

# **Fault Tolerance in Distributed Systems**



**ICS 230**

**Prof. Nalini Venkatasubramanian**

**(with some slides modified from Prof.  
Ghosh, University of Iowa and  
Indranil Gupta, UIUC)**

# Fundamentals



- **What is fault?**
  - A fault is a blemish, weakness, or shortcoming of a particular hardware or software component.
  - Fault, error and failures
- **Why fault tolerant?**
  - Availability, reliability, dependability, ...
- **How to provide fault tolerance ?**
  - Replication
  - Checkpointing and message logging
  - Hybrid

# Reliability

- Reliability is an emerging and critical concern in traditional and new settings
  - Transaction processing, mobile applications, cyberphysical systems
- New enhanced technology makes devices vulnerable to errors due to high complexity and high integration
  - Technology scaling causes problems
    - Exponential increase of soft error rate
  - Mobile/pervasive applications running close to humans
    - E.g Failure of healthcare devices cause serious results
  - Redundancy techniques incur high overheads of power and performance
    - TMR (Triple Modular Redundancy) may exceed 200% overheads without optimization [Nieuwland, 06]
- Challenging to optimize multiple properties (e.g., performance, QoS, and reliability)

# Classification of failures



Crash failure

Security failure

Omission failure

Temporal failure

Transient failure

Byzantine failure

Software failure

Environmental perturbations

# Crash failures



Crash failure = the process halts. It is *irreversible*.

In synchronous system, it is easy to detect crash failure (using *heartbeat signals* and *timeout*). But in asynchronous systems, it is never accurate, since it is *not possible* to distinguish between a process that has crashed, and a process that is running *very slowly*.

Some failures may be complex and nasty. **Fail-stop failure** is a *simple abstraction* that *mimics* crash failure when program execution becomes arbitrary. Implementations help detect which processor has failed. If a system cannot tolerate fail-stop failure, then it cannot tolerate crash.

# Transient failure



**(Hardware)** Arbitrary perturbation of the global state. May be induced by power surge, weak batteries, lightning, radio-frequency interferences, cosmic rays etc.

**(Software)** Heisenbugs are a class of temporary internal faults and are intermittent. They are essentially permanent faults whose conditions of activation occur rarely or are not easily reproducible, so they are harder to detect during the testing phase.

Over 99% of bugs in IBM DB2 production code are non-deterministic and transient (Jim Gray)

# Temporal failures



**Inability to meet deadlines** – correct results are generated, but too late to be useful. Very important in real-time systems.

May be caused by poor algorithms, poor design strategy or loss of synchronization among the processor clocks

# Byzantine failure

Anything goes! Includes every conceivable form of erroneous behavior. The weakest type of failure

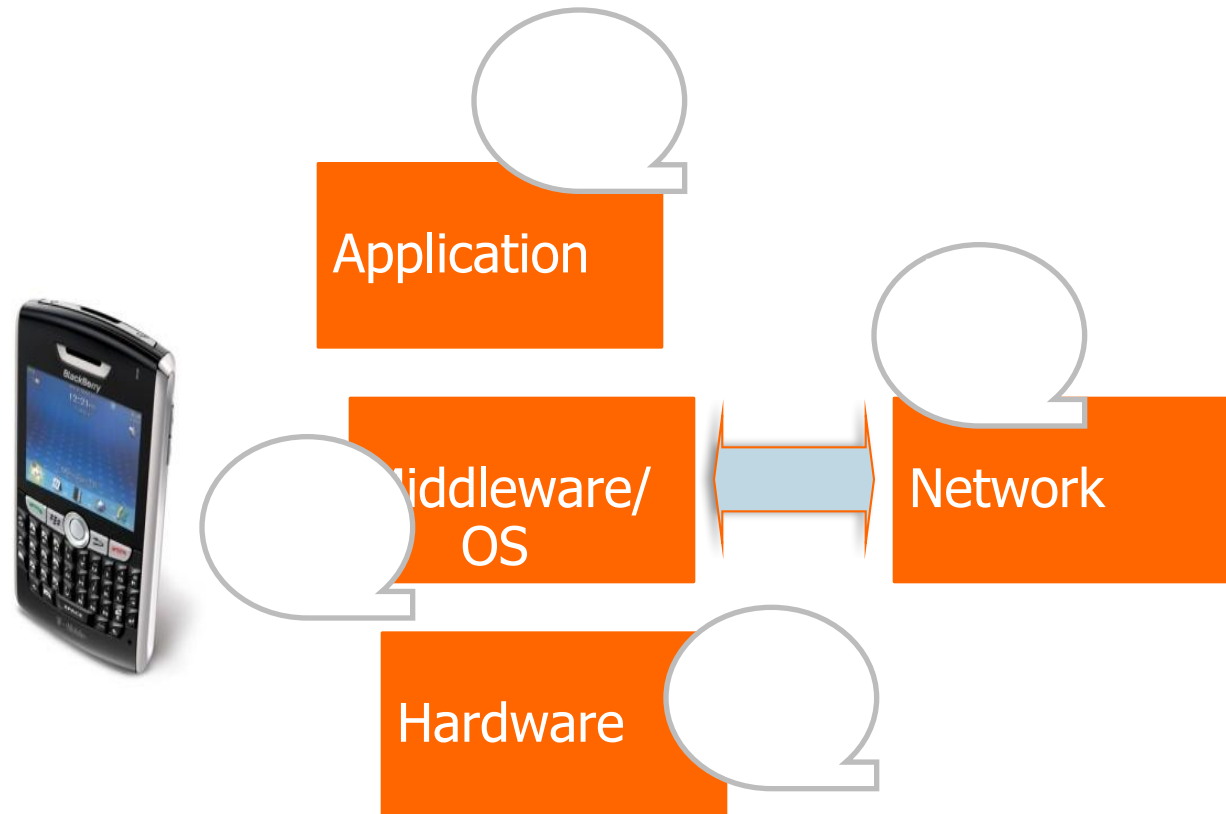
Numerous possible causes. Includes **malicious behaviors** (like a process executing a different program instead of the specified one) **too**.

Most difficult kind of failure to deal with.



# Errors/Failures across system layers

- Faults or Errors can cause Failures



# Hardware Errors and Error Control Schemes

Failures	Causes	Metric s	Traditional Approaches
Soft Errors, Hard Failures, System Crash	External Radiations, Thermal Effects, Power Loss, Poor Design, Aging	FIT, MTTF, MTBF	Spatial Redundancy (TMR, Duplex, RAID-1 etc.) and Data Redundancy (EDC, ECC, RAID-5, etc.)

- Hardware failures are increasing as technology scales

- (e.g.) SER increases by up to 1000 times [Mastipuram, 04]

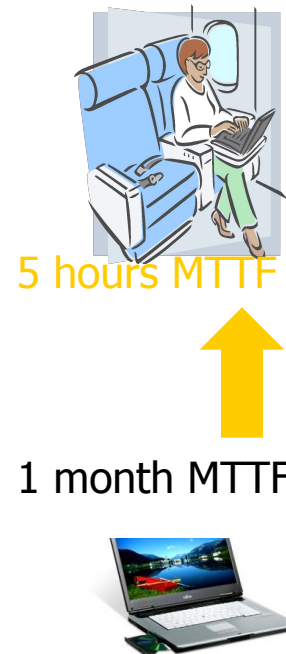
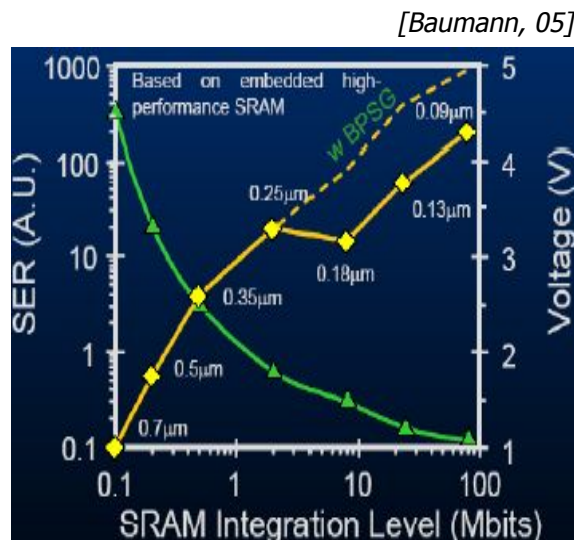
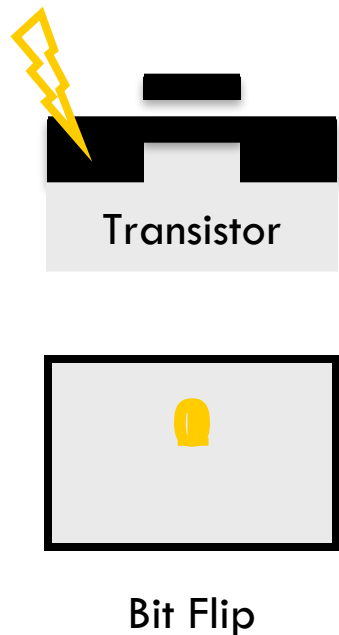
- Redundancy techniques are expensive

