# Cloud Computing
# #6 - Distributed computing topics & cloud

ERICSSON

# Homework #1

## The prisoner problem

There are $N$ prisoners. At random times, the guards will randomly select one prisoner to visit a room in which a 2-way switch is located. A prisoner that is visiting the room can operate the switch at their will. The warden asks if the prisoners can tell him when all prisoners have been in the room at least one time each. If they are correct he will let them free, if they are wrong they will all be executed. Prisoners have no way of communicating with each other apart from during a brief session before the process starts during which they will have to agree upon a scheme for how to solve the task. Moreover, it is not possible to detect the current status or change of the switch's state from outside the room (i.e., no lamp light is visible through a window or through the door, you cannot hear the sound of the switch changing, *et cetera*).

One prisoner is chosen to be the counter.
He keeps a count, starting at 0.

If the counter goes into the room:

If the switch is up:

Move it down.
Increment the count.
When the count reaches 99 (100 – himself),
tell the warden everyone has been
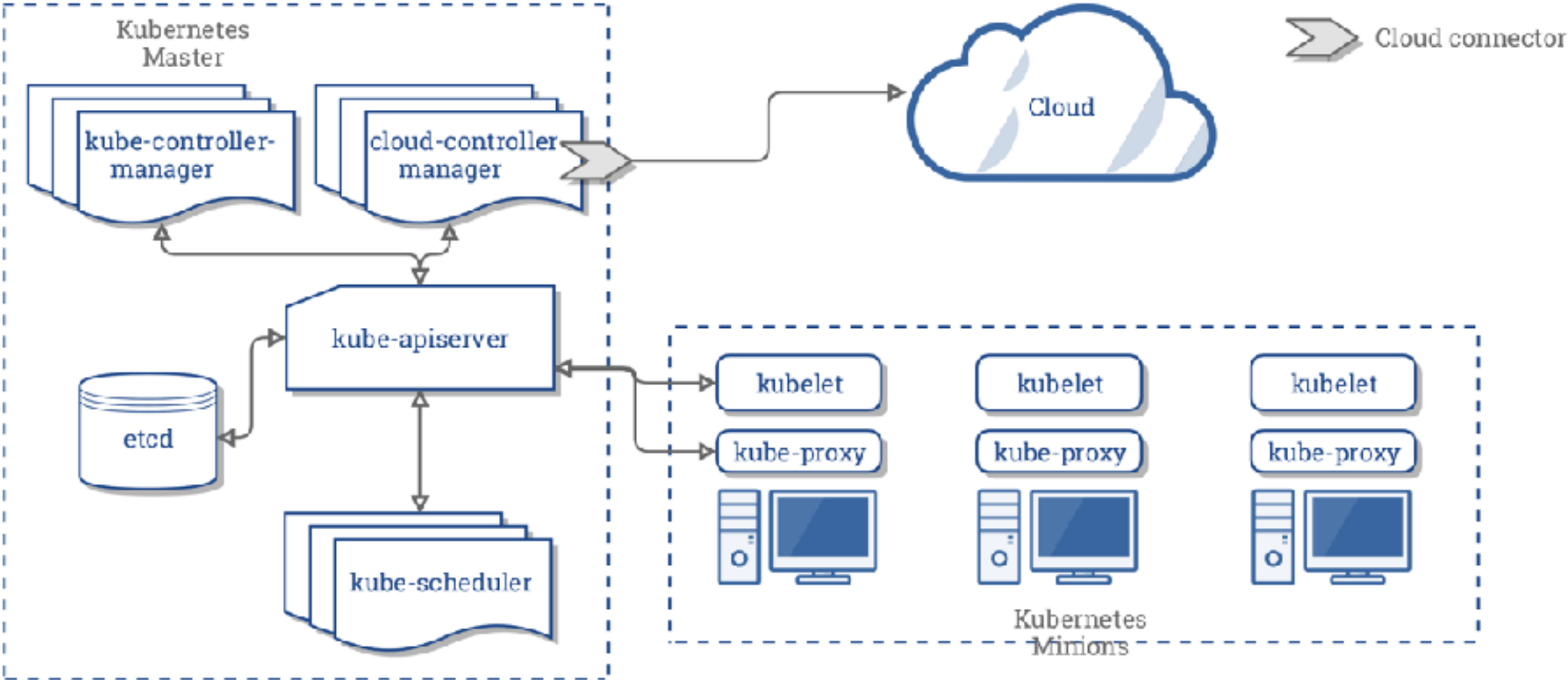in the room.

If the switch is down:
Leave it.

If any other prisoner goes into the room:

If the switch is down and he has not
previously moved it up:
Move it up.

Otherwise leave it.

# Last week: k8 architecture

# Distributed Computing

A distributed system is one in which the failure of a computer you didn't even know existed can render your own computer unusable.

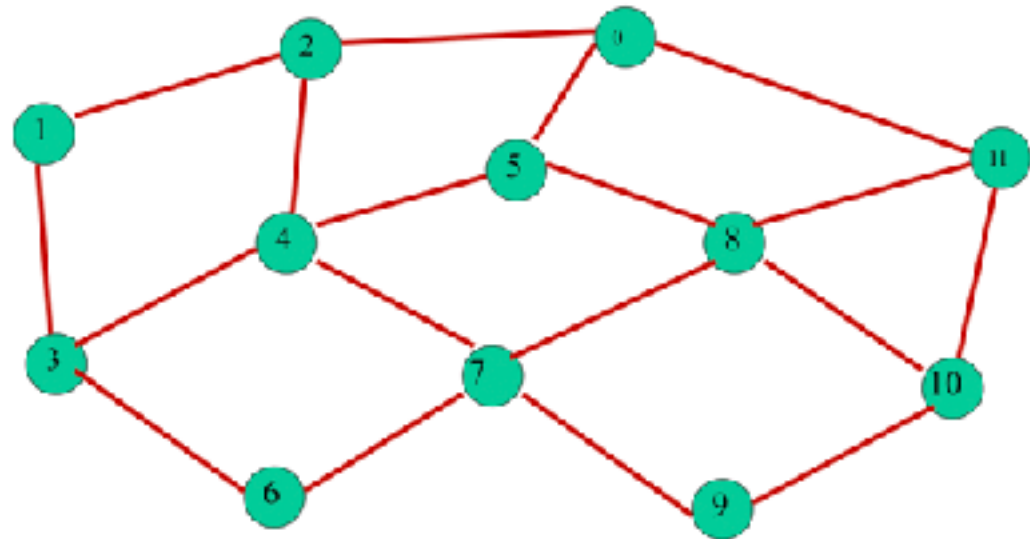Leslie Lamport
email communication
1987

Distributed computations are concurrent programs
in which processes communicate by message passing.

