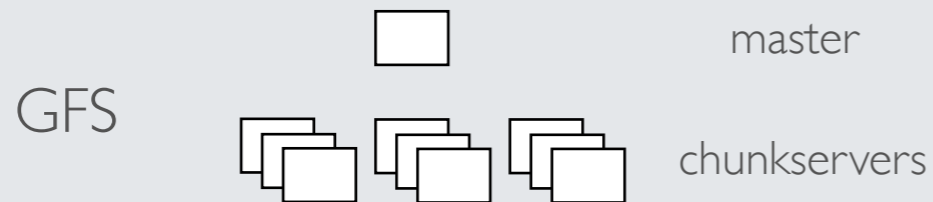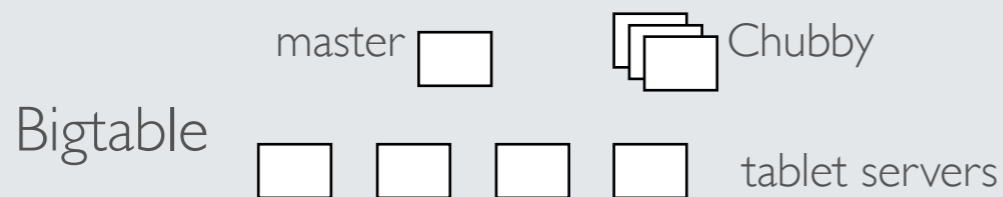ZooKeeper
CORFU

**Others**
Zyzzyva
Falcon
Mencius
IronFleet

# LARGE SCALE STORAGE SYSTEMS

**Eventual and causal consistency**
Bayou
Dynamo
COPS

**Google:**

Megastore ☐ ☐ ☐   entity group replicas
(across datacenters)

master ☐   ☐☐☐ Chubby

Bigtable ☐ ☐ ☐ ☐   tablet servers

GFS ☐   master

chunkservers

Spanner

other group's participant leader ⟷ participant leader ⟷ other group's participant leader

replica — Paxos — replica — Paxos — replica

tablet   tablet   tablet

Colossus   ????

# LARGE SCALE COMPUTATION SYSTEMS

MapReduce

input

DOG CAT RAT
CAR CAR RAT
DOG CAR CAT

split

DOG CAT RAT

CAR CAR RAT

DOG CAR CAT

map

DOG,1
CAT,1
RAT,1

CAR,1
CAR,1
RAT,1

DOG,1
CAR,1
CAT,1

shuffle

DOG,1
DOG,1

CAR,1
CAR,1
CAR,1

CAT,1
CAT,1

RAT,1
RAT,1

reduce

DOG,2

CAR,3

CAT,2

RAT,2

final
result

DOG,2
CAR,3
CAT,2
RAT,2

Spark
- No fixed graph - more expressive
- Coarse-grained transformations
- Much faster

# CRYPTOCURRENCIES

Bitcoin, Ethereum, Hyperledger Fabric, Algorand

# THINGS I HOPE YOU WILL REMEMBER

1. Consistency-performance tradeoff

2. Read papers critically

3. Read papers often

# Presentations

- Motivation, motivation, motivation!
- Keep it simple
  - Give the high-level intuition
  - Don't go too deep
- Avoid the "wall of text"
- Speak normally, with changes to your inflection
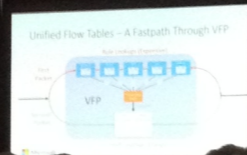- Practice, practice, practice!

# PRESENTATIONS (FINE-GRAINED)

- Your talk is a *story*, not a sequence of slides
- Look at the audience, not your laptop
- Use an outline, refer back to it frequently
  - reconnect back to your story
- Use *examples* to clarify your points
- Make sure everyone can see your text

# QUESTIONS?

- On EECS591

- On distributed systems

- On computer science

- On research

- On Life, the Universe and Everything…

# THANK YOU FOR ATTENDING EECS591!



I made it through!
Survive through the EECS591 lectures
50
Reward: Wisdom, knowledge, power!