- (e.g.) ECC-based protection in caches can incur 95% performance penalty [Li, 05]

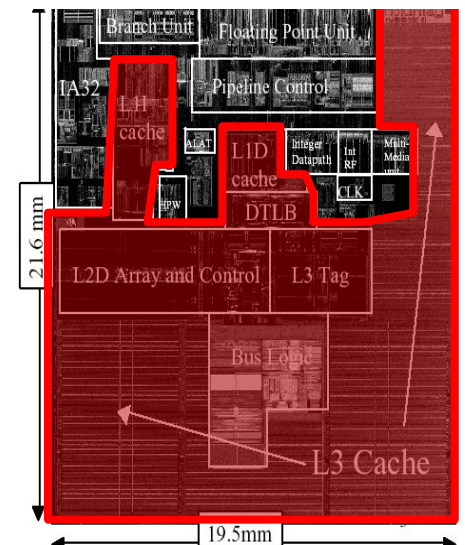
- FIT: Failures in Time ( $10^9$  hours)
- MTTF: Mean Time To Failure
- MTBF: Mean Time b/w Failures
- TMR: Triple Modular Redundancy
- EDC: Error Detection Codes
- ECC: Error Correction Codes
- RAID: Redundant Array of Inexpensive Drives

# Soft Errors (Transient Faults)

- SER increases exponentially as technology scales
- Integration, voltage scaling, altitude, latitude
- Caches are most hit due to:
  - Larger portion in processors (more than 50%)
  - No masking effects (e.g., logical masking)



Intel Itanium II Processor



•MTTF: Mean time To Failure

# Soft errors

	SER (FIT)	MTTF	Reason
1 Mbit @ 0.13 $\mu$ m	1000	<b>104 years</b>	
<b>64 MB</b> @ 0.13 $\mu$ m	<b>64x8x1000</b>	<b>81 days</b>	High Integration
<b>128 MB</b> @ <b>65 nm</b>	<b>2x1000x64x8x1000</b>	<b>1 hour</b>	Technology scaling and Twice Integration
<b>A system</b> @ 65 nm	<b>2x2x1000x64x8x1000</b>	<b>30 minutes</b>	Memory takes up 50% of soft errors in a system
A system <b>with voltage scaling</b> @ 65 nm	<b>100x2x2x1000x64x8x1000</b>	<b>18 seconds</b>	Exponential relationship b/w SER & Supply Voltage
A system with voltage scaling @ <b>flight (35,000 ft)</b> @ 65 nm	<b>800x100x2x2x1000x64x8x1000</b> FIT	<b>0.02 seconds</b>	High Intensity of Neutron Flux at flight (high altitude)
Soft Error Rate (SER) – FIT (Failures in Time) = number of errors in $10^9$ hours			

# Software Errors and Error Control Schemes

Failures	Causes	Metrics	Traditional Approaches
Wrong outputs, Infinite loops, Crash	Incomplete Specification, Poor software design, Bugs, Unhandled Exception	Number of Bugs/Klines, QoS, MTTF, MTBF	Spatial Redundancy (N-version Programming, etc.), Temporal Redundancy (Checkpoints and Backward Recovery, etc.)

- ❑ Software errors become dominant as system's complexity increases
  - ❑ (e.g.) Several bugs per kilo lines
- ❑ Hard to debug, and redundancy techniques are expensive
  - ❑ (e.g.) Backward recovery with checkpoints is inappropriate for real-time applications

# Software failures



## Coding error or human error

On September 23, 1999, NASA lost the \$125 million Mars orbiter spacecraft because **one engineering team used metric units** while **another used English units** leading to a navigation fiasco, causing it to burn in the atmosphere.

## Design flaws or inaccurate modeling

Mars pathfinder mission landed flawlessly on the Martial surface on July 4, 1997. However, later its communication failed due to a design flaw in the real-time embedded software kernel VxWorks. The problem was later diagnosed to be caused due to **priority inversion**, when a medium priority task could preempt a high priority one.

# Software failures



## Memory leak

Processes fail to entirely free up the physical memory that has been allocated to them. This effectively reduces the size of the available physical memory over time. When this becomes smaller than the minimum memory needed to support an application, it crashes.

## Incomplete specification (example Y2K)

Year = 99 (1999 or 2099)?

***Many failures (like crash, omission etc) can be caused by software bugs too.***

# Network Errors and Error Control Schemes

Failures	Causes	Metrics	Traditional Approaches
Data Losses, Deadline Misses, Node (Link) Failure, System Down	Network Congestion, Noise/Interference, Malicious Attacks	Packet Loss Rate, Deadline Miss Rate, SNR, MTTF, MTBF, MTTR	Resource Reservation, Data Redundancy (CRC, etc.), Temporal Redundancy (Retransmission, etc.), Spatial Redundancy (Replicated Nodes, MIMO, etc.)

- Omission Errors – lost/dropped messages
- Network is unreliable (especially, wireless networks)
  - Buffer overflow, Collisions at the MAC layer, Receiver out of range
- Joint approaches across OSI layers have been investigated for minimal costs [Vuran, 06][Schaar, 07]

- SNR: Signal to Noise Ratio
- MTTR: Mean Time To Recovery
- CRC: Cyclic Redundancy Check
- MIMO: Multiple-In Multiple-Out



# Classifying fault-tolerance



## Masking tolerance.

Application runs as it is. The failure does not have a visible impact. All properties (both liveness & safety) continue to hold.

## Non-masking tolerance.

Safety property is *temporarily affected*, but not liveness.

**Example 1.** Clocks lose synchronization, but recover soon thereafter.

**Example 2.** Multiple processes temporarily enter their critical sections, but thereafter, the normal behavior is restored.

# Classifying fault-tolerance



## Fail-safe tolerance

Given safety predicate is preserved, but liveness may be affected

**Example.** Due to failure, no process can enter its critical section for an indefinite period. In a traffic crossing, failure changes the traffic in both directions to red.

## Graceful degradation

Application continues, but in a “degraded” mode. Much depends on what kind of degradation is acceptable.

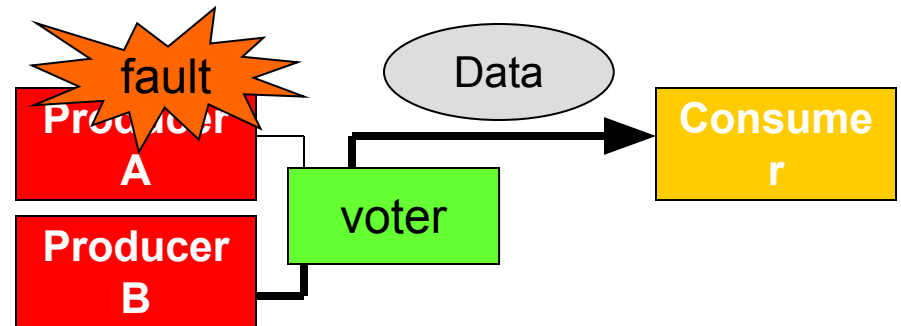
**Example.** Consider message-based mutual exclusion. Processes will enter their critical sections, but not in timestamp order.

# Conventional Approaches

- Build redundancy into hardware/software
  - Modular **Redundancy, N-Version** Programming Conventional TRM (Triple Modular Redundancy) can incur 200% overheads without optimization.
  - Replication of tasks and processes may result in overprovisioning
  - Error Control Coding
- Checkpointing and rollbacks
  - Usually accomplished through logging (e.g. messages)
  - Backward Recovery with Checkpoints cannot guarantee the completion time of a task.
- Hybrid
  - Recovery Blocks

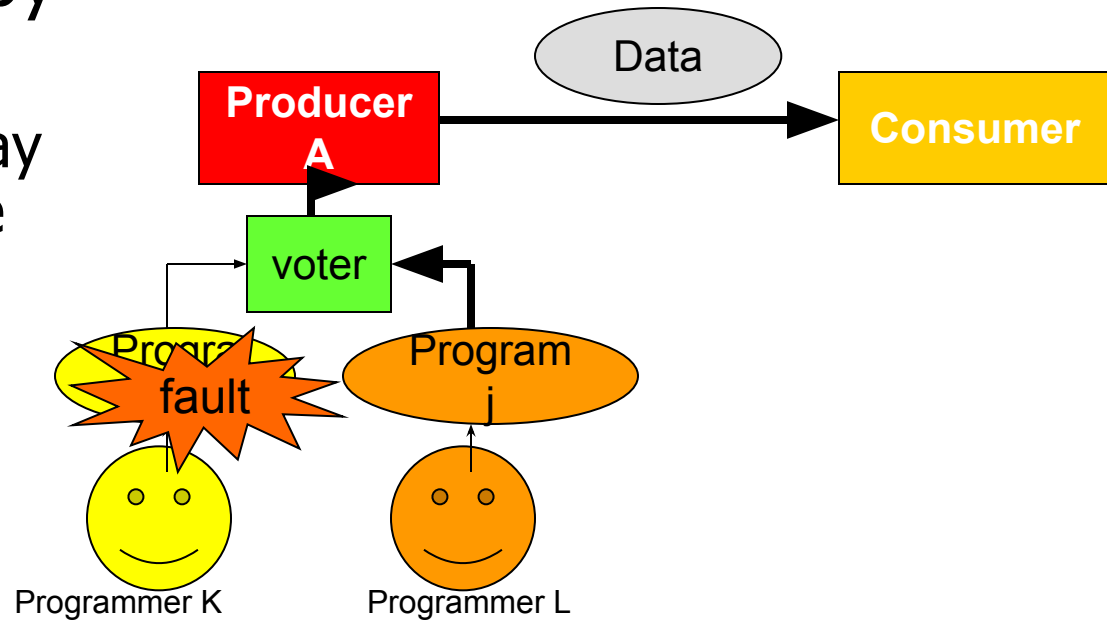
# 1) Modular Redundancy

- Modular Redundancy
  - Multiple **identical** replicas of hardware modules
  - **Voter** mechanism
    - Compare outputs and select the correct output
- Tolerate most hardware faults
- Effective but expensive