Gregory R. Andrews
"Paradigms for Process Interaction in Distributed Programs"
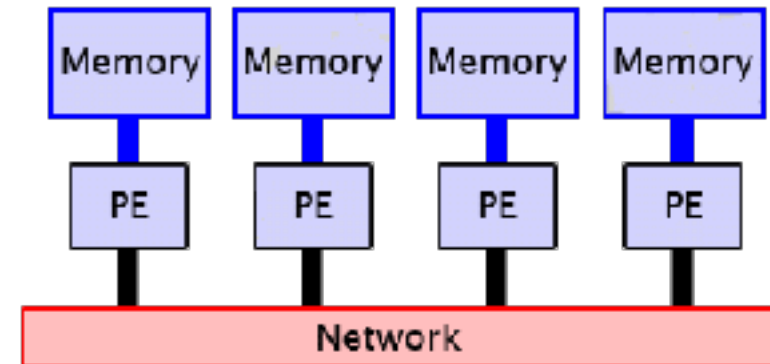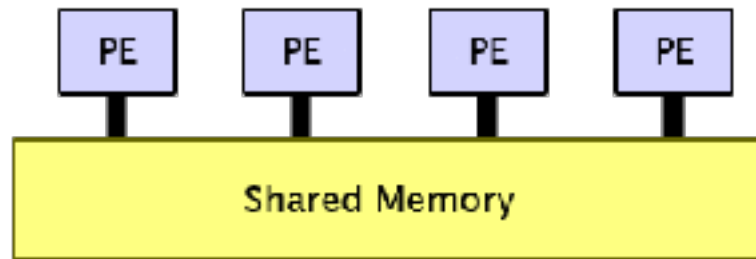ACM Computing Surveys 23(1), 1991

# A Distributed System

— Distributed system is composed of n processes

— A process executes a sequence of events

  — Local computation

  — Sending a message m

  — Receiving a message m



— A distributed algorithm is an algorithm that runs on more than one process

# Inter-Process Communication Models

# Distributed Computing

**Things being looked at**

- the algorithm(s) of the processes
- the messages
- order, causality
- whether delivery is reliable
- whether processes crash (and how)
- whether processes are "nice"

**Things that usually aren't**

- the nature of the interconnect
- time / speed
- location of processes
- data formats

# Examples of distributed algorithms

—Synchronizers (time & order)

—Resource allocation, mutual exclusion

—Data Consistency

—Failure detectors

—Consensus and agreement

——Leader election

# Synchronous vs Asynchronous

**Synchronous systems:**
known upper bounds on time for computation and message delivery
or access to global clock
or execution in synchronized rounds (not realistic, unfortunately)

**Asynchronous systems:**
no upper bounds on time for computation and message delivery

**Partially synchronous systems:**
anything in between, e.g.
- unknown upper bounds on time for computation and message delivery
- almost-synchronized clocks
- bounded-drift local clocks
- approximate bounds (on execution/message delivery time)
- bound on message delay, bound on relative process speeds
- bound on the delay ratio between fastest and slowest message at any time

# etcd

—Distributed key-value store

—**The** core component in Kubernetes
  —Single source of truth
  —Stores state of all API objects and all events that occur

—Based on the Raft consensus protocol

—Leader selection
  —only one controller-manager, scheduler, apiserver active

Why is this even hard?

# It is not only hard, it is in general impossible

Enter: the FLP theorem

## Impossibility of Distributed Consensus with One Faulty Process

MICHAEL J. FISCHER

*Yale University, New Haven, Connecticut*

NANCY A. LYNCH

*Massachusetts Institute of Technology, Cambridge, Massachusetts*

AND

MICHAEL S. PATERSON

*University of Warwick, Coventry, England*

Abstract. The consensus problem involves an asynchronous system of processes, some of which may be unreliable. The problem is for the reliable processes to agree on a binary value. In this paper, it is shown that every protocol for this problem has the possibility of nontermination, even with only one faulty process. By way of contrast, solutions are known for the synchronous case, the "Byzantine Generals" problem.

### 1. Introduction

The problem of reaching agreement among remote processes is one of the most fundamental problems in distributed computing and is at the core of many

# etcd

—Pod activation example:

—All persistent data is stored in etcd

— Distributed data base



Why is this even hard?
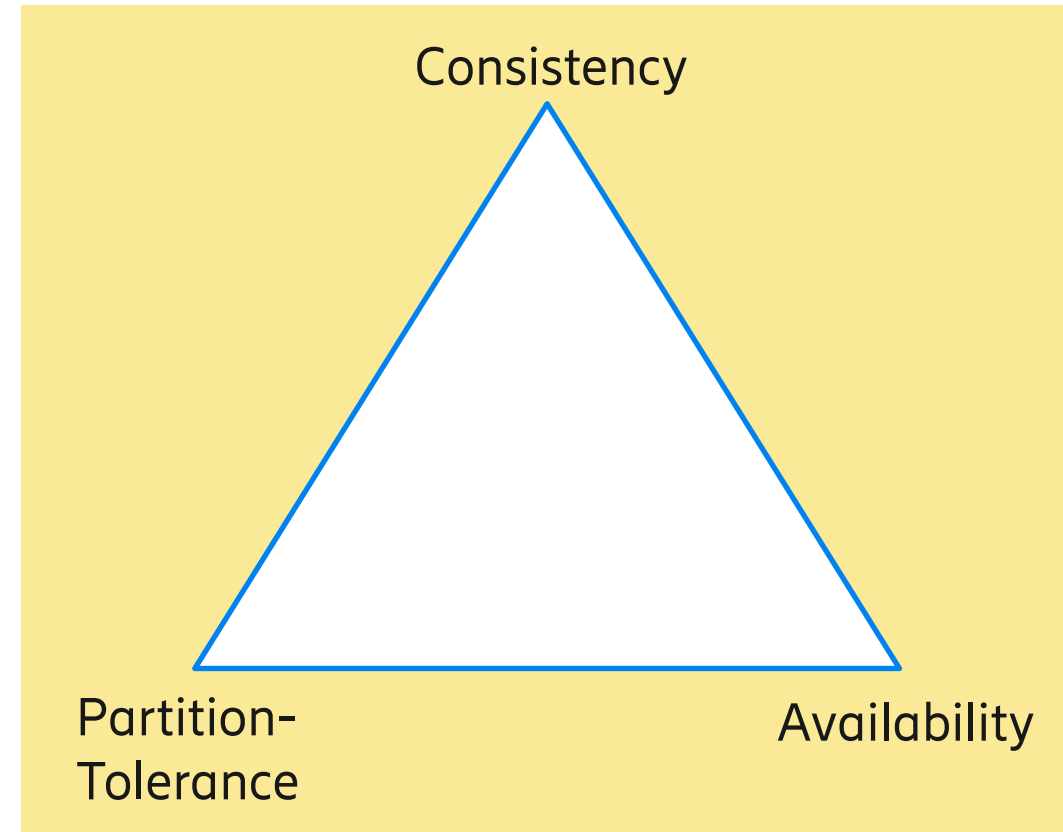
# More bad news: The CAP Theorem