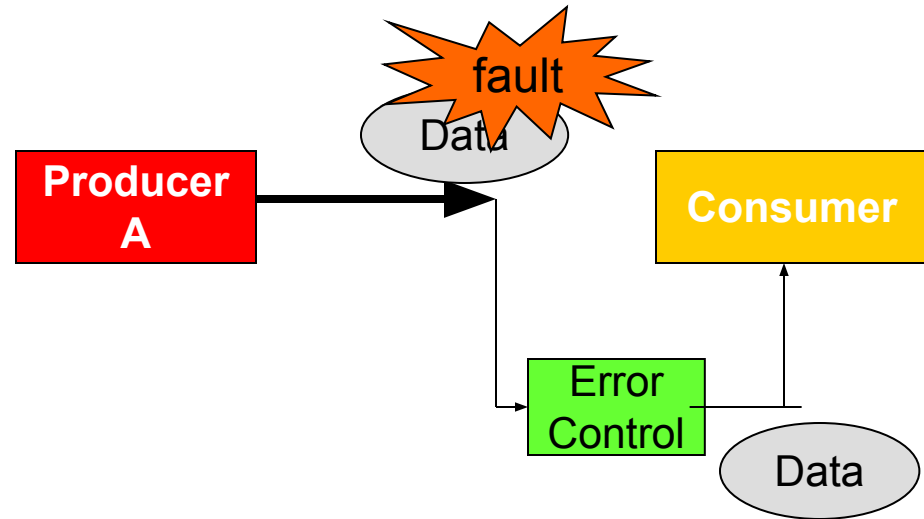
# 2) N-version Programming

- N-version Programming
  - Different versions by different teams
    - Different versions may not contain the same bugs
  - Voter mechanism
  - Tolerate some software bugs



# 3) Error-Control Coding

- Error-Control Coding
  - Replication is effective but **expensive**
  - Error-**Detection** Coding and Error-**Correction** Coding
    - (example) Parity Bit, Hamming Code, CRC
  - Much **less redundancy** than replication



# Concept: Consensus



## Reaching Agreement is a fundamental problem in distributed computing

- Mutual Exclusion
  - processes agree on which process can enter the critical section
- Leader Election
  - processes agree on which is the elected process
- Totally Ordered Multicast
  - the processes agree on the order of message delivery
- Commit or Abort in distributed transactions
- Reaching agreement about which processes have failed
- Other examples
  - Air traffic control system: all aircrafts must have the same view
  - Spaceship engine control – action from multiple control processes( “proceed” or “abort” )
  - Two armies should decide consistently to attack or retreat.

# Defining Consensus



N processes

- Every process contributes a value
- Goal: To have all processes decide on the same (some) value
  - Once made, the decision cannot be changed.

Each process  $p$  has

- input variable  $x_p$  : initially either 0 or 1
- output variable  $y_p$  : initially  $b$  ( $b$ =undecided) – can be changed only once

**Consensus problem:** design a protocol so that either

1. all non-faulty processes set their output variables to 0
2. Or non-faulty all processes set their output variables to 1
3. There is at least one initial state that leads to each outcomes 1 and 2 above



# Consensus Properties/Terms



- **Termination**
  - Every non-faulty process must eventually decide.
- **Integrity**
  - The decided value must have been proposed by some process
- **Validity**
  - If every non-faulty process proposes the same value  $v$ , then their final decision must be  $v$ .
- **Agreement**
  - The final decision of every non-faulty process must be identical.
- **Non-triviality**
  - There is at least one initial system state that leads to each of the all-0's or all-1's outcomes

# Variant of Consensus Problem



- Consensus Problem (C)
  - Each process proposes a value
  - All processes agree on a single value
- Byzantine General Problem (BG)
  - Process fails arbitrarily, byzantine failure
  - Still processes need to agree
- Interactive Consistency (IC)
  - Each process propose its value
  - All processes agree on the vector

# Solving Consensus



- No failures – trivial
  - All-to-all broadcast followed by applying a choice function
- With failures
  - One Assumption: Processes fail only by *crash-stopping*
- Synchronous system: Possible?
- Asynchronous system: ???

What about other failures??

- Omission Failures
- Byzantine Failures

# Consensus

## Synchronous vs. Asynchronous Models

### Synchronous Distributed System

- Drift of each process' local clock has a known bound
- Each step in a process takes  $lb < \text{time} < ub$
- Each message is received within bounded time

*Consensus is possible in the presence of failures!!*

- **Asynchronous** Distributed System
  - No bounds on process execution
  - The drift rate of a clock is arbitrary
  - No bounds on message transmission delays

*Consensus is impossible with the possibility of even 1 failure!!*

# Consensus in a Synchronous System

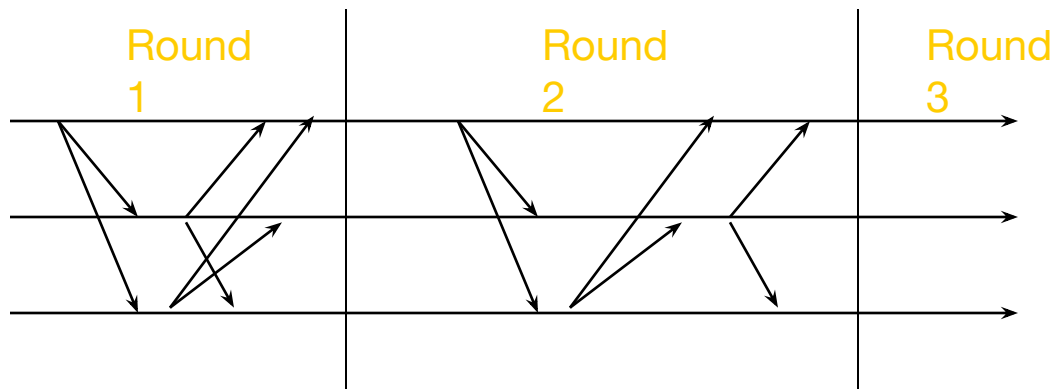


- Possible
  - With one or more faulty processes
- Solution - Basic Idea:
  - all processes exchange (multicast) what other processes tell them in several rounds
- To reach consensus with  $f$  failures, the algorithm needs to run in  $f + 1$  rounds.

# Consensus in Synchronous Systems

For a system with at most  $f$  processes crashing

- All processes are synchronized and operate in “rounds” of time.
  - Round length  $\gg$  max transmission delay.
- The algorithm proceeds in  $f+1$  rounds (with timeout), using reliable communication to all members
- $Values_i^r$ : the set of proposed values known to  $p_i$  at the beginning of round  $r$ .



# Consensus with at most $f$ failures : Synchronous Systems

achieved

For a system with at most  $f$  processes crashing

- All processes are synchronized and operate in “rounds” of time
- the algorithm proceeds in  $f+1$  rounds (with timeout), using reliable communication to all members. Round length  $\gg$  max transmission delay.
- $Values_i^r$ : the set of proposed values known to  $p_i$  at the beginning of round  $r$ .

Initially  $Values_i^0 = \{\}$  ;  $Values_i^1 = \{v_i\}$

for round = 1 to  $f+1$  do

**multicast** ( $Values_i^r - Values_i^{r-1}$ ) // iterate through processes, send each a message

$Values_i^{r+1} \sqsupseteq Values_i^r$

for each  $V_j$  received

$Values_i^{r+1} = Values_i^{r+1} \cup V_j$

end

end

$d_i = \mathbf{minimum}(Values_i^{f+1})$  // consistent minimum based on say, id (not minimum value)

# Proof: Consensus in Synchronous Systems (extra)

After  $f+1$  rounds, all non-faulty processes would have received the same set of Values.

## Proof by contradiction.

- Assume that two non-faulty processes, say  $p_i$  and  $p_j$ , differ in their final set of values (i.e., after  $f+1$  rounds)
- Assume that  $p_i$  possesses a value  $v$  that  $p_j$  does not possess.
  - $p_i$  must have received  $v$  in the **very last** round
    - Else,  $p_i$  would have sent  $v$  to  $p_j$  in that last round
    - So, in the last round: a third process,  $p_k$ , must have sent  $v$  to  $p_i$ , but then crashed before sending  $v$  to  $p_j$ .
    - Similarly, a fourth process sending  $v$  in the **last-but-one round** must have crashed; otherwise, both  $p_k$  and  $p_j$  should have received  $v$ .
    - Proceeding in this way, we infer at least one (unique) crash in each of the preceding rounds.
    - This means a total of  $f+1$  crashes, while we have assumed at most  $f$  crashes can occur => contradiction.



# Asynchronous Consensus

- Messages have arbitrary delay, processes arbitrarily slow
- **Impossible to achieve!**
  - a slow process indistinguishable from a crashed process
- Result due to Fischer, Lynch, Patterson (commonly known as FLP 85).

## 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.

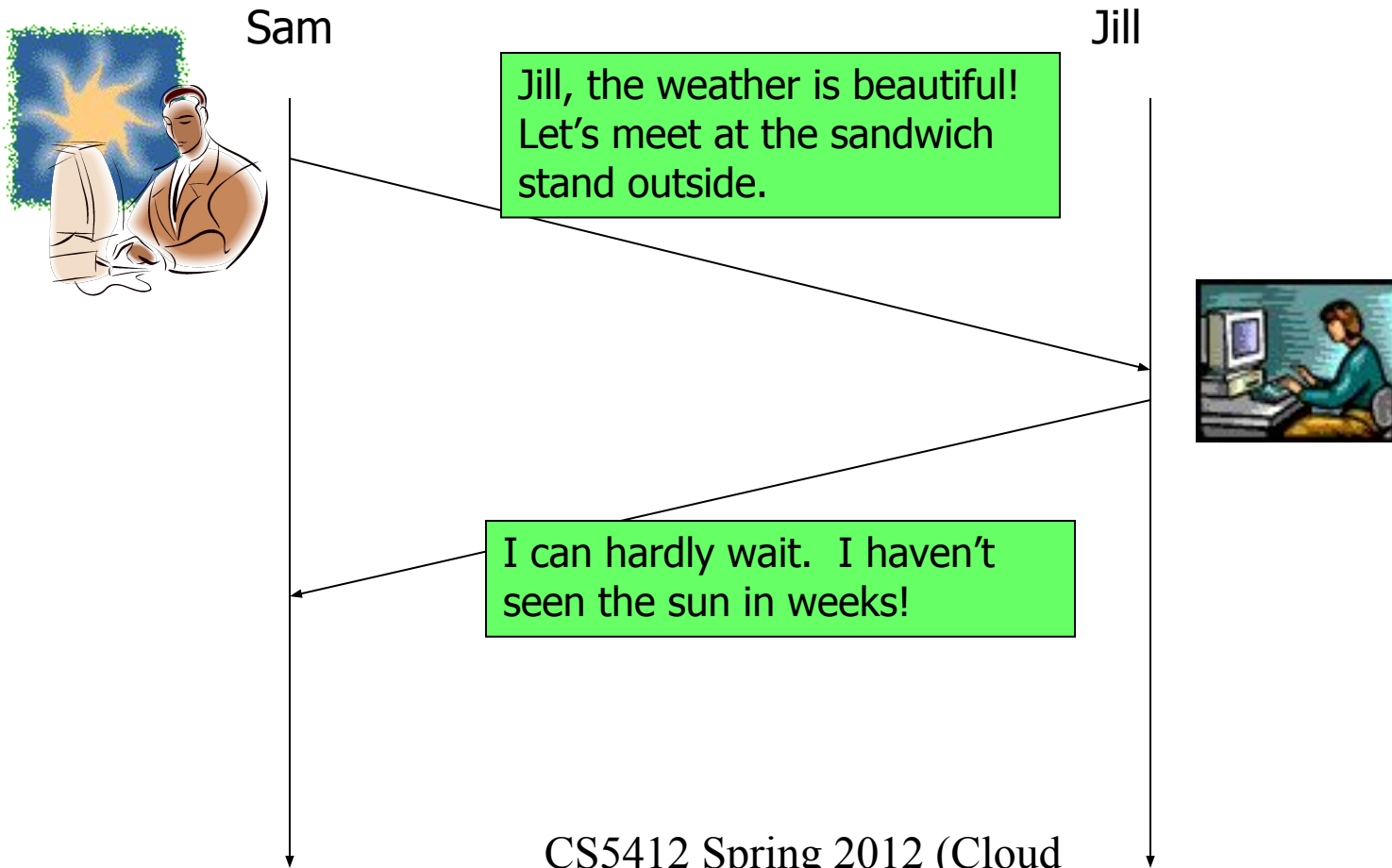
**Categories and Subject Descriptors:** C.2.2 [Computer-Communication Networks]: Network Protocols-protocol architecture; C.2.4 [Computer-Communication Networks]: Distributed Systems-distributed applications; distributed databases; network operating systems; C.4 [Performance of Systems]: Reliabil-

**Theorem:** In a purely asynchronous distributed system, the consensus problem is impossible to solve if even a single process crashes.

# Intuition Behind FLP Impossibility Theorem

- Jill and Sam will meet for lunch. They'll eat in the cafeteria unless both are sure that the weather is good
  - Jill's cubicle is inside, so Sam will send email
  - Both have lots of meetings, and might not read email. So she'll acknowledge his message.
  - They'll meet inside if one or the other is away from their desk and misses the email.
- Sam sees sun. Sends email. Jill acks's. Can they meet outside?

# Sam and Jill

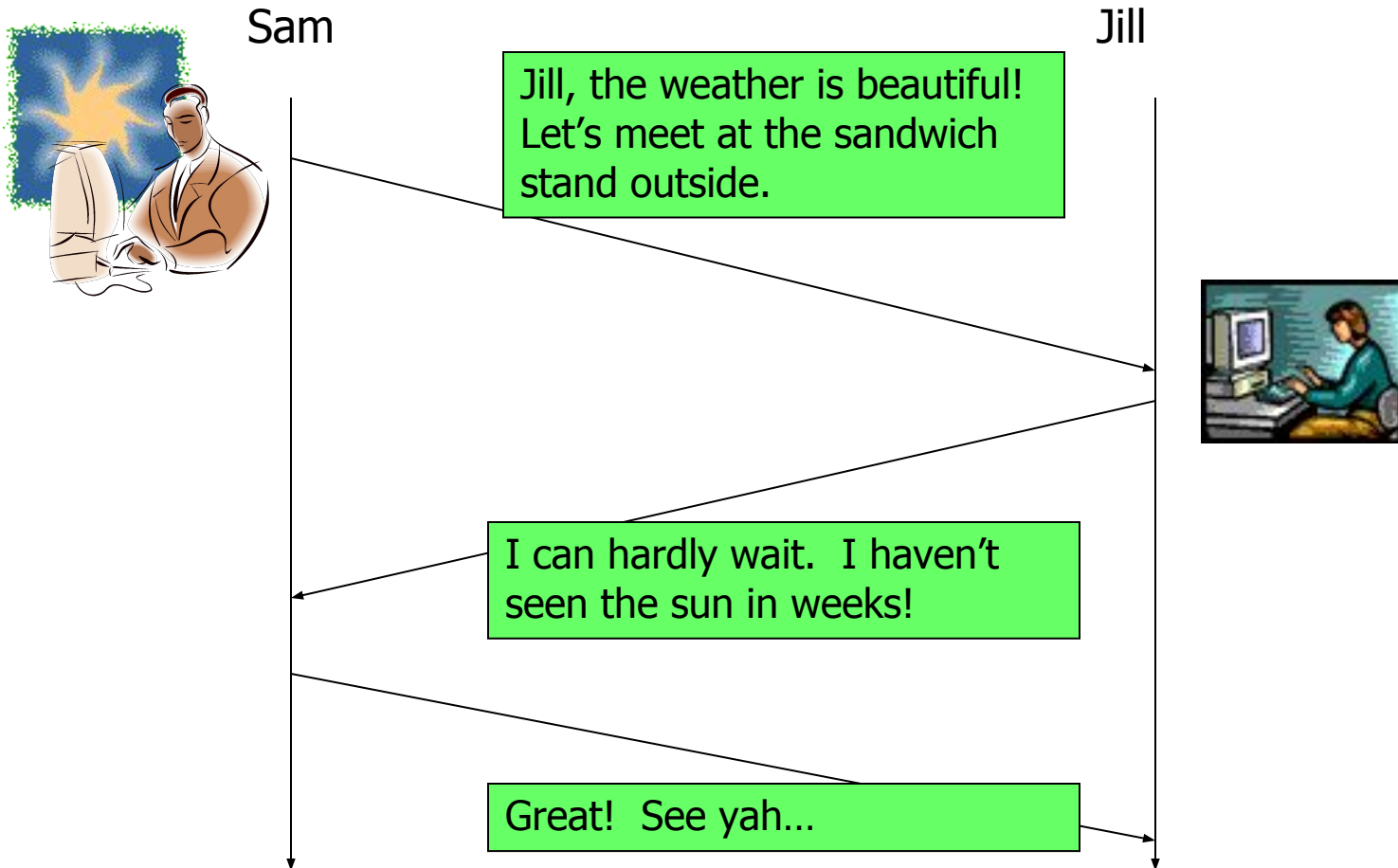


# **They eat inside! Sam reasons:**




- “Jill sent an acknowledgement but doesn’t know if I read it
- “If I didn’t get her acknowledgement I’ll assume she didn’t get my email
- “In that case I’ll go to the cafeteria
- “She’s uncertain, so she’ll meet me there

# Sam had better send an Ack



# Why didn't this help?



- Jill got the ack... but she realizes that Sam won't be sure she got it
- Being unsure, he's in the same state as before
- So he'll go to the cafeteria, being dull and logical. And so she meets him there.