**Consistency**
All clients shall see the same data at any given time.
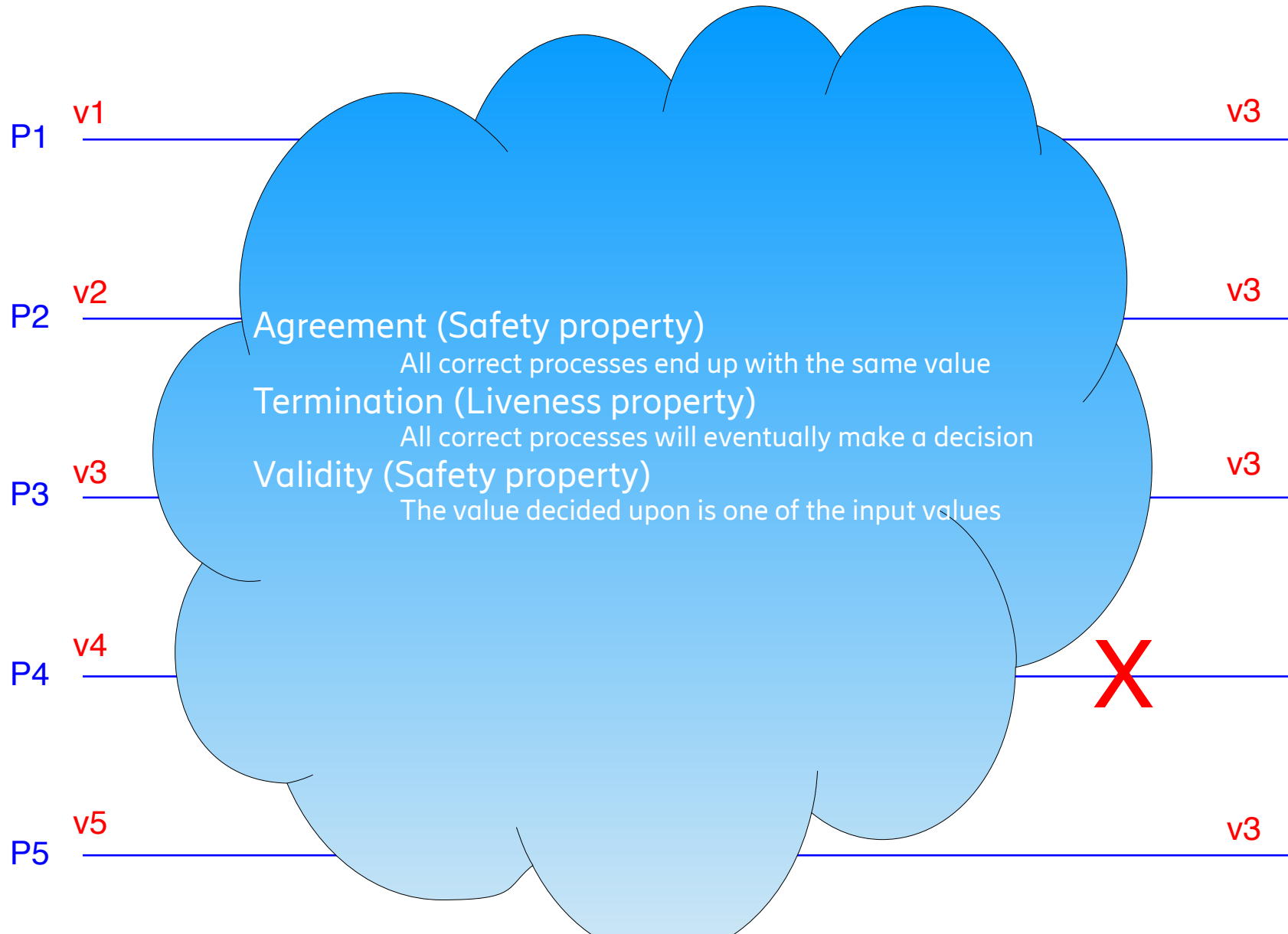A read must return the latest written value by any client.

**Availability**
The system allows read and write operations all the time, and these operations return within a reasonable time.

**Partition-Tolerance**
The system continues to work normally even if network partitions occur.



The CAP theorem was proposed by Eric Brewer from UC Berkeley, and was proved theoretically by Gilbert and Lynch from NUS and MIT.

# The Consensus Problem

P1   v1 → v3

P2   v2 → v3

P3   v3 → v3

P4   v4 → X

P5   v5 → v3

**Agreement (Safety property)**
All correct processes end up with the same value

**Termination (Liveness property)**
All correct processes will eventually make a decision

**Validity (Safety property)**
The value decided upon is one of the input values

# A Naïve Protocol

— Collect votes from all N processes
— At most one is faulty, so if one doesn't respond, count that vote as 0
— Compute majority
— Tell everyone the outcome
— They "decide" (they accept outcome)
— ... but this has a problem! Why?

# A Naïve Protocol

— In an asynchronous environment, we can't detect failures reliably

— A faulty process stops sending messages but a "slow" message might confuse us

— When the vote is nearly a tie, this confusing situation really matters

# Failure models

**Fail-Stop failures**
A faulty process halts execution prematurely.

**Byzantine failures**
*No assumption* about behavior of a faulty process.

## Some failure models (not a complete list)

— Fail-Stop (Crash-stop)

  — A processor stops, and never starts again

— Byzantine

  — A processor behaves adversarially, maliciously.

— Crash-recover

  — Well, it crashes and then restarts sometime later

— Omission

  — Doesn't respond to input (or infinitely late)

— Timing

  — Correct response, but outside required time window

# 2PC (Two-Phase Commit)



phase 1
propose & vote

phase 2
commit or abort

v

P1

propose(v)

commit()/abort()

P2

agreement?
validity?
termination?

P3

vote(Y|N)

robustness?

P4

P5

# 2PC (Two-Phase Commit)



Pi

propose()

S0

vote(Y)          vote(N)

What if proposer fails here?

S2          abort()          S1

commit()

What if both Pi & the proposer fails here?

S3

S1 - **aborted**: Voted N, received abort()

S2 - **uncertain**: Voted Y, not received commit()

S3 - **committed**: Received commit()

# 3PC (Three-Phase Commit)



*phase 1*
*propose & vote*

*phase 2*
*prepare to commit or abort*

*phase 3*
*commit*

P1

v

propose(v)

prepare()
| abort()

commit()

agreement?
validity?
termination?

P2

P3

how does robustness
differ from 2PC?

vote(Y|N)

P4

ack()

done()

P5

# 3PC (Three-Phase Commit)

Pi

propose()

S0

vote(Y)        vote(N)

S2        abort()        S1

prepare()

S3

commit()

S4

What if Pi & the proposer fails here?

S1 - **aborted**: voted N, received abort()

S2 - **uncertain**: Voted Y, not received prepare()

S3 - **committable**: Received prepare(), not commit()

S4 - **committed**: Received commit()

# 3PC (Three-Phase Commit)

Recovery rules (run by a elected node)

If some process in state **aborted**
    send abort() to all
else if some process in state **committed**
    send commit() to all
else if all processes in state **uncertain**
    send abort() to all
else if some process in state **commitable**
    send prepare() to all process in state **uncertain**
    wait for ack() and then send commit() to all