# New and improved protocol

- Jill sends an ack. Sam acks the ack. Jill acks the ack of the ack....
- Suppose that noon arrives and Jill has sent her 117'th ack.
  - Should she assume that lunch is outside in the sun, or inside in the cafeteria?

# How Sam and Jill's romance ended



Jill, the weather is beautiful!  
Let's meet at the sandwich  
stand outside.

I can hardly wait. I haven't seen the sun  
in weeks!

Great! See yah...

Yup...

Got that...

...

Oops, too late for lunch

Maybe tomorrow?





# Things we just can't do



- We can't detect failures in a trustworthy, consistent manner
- We can't reach a state of "common knowledge" concerning something not agreed upon in the first place
- We can't guarantee agreement on things (election of a leader, update to a replicated variable) in a way certain to tolerate failures

# But what does it mean?



- In formal proofs, an algorithm is totally correct if
  - It computes the right thing
  - And it always terminates
- When we say something is possible, we mean “there is a totally correct algorithm” solving the problem
- FLP proves that any fault-tolerant algorithm solving consensus has runs that never terminate
  - These runs are extremely unlikely (“probability zero”)
  - Yet they imply that we can’t find a totally correct solution
  - And so “consensus is impossible” (“not always possible”)
- In practice, fault-tolerant consensus is ..
  - Definitely possible.
  - E.g. **Paxos [Lamport 1998, 2001] that has become quite popular – discussed later!**

# FLP Proof Sketch (extra):

## Terms

- **Bivalent and Univalent states:** A decision state is bivalent, if starting from that state, there exist two distinct executions leading to two distinct decision values 0 or 1. Otherwise it is univalent.  
Bivalent ---> outcome is unpredictable
- **Process:** has state
- **Network:** Global buffer (processes put and get messages)
- **Configuration** -- global state (state for each process + state of global buffer)
- **Atomic Events** -- receipt of message by process p, processing of message (may change state), send out all needed messages from p
- **Schedule:** sequence of atomic events

**Lemma 1: Schedules are commutative**

**Lemma 2: There exists an initial configuration that is bivalent.**

**Lemma 3: Starting from a bivalent config., there is always another bivalent config. that is reachable**

# The PAXOS Algorithm

## -- Towards a Practical Approach to Consensus

Landmark paper by Leslie Lamport (1998)

- Does not solve pure consensus problem (impossibility);  
But, provides consensus with a twist
- Paxos provides safety and eventual liveness
  - Safety: Consensus is not violated
  - Eventual Liveness: If things go well sometime in the future (messages, failures, etc.), there is a good chance consensus will be reached. But there is no guarantee.
- FLP result still applies: Paxos is not *guaranteed* to reach Consensus (ever, or within any bounded time)
- Used in Zookeeper (Yahoo!), Google Chubby, and many other companies

# The Paxos Strategy



- Paxos has **rounds**; each round has a unique ballot id
- Rounds are asynchronous
  - Time synchronization not required
  - If you're in round  $j$  and hear a message from round  $j+1$ , abort everything and move over to round  $j+1$
  - Use timeouts; may be pessimistic
- Each round itself broken into phases (which are also asynchronous)
  - Phase 1: A leader is elected (**Election**)
  - Phase 2: Leader proposes a value, processes ack (**Bill**)
  - Phase 3: Leader multicasts final value (**Law**)
- <http://research.microsoft.com/en-us/um/people/lamport/pubs/paxos-simple.pdf>

***MORE DETAILS LATER!!***

# Failure detection



The design of fault-tolerant algorithms will be simple if processes can detect failures.

- Impossibility results assume failures cannot be observed.
- In synchronous systems with bounded delay channels, crash failures can **definitely be detected** using timeouts.
- In asynchronous distributed systems, the detection of **crash failures** is imperfect.



# Designing failure detectors



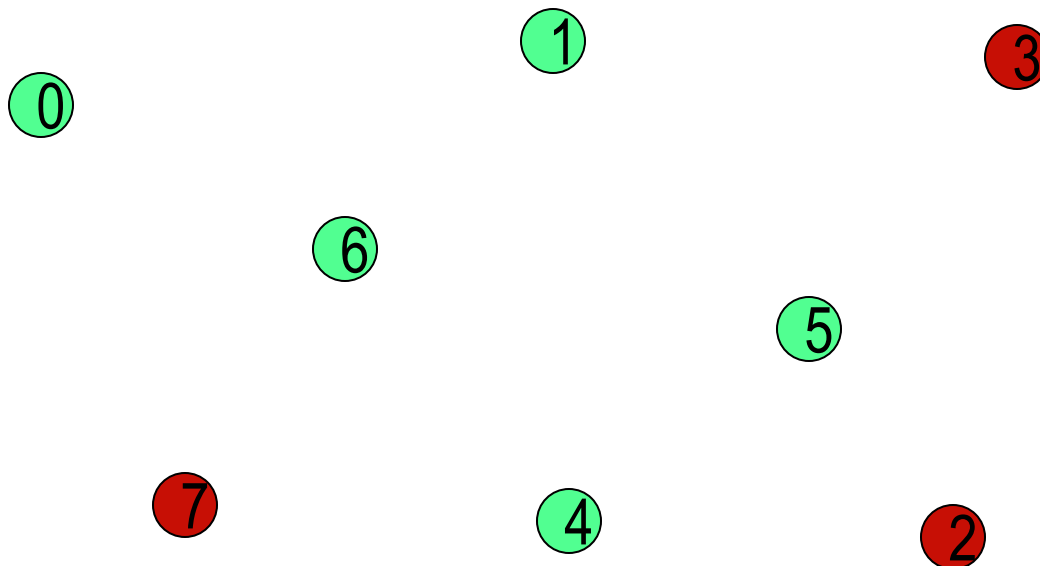
Processes carry a *Failure Detector* to detect crashed processes.

## *Desirable Properties of a failure detector:*

- *Completeness* – Every crashed process is suspected
- *Accuracy* – No correct process is suspected.
- *Other factors*
  - Speed -- time to first detection of a failure
  - Overhead -- load on member process, network message load

# Example

---



0 suspects  $\{1,2,3,7\}$  to have failed.

Does this satisfy **completeness**?

Does this satisfy **accuracy**?



# Classification of completeness



- **Strong completeness.** Every crashed process is eventually suspected by *every correct process*, and remains a suspect thereafter.
- **Weak completeness.** Every crashed process is eventually suspected by *at least one* correct process, and remains a suspect thereafter.

*Note that we don't care what mechanism is used for suspecting a process.*

# Classification of accuracy



- **Strong accuracy.** No correct process is ever suspected.
- **Weak accuracy.** There is at least one correct process that is never suspected.

# Eventual accuracy



A failure detector is *eventually strongly accurate*, if there exists a time **T** after which no correct process is suspected.

*(Before that time, a correct process be added to and removed from the list of suspects any number of times)*

A failure detector is *eventually weakly accurate*, if there exists a time **T** after which **at least one process** is no more suspected.