# Paxos and Chubby

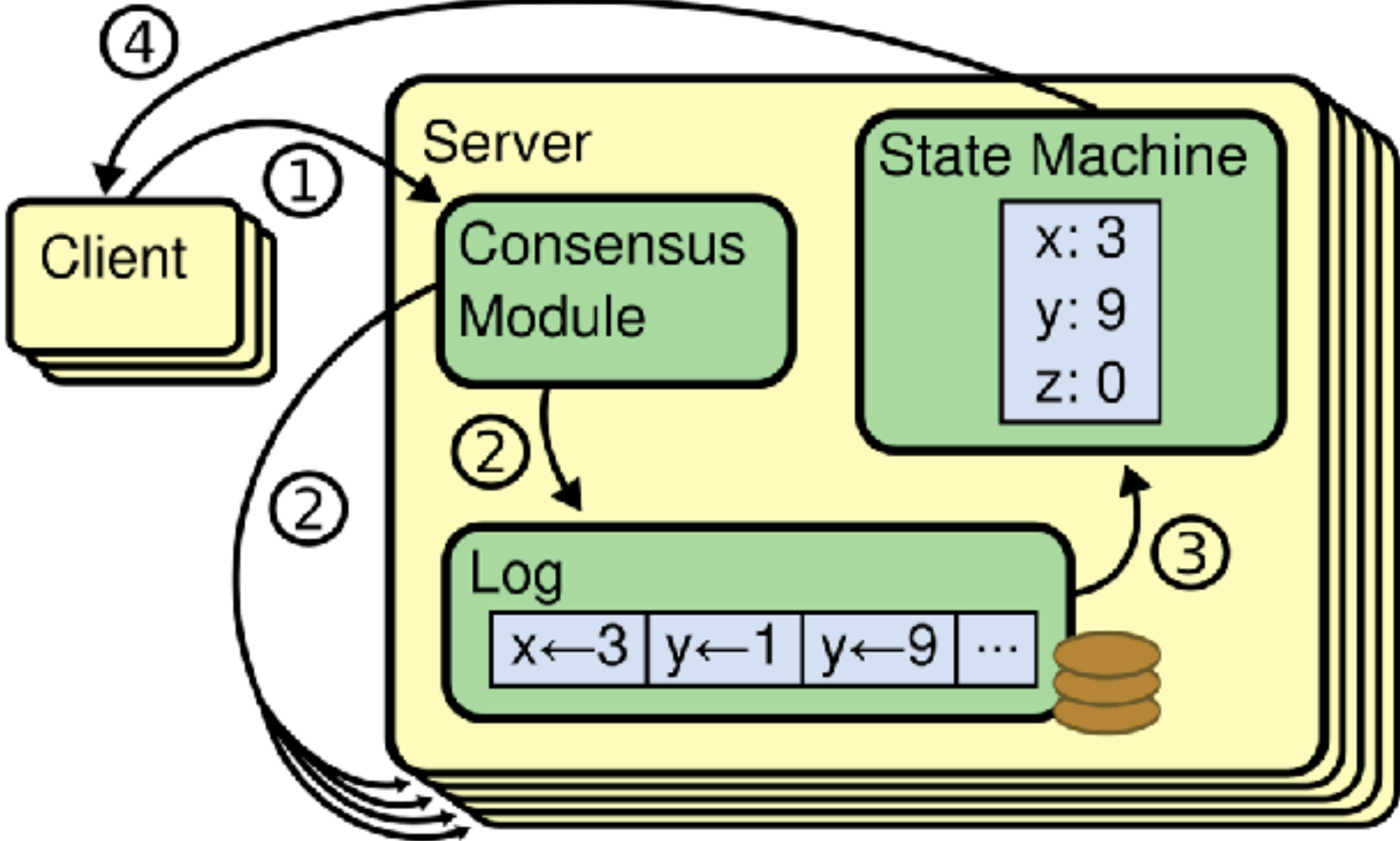P4XOS, dude

The Part-Time Parliament
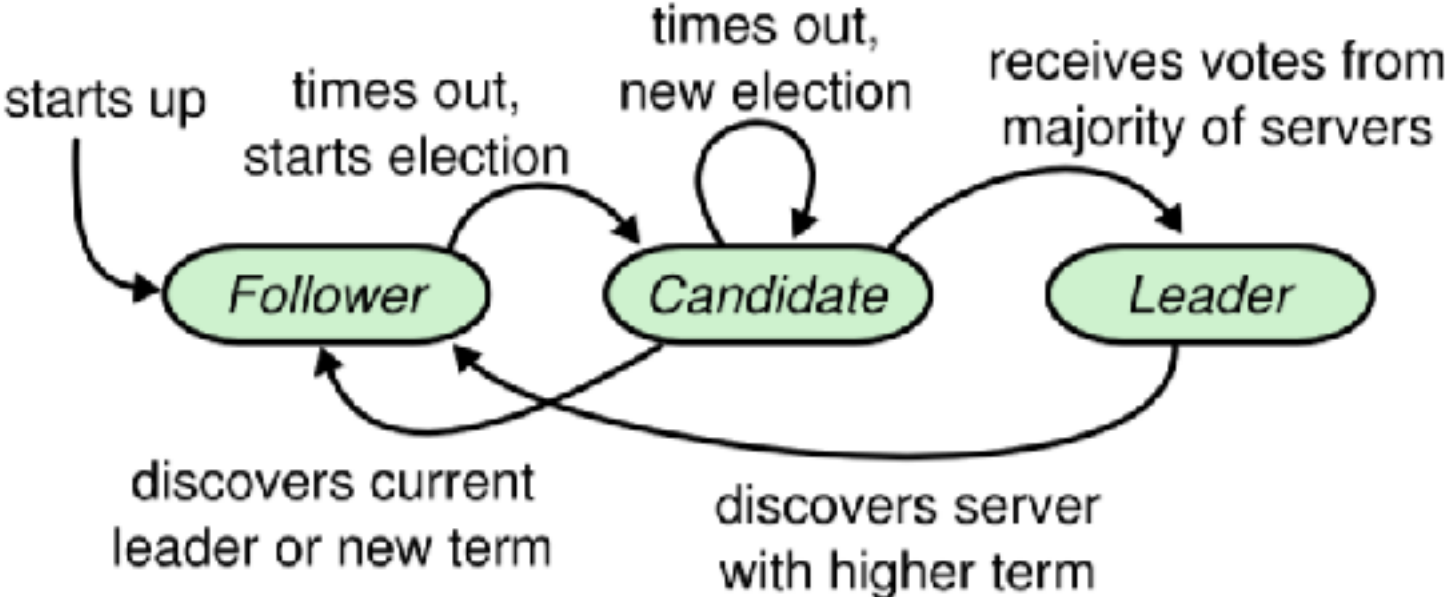
Took 8 years to review…