# Classifying failure detectors



**Perfect P.** (Strongly) Complete and strongly accurate

**Strong S.** (Strongly) Complete and weakly accurate

**Eventually perfect  $\diamond P$ .**

(Strongly) Complete and eventually strongly accurate

**Eventually strong  $\diamond S$**

(Strongly) Complete and eventually weakly accurate

Other classes are feasible: W (weak completeness) and weak accuracy) and  $\diamond W$

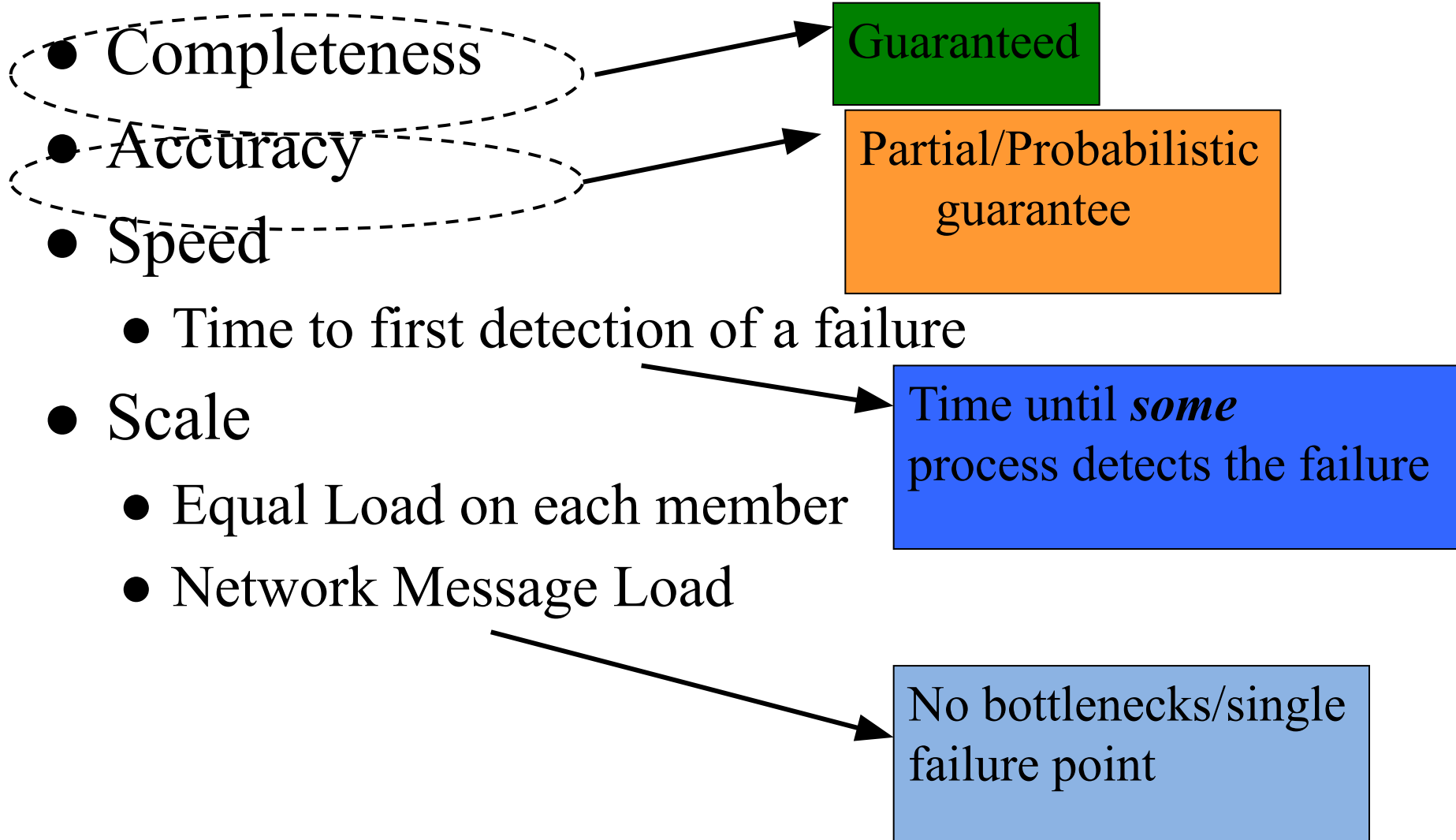
# Distributed Failure Detectors: Desired Properties

- Completeness
- Accuracy
- Speed
  - Time to first detection of a failure
- Scale
  - Equal Load on each member
  - Network Message Load

Completeness and Accuracy impossible together in lossy networks [Chandra and Toueg]

If possible, then can solve consensus! (but consensus is known to be unsolvable in asynchronous systems)

# Real Failure Detectors



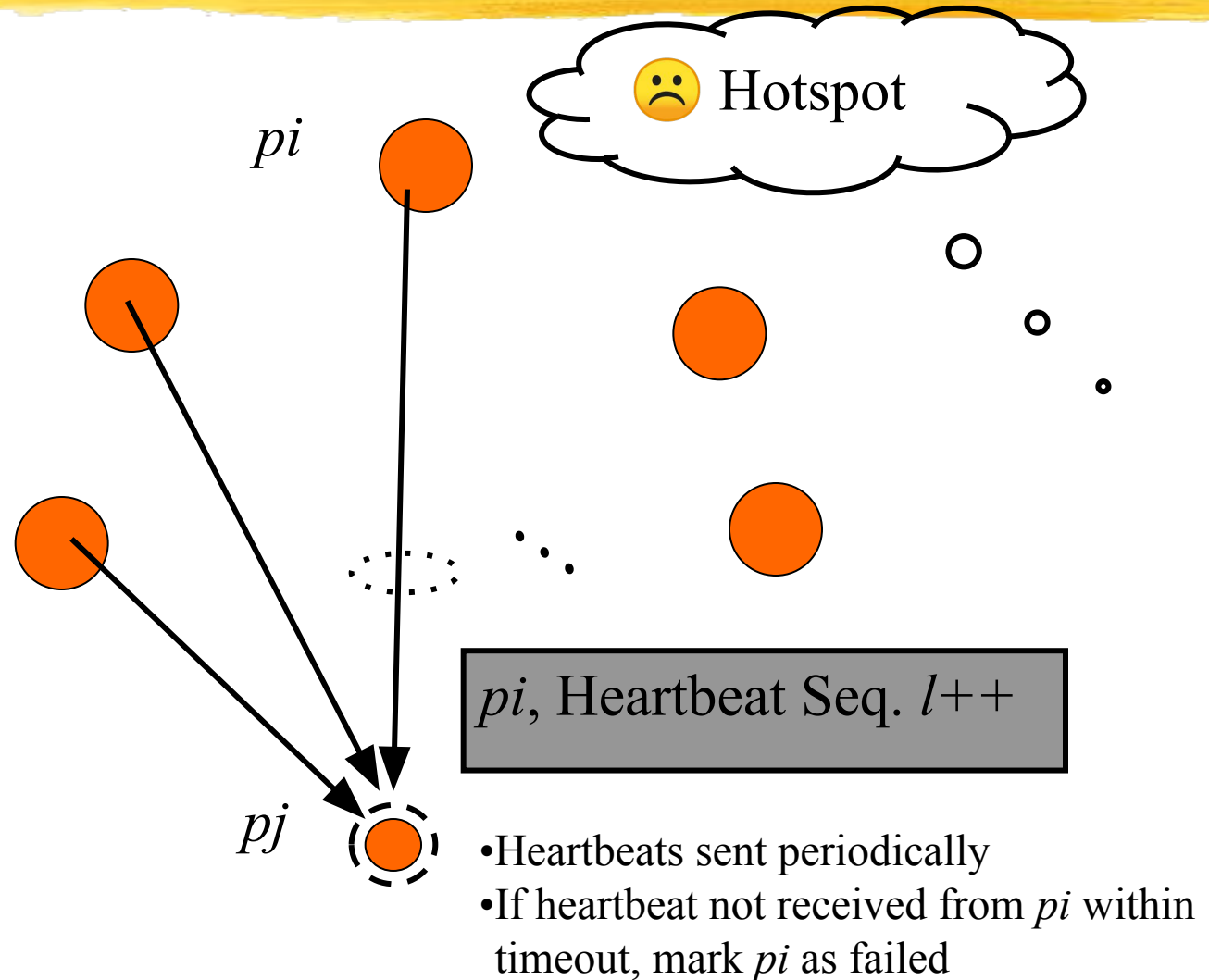
# Detection of crash failures



Failure can be detected using **heartbeat messages** (periodic “**I am alive**” broadcast) and **timeout**