# Raft

— Designed to be understood

# Raft

# Consensus using etcd

```python
"""etcd3 Leader election."""
import sys
import time
from threading import Event

import etcd3


LEADER_KEY = '/leader'
LEASE_TTL = 5
SLEEP = 1


def put_not_exist(client, key, value, lease=None):
    status, _ = client.transaction(
        compare=[client.transactions.version(key) == 0],
        success=[client.transactions.put(key, value, lease)],
        failure=[],
    )
    return status


def leader_election(client, me):
    try:
        lease = client.lease(LEASE_TTL)
        status = put_not_exist(client, LEADER_KEY, me, lease)
        return status, lease
    except Exception:
        status = False
        return status, None
```

```python
def main(me):
    client = etcd3.client()
    while True:
        print('leader election')
        leader, lease = leader_election(client, me)

        if leader:
            print(me + ': leader')
            try:
                while True:
                    # do work
                    lease.refresh()
                    time.sleep(SLEEP)
            except (Exception, KeyboardInterrupt):
                return
            finally:
                lease.revoke()
        else:
            print('follower; standby')

            election_event = Event()
            def watch_cb(event):
                if isinstance(event, etcd3.events.DeleteEvent):
                    election_event.set()
            watch_id = client.add_watch_callback(LEADER_KEY, watch_cb)

            try:
                while not election_event.is_set():
                    time.sleep(SLEEP)
                print('new election')
            except (Exception, KeyboardInterrupt):
                return
            finally:
                client.cancel_watch(watch_id)

if __name__ == '__main__':
    me = sys.argv[1]
    main(me)
```
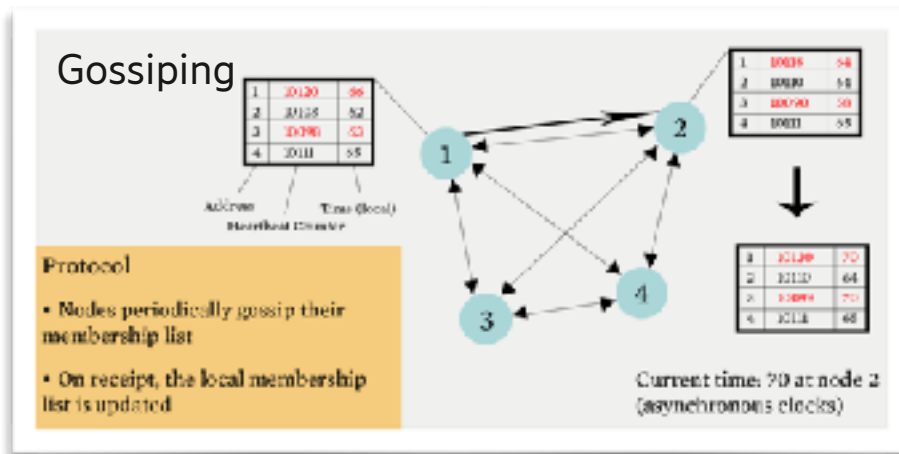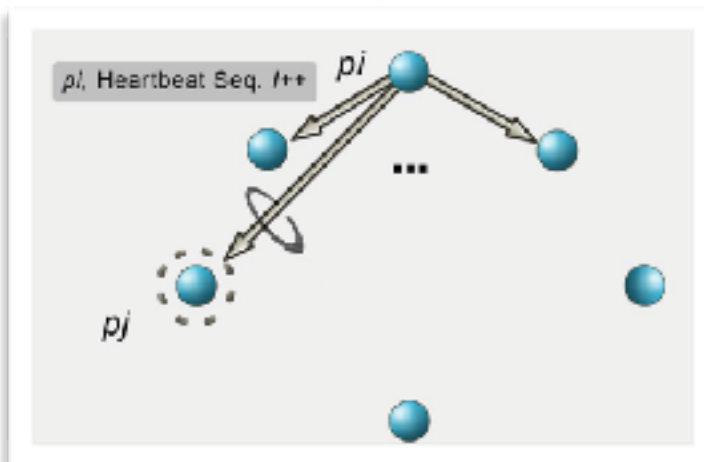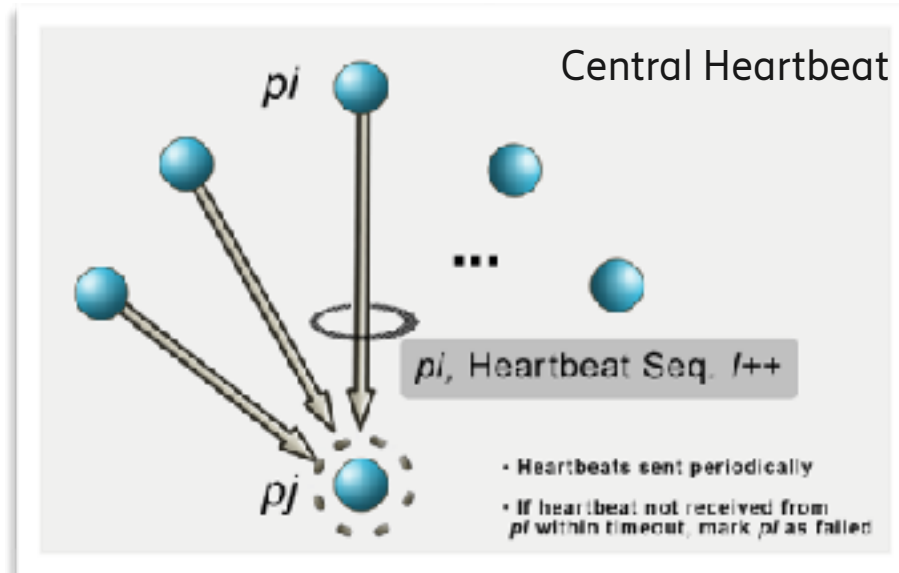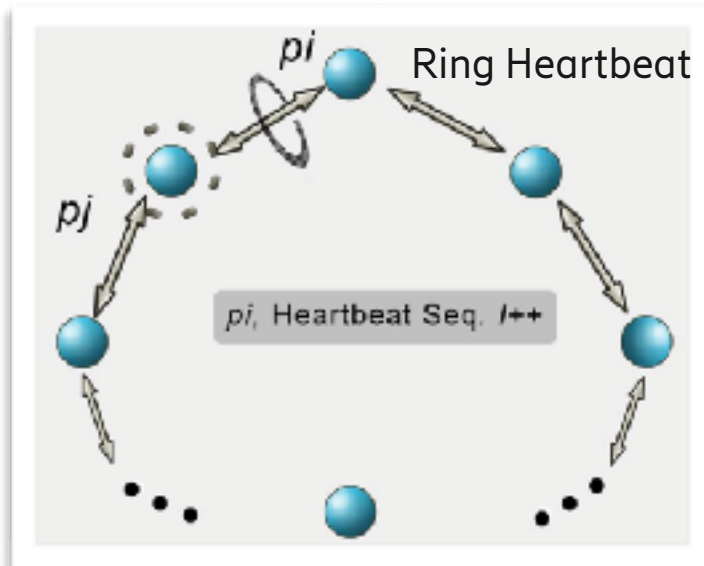
# Heartbeating
Some simple failure detectors



Ring Heartbeat

$pi$

$pj$

$pi$, Heartbeat Seq. $l$++



Central Heartbeat

$pi$

$pi$, Heartbeat Seq. $l$++

$pj$

- Heartbeats sent periodically
- If heartbeat not received from $pi$ within timeout, mark $pi$ as failed



$pi$, Heartbeat Seq. $l$++

$pi$

$pj$



Gossiping

Address        Time (local)
Heartbeat Counter

**Protocol**

- Nodes periodically gossip their membership list

- On receipt, the local membership list is updated

Current time: 70 at node 2 (asynchronous clocks)

# The Byzantine Generals

— How can we handle faulty, malicious or incomplete messages?

  — Failed broadcasts for example

  — May actively try to trick other processes, eg fake message or not sending

— Is synchronous and we can detect missing message

# The Byzantine Generals

— A city under siege by several divisions of the Byzantine army

— Each division is commanded by its own general.

— The generals communicate only using messengers.

— After observing the enemy, they must decide upon a common plan of action.

— However, some of the generals may be traitors, trying to prevent the loyal generals from reaching agreement.

— We will assume that there is a single commanding general (Commander), and the rest of the generals are his subordinates (Lieutenants)

# Objective

- All loyal generals decide upon the same plan of action
- A small number of traitors will not cause the loyal generals to adopt a bad plan

Formally rephrased as

Byzantine General Problems

1. All loyal lieutenants obey the same order
2. If the commander is loyal, then every loyal lieutenant obeys the order he sends
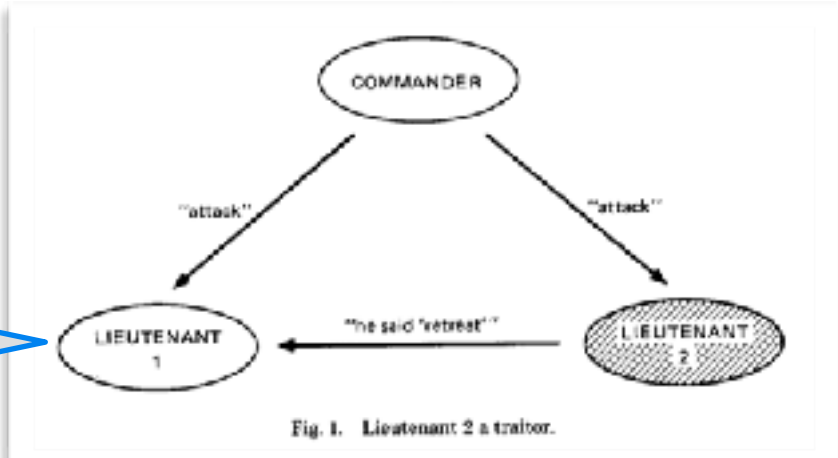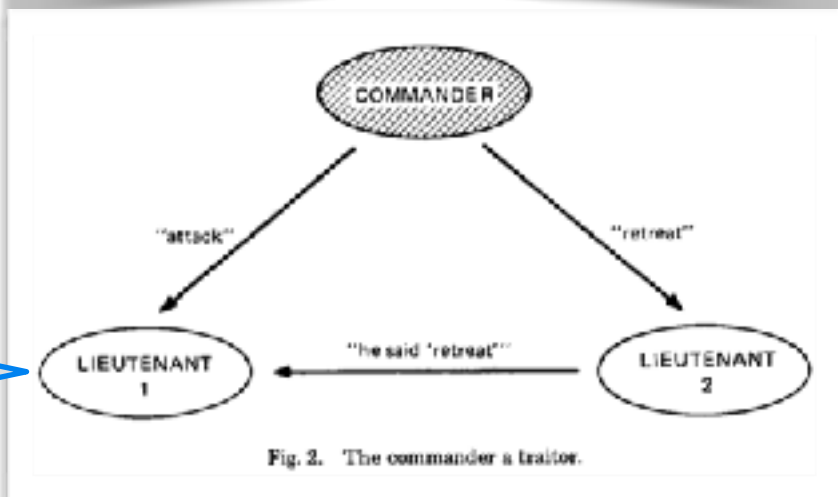
# The Byzantine Generals

1. All loyal lieutenants obey the same order
2. If the commander is loyal, then every loyal lieutenant obeys the order he sends

Desired behaviour (condition #2) is
**Attack!**

Desired behaviour (condition #1) is
**Retreat!**



Fig. 1. Lieutenant 2 a traitor.

Fig. 2. The commander a traitor.

Impossible for L1 to distinguish between the two cases! ⇒ impossible to fulfil requirements 1 & 2

# Minimal Bound on Traitors

There is no algorithm to reach consensus unless more than two thirds of the generals are loyal. In other words, impossible if $n \leq 3m$ for n processes, m of which are faulty
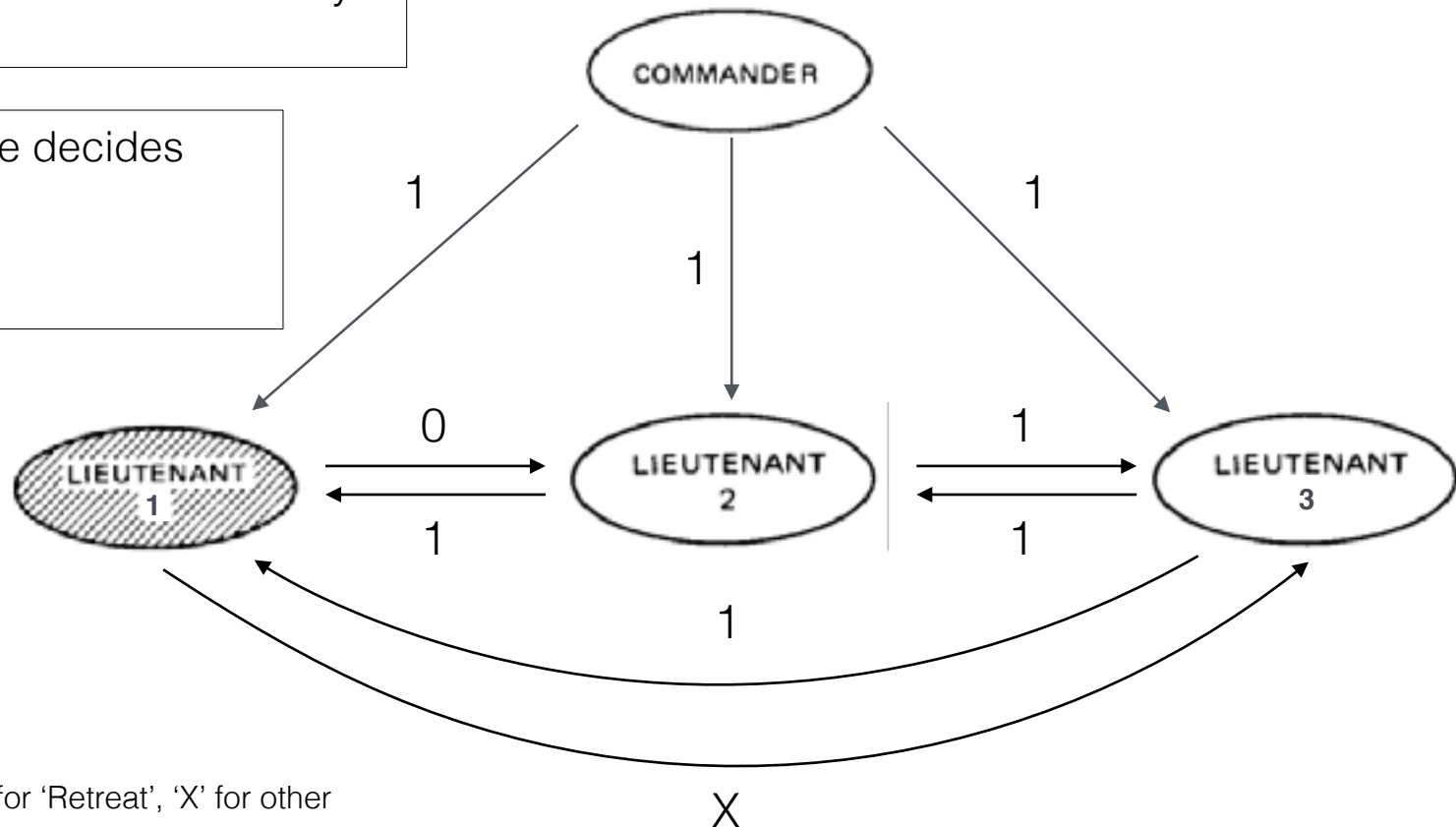
⇒ n > 3m + 1 for all algorithms that solve the

Byzantine Generals problem

# A Loyal Commander and One Traitor Lieutenant

Step 1: Commander sends same value *1* to all

Step 2: Each of L1, L2, L3 forwards the message, but L1 sends arbitrary values

Step 3: Each node decides
L2 has {1,1,0},
L3 has {1,1,X},
Both choose 1.



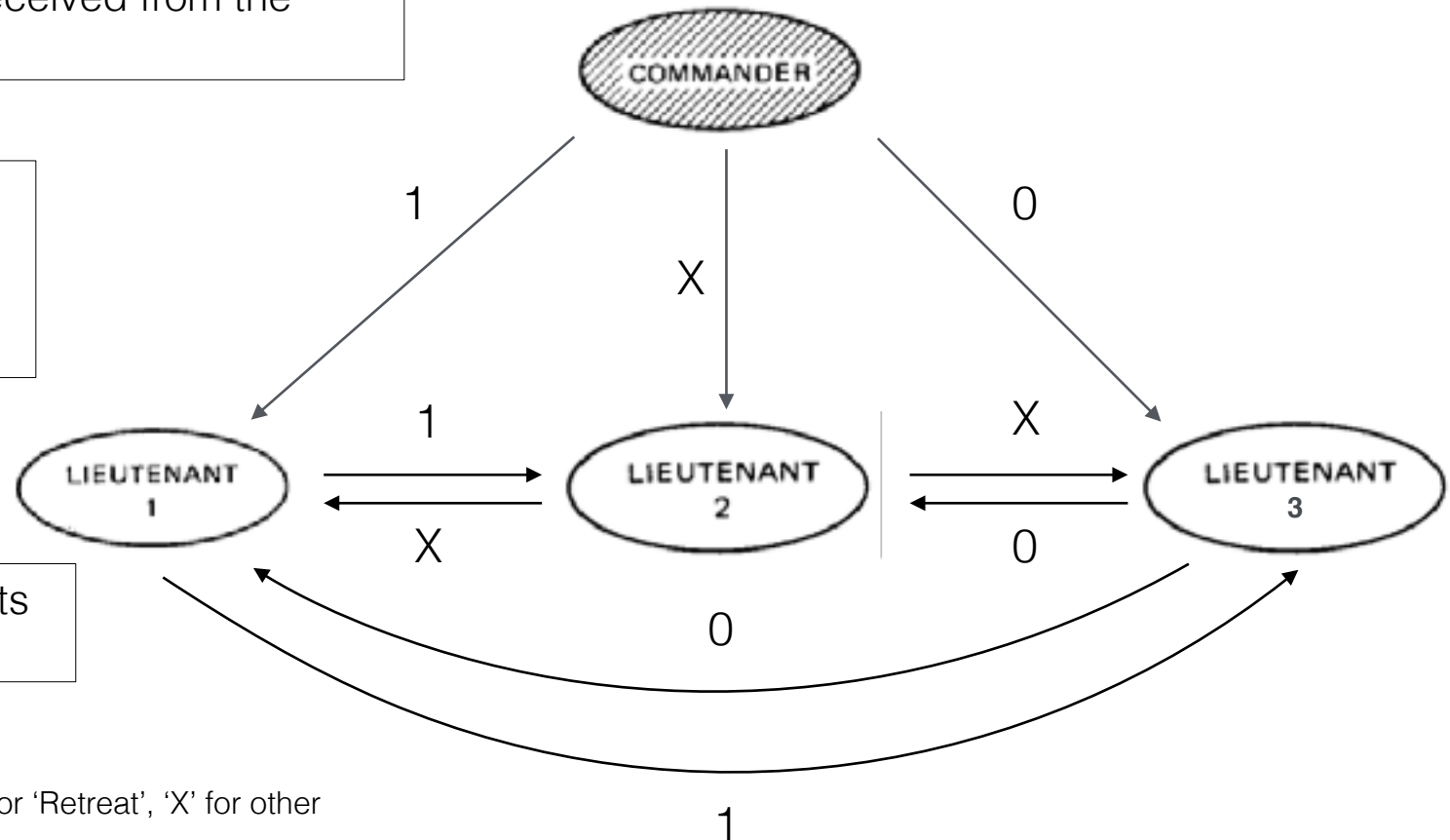Use '1' for 'Attack', or '0' for 'Retreat', 'X' for other