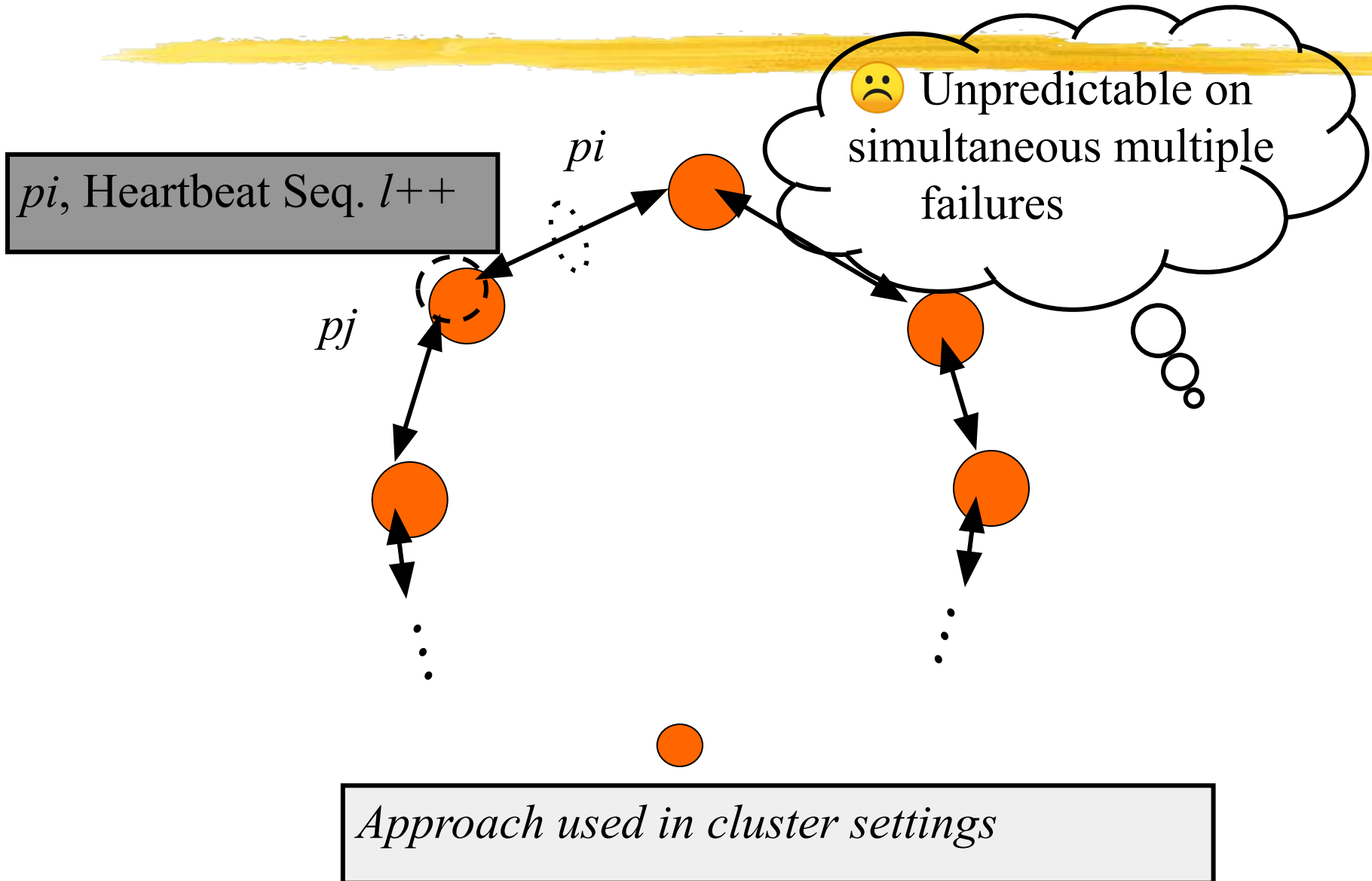
- if processor speed has a known lower bound
- channel delays have a known upper bound.

# Centralized Heartbeating

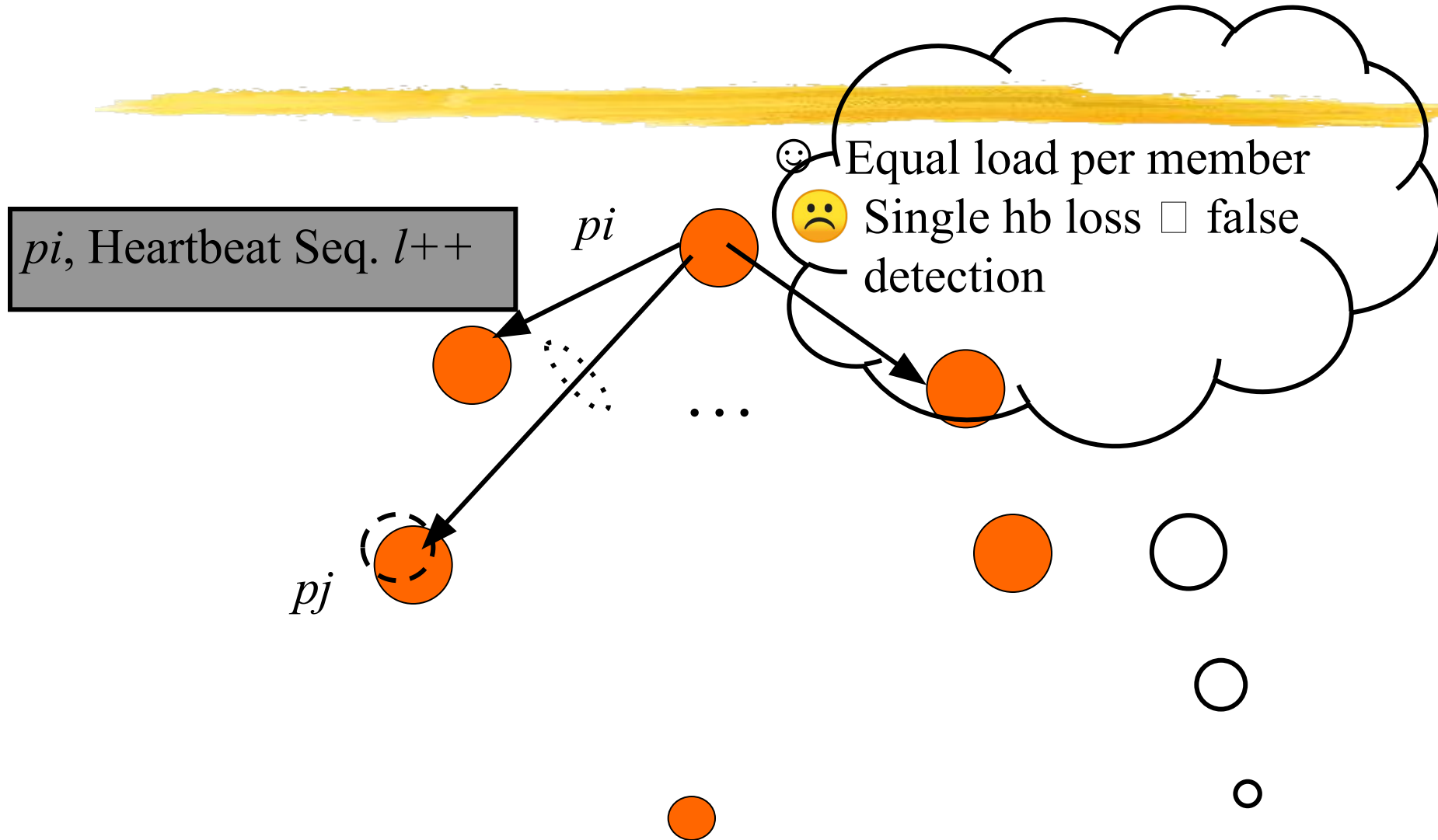




# Ring Heartbeating



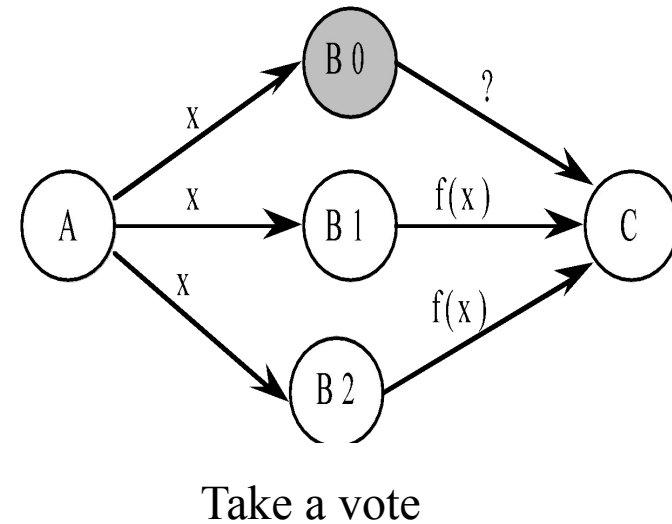
# All-to-All Heartbeating



Variant - ***gossip style heartbeating*** (heartbeats with a member subset) -- AWS???  
Determine gossip-period; send  $o(N)$  heartbeats to a subset every gossip period

# Tolerating crash failures

Triple modular redundancy (TMR) for masking any single failure. *N*-modular redundancy masks up to *m* failures, when  $N = 2m + 1$ .



**What if the voting unit fails?**

# Detection of omission failures



For **FIFO** channels: Use **sequence numbers** with messages.

(1, 2, 3, 5, 6 ... )  $\Rightarrow$  message 4 is missing

**Non-FIFO bounded delay channels** - use **timeout**

What about non-FIFO channels for which the **upper bound of the delay is not known?**

Use *unbounded sequence numbers* and **acknowledgments**.

But acknowledgments may be lost too!