# A Traitor Commander with Loyal Lieutenants

Step 1: Commander sends different values to all

Step 2: Each of L1, L2, L3 forwards the values they received from the commander

Step 3: Decide
L1 has {1, 0, X},
L2 has {1, 0, X},
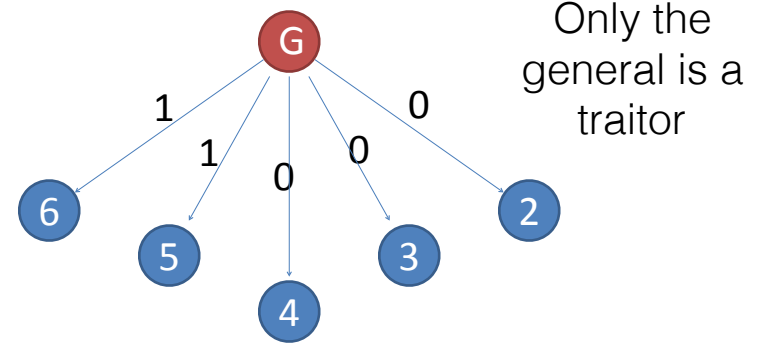L3 has {1, 0, X}

All loyal lieutenants get same result!

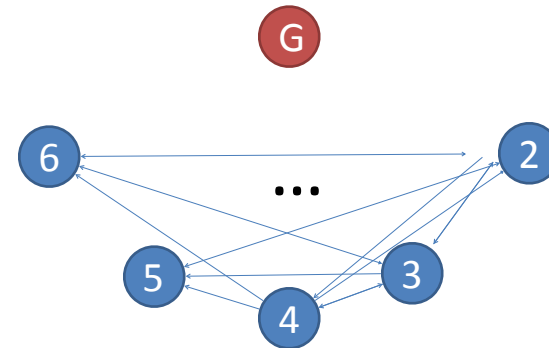Use '1' for 'Attack', or '0' for 'Retreat', 'X' for other

Let's try this live

The commander sends a message to all lieutenants

Only the general is a traitor

Use a single bit: '1' for 'Attack', or '0' for 'Retreat'.

Each lieutenant sends the message he received to all other lieutenants

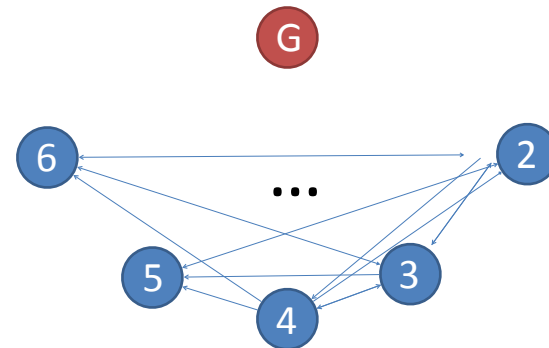| Sender=$P_2$ | Sender=$P_3$ | Sender=$P_4$ | Sender=$P_5$ | Sender=$P_6$ |
|---|---|---|---|---|
| {0,12} | {0,13} | {0,14} | {1,15} | {1,16} |

*The message {0,12} is sent to all other nodes by node 2, and so on*

Each lieutenant forwards all the messages he received to all other lieutenants

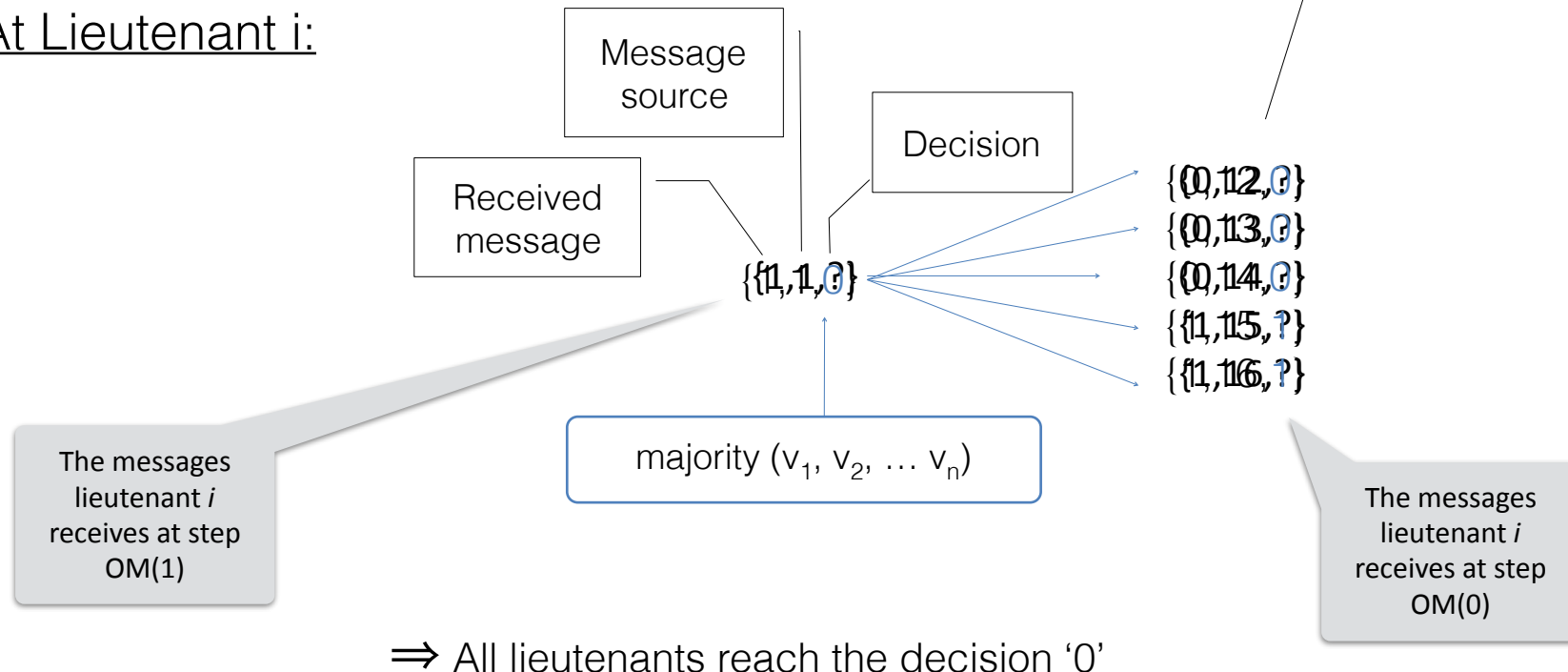| Sender=$P_2$ | Sender=$P_3$ | Sender=$P_4$ | Sender=$P_5$ | Sender=$P_6$ |
|---|---|---|---|---|
| {0,132} | {0,123} | {0,124} | {0,125} | {0,126} |
| {0,142} | {0,143} | {0,134} | {0,135} | {0,136} |
| {1,152} | {1,153} | {1,154} | {0,145} | {0,146} |
| {1,162} | {1,163} | {1,164} | {1,165} | {1,156} |

(This round is not necessary when we only have one traitor)

# The Decision Making

All received messages are kept and
organised in a tree structure

The message source is a
concatenation of all involved
node ids, i.e. the nodes
append their id when
forwarding a message

At Lieutenant i:

Message
source

Decision

Received
message

$\{0,12,0\}$
$\{0,13,0\}$
$\{0,14,0\}$
$\{1,15,?\}$
$\{1,16,?\}$

$\{1,1,0\}$

majority $(v_1, v_2, \dots v_n)$

The messages
lieutenant *i*
receives at step
OM(1)

The messages
lieutenant *i*
receives at step
OM(0)

$\Rightarrow$ All lieutenants reach the decision '0'

# Logical Clocks

— "Time, Clocks, and the Ordering of Events in a Distributed System", Leslie Lamport, 1978

— A distributed algorithm for find a total order of events in a distributed system

# The 'Happened-before' Relation



Fig. 1.

For example:

$p_1 \longrightarrow r_4$
path exists

$p_2$, $q_3$ : concurrent

$p_3$, $q_3$ : concurrent

Permits out of order message arrival

# Logical Clocks

— Logical clocks – abstract way of assigning a number to an event where the number denotes the time of occurrence of the event

— A clock $C_i\langle a \rangle$ assigns a number to an event $a$ in process $P_i$

   — Could be a number or actual time

      — The Clock Condition: $a \longrightarrow b \Longrightarrow C\langle a \rangle < C\langle b \rangle$      (the opposite is not true, why?)

# Total Ordering

— Use system of clocks satisfying the Clock Condition to place a total ordering of all events denoted "$\Rightarrow$"

    — Simply order events by the their time $C\langle a\rangle$ and break ties with any arbitrary total order (alphabetical, etc)

— We can define a total ordering on the set of all system events
    — $a \Rightarrow b$ if either $C_i\langle a\rangle < C_j\langle b\rangle$ or $C_i\langle a\rangle = C_j\langle b\rangle$ and $P_i < P_j$

— This ordering is not unique, but well defined

# Data Replication

— There are two primary reasons for replicating data

   — Reliability

   — Performance.

— Three classic approaches to replicated data

   — primary copy

   — multi-master

   — quorum consensus



W1. Write request
W2. Forward request to primary
W3. Tell backups to update
W4. Acknowledge update
W5. Acknowledge write completed

R1. Read request
R2. Response to read

# Example: Facebook

- Read from closest server
- Write to California
- Other servers update cache every 15 minutes
- After write: read from CA for 15 minutes



Slide: Marc Shapiro

# The Consistency Problem

— How is data replicated across a distributed system?

— The client side view :
- — A set of process that reads & writes
- — A distributes data store that is treated as a black box
- — How & when do updates become observable?

— The server side view:
- — The inside of the distributed data store
- — How is data propagated and replicated between storage nodes?



Distributed data store

# Strict consistency

- All writes are instantaneously visible to all processes and absolute global time order is maintained

- Hard to achieve!

| | | |
|---|---|---|
| P1: | W(x)a | |
| P2: | | R(x)a |
| | (a) | |

correct

| | | | |
|---|---|---|---|
| P1: | W(x)a | | |
| P2: | | R(x)NIL | R(x)a |
| | (b) | | |

incorrect

# Sequential consistency

— Operations on each process must be in order
— All processes see the same interleaving set of operations, regardless of what that interleaving is.
— Cares about **program order,** not time

| | | | |
|---|---|---|---|
| P1: W(x)a | | | |
| P2: | W(x)b | | |
| P3: | | R(x)b | R(x)a |
| P4: | | | R(x)b R(x)a |

(a)

correct

| | | | |
|---|---|---|---|
| P1: W(x)a | | | |
| P2: | W(x)b | | |
| P3: | | R(x)b | R(x)a |
| P4: | | | R(x)a R(x)b |

(b)

incorrect

# Causal consistency

Writes that are *potentially* causally related must be seen by all processes in the same order. Concurrent writes may be seen in a different order by different processes.



incorrect



correct

# Eventual Consistency

— The only requirement is that all replicas will eventually be the same.

— DNS is a well know example. Updates to a name are distributed according to a configured pattern and in combination with time-controlled caches; eventually, all clients will see the update.

# Server Side Consistency

The algorithms that implement the consistency models of choice

— N = the number of replicas

— W = the write set

— R = the read set



Process        Process        Process

Local copy

Distributed data store

# Server Configurations

— W+R > N ⇒ guarantee strong consistency.

— N=2, W=2, and R=1 (primary-copy scenario with synchronous replication).
  — No matter from which replica the client reads, it will always get a consistent answer.

— N=2, W=1, and R=1 (primary-copy scenario with asynchronous replication).
  — In this case R+W=N, and consistency cannot be guaranteed.

— With N=3 and W=3 and only two nodes available, the system will fail to write.

# AWS S3

```
PUT /key-prefix/cool-file.jpg 200
GET /key-prefix/cool-file.jpg 200
```

```
PUT /key-prefix/cool-file.jpg 200
PUT /key-prefix/cool-file.jpg 200 (new content)
GET /key-prefix/cool-file.jpg 200
```

```
GET /key-prefix/cool-file.jpg 404
PUT /key-prefix/cool-file.jpg 200
GET /key-prefix/cool-file.jpg 404
```

WTF!?

## Amazon S3 Data Consistency Model

Amazon S3 provides read-after-write consistency for PUTS of new objects in your S3 bucket in all Regions with one caveat. The caveat is that if you make a HEAD or GET request to the key name (to find if the object exists) before creating the object, Amazon S3 provides eventual consistency for read-after-write.

Amazon S3 offers eventual consistency for overwrite PUTS and DELETES in all Regions.

Updates to a single key are atomic. For example, if you PUT to an existing key, a subsequent read might return the old data or the updated data, but it never returns corrupted or partial data.

Consistency

Availability

CA

CP

AP

Partition
Tolerance

*cassandra*

```python
from cassandra import ConsistencyLevel
from cassandra.query import SimpleStatement

query = SimpleStatement(
    "INSERT INTO users (name, age) VALUES (%s, %s)",
    consistency_level=ConsistencyLevel.QUORUM)
session.execute(query, ('John', 42))
```

# Summary

— Overview of the foundations for cloud

— Distributed systems play an crucial role in getting things to work

— Take Jorn Janneck's course on distributed algorithms and learn the details.

```
Distributed Systems (7.5hp)

To give an introduction to the fundamental concepts of distributed
systems, their properties and application in practice.

display basic knowledge of:
different types of distributed systems and their properties,
failure and recovery in distributed systems,
models and abstractions for distributed systems,
distributed models of logical time,
distributed algorithms and protocols,
and distributed state and computing.

be able to reason about properties of distributed systems,
be able to use concepts and abstractions to model distributed systems
and to express their behavior,
be able to use fundamental distributed algorithms and protocols in
managing resources, sharing state, maintaining distributed state, and
coordinate distributed computation,
be able to apply the conceptual knowledge to the implementation of
distributed algorithms on a variety of platforms.

be able to judge the suitability of models and platforms for distributed
systems for a given problem,
display a basic understanding of the tradeoffs and limits of the
concepts and techniques in distributed system design.
```