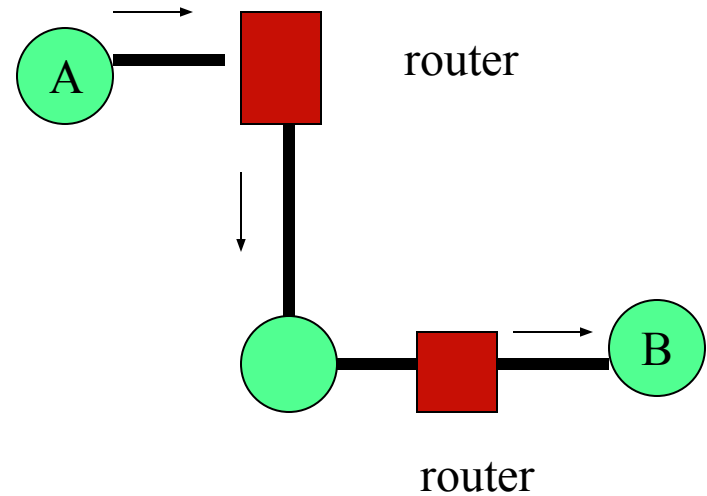
# Tolerating omission failures

## A real example

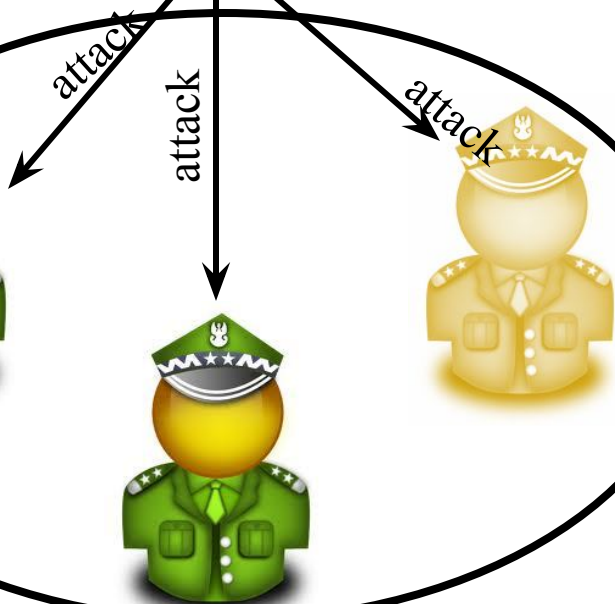
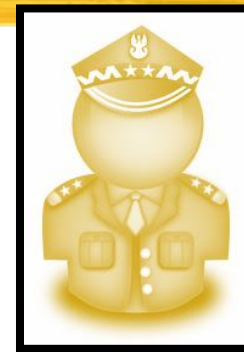
### *A central issue in networking*

Routers may drop messages, but **reliable end-to-end transmission** is an important requirement. If the sender does not receive an **ack** within a time period, it retransmits (it may so happen that the was not lost, so a duplicate is generated).

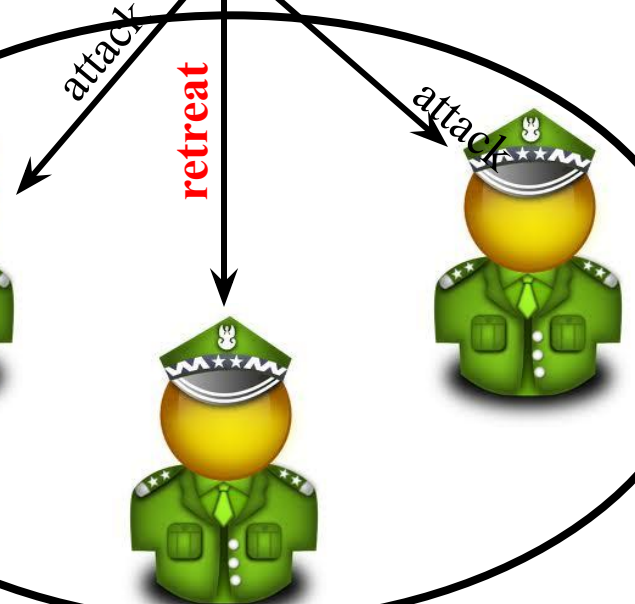
This implies, the communication must tolerate **Loss, Duplication, and Re-ordering** of messages



# Byzantine Generals Problem



Lieutenants agree on what the commander says



Lieutenants agree on what the commander says

# Byzantine Generals Problem

**The Byzantine generals problem** • In the informal statement of the *Byzantine generals problem* [Lamport *et al.* 1982], three or more generals are to agree to attack or to retreat. One, the commander, issues the order. The others, lieutenants to the commander, must decide whether to attack or retreat. But one or more of the generals may be ‘treacherous’ – that is, faulty. If the commander is treacherous, he proposes attacking to one general and retreating to another. If a lieutenant is treacherous, he tells one of his peers that the commander told him to attack and another that they are to retreat.

The Byzantine generals problem differs from consensus in that a distinguished process supplies a value that the others are to agree upon, instead of each of them proposing a value. The requirements are:

*Termination:* Eventually each correct process sets its decision variable.

*Agreement:* The decision value of all correct processes is the same: if  $p_i$  and  $p_j$  are correct and have entered the *decided* state, then  $d_i = d_j$  ( $i, j = 1, 2, \dots, N$ ).

*Integrity:* If the commander is correct, then all correct processes decide on the value that the commander proposed.

# Replication

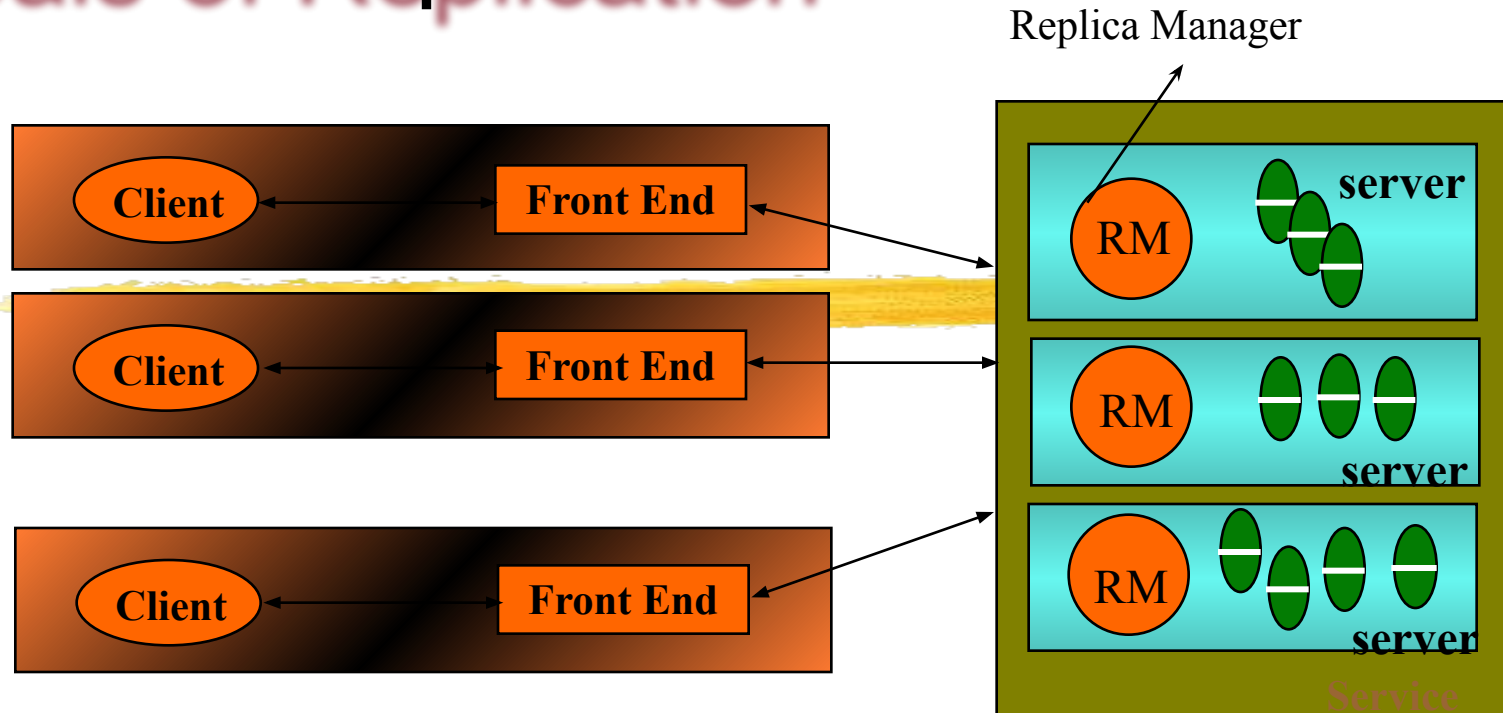
- ❖ Enhances a service by replicating data
  - ❖ Increased Availability
    - ❖ Of service. When servers fail or when the network is partitioned.
  - ❖ Fault Tolerance
    - ❖ Under the fail-stop model, if up to  $f$  of  $f+1$  servers crash, at least one is alive.
  - ❖ Load Balancing
    - ❖ One approach: Multiple server IPs can be assigned to the same name in DNS, which returns answers round-robin.

$P$ : probability that one server fails =  $1 - P$  = availability of service.  
e.g.  $P = 5\% \Rightarrow$  service is available 95% of the time.

$P^n$ : probability that  $n$  servers fail =  $1 - P^n$  = availability of service.  
e.g.  $P = 5\%$ ,  $n = 3 \Rightarrow$  service available 99.875% of the time



# Goals of Replication



## ❖ Replication Transparency

User/client need not know that multiple physical copies of data exist.

## ❖ Replication Consistency

Data is consistent on all of the replicas (or is converging towards becoming consistent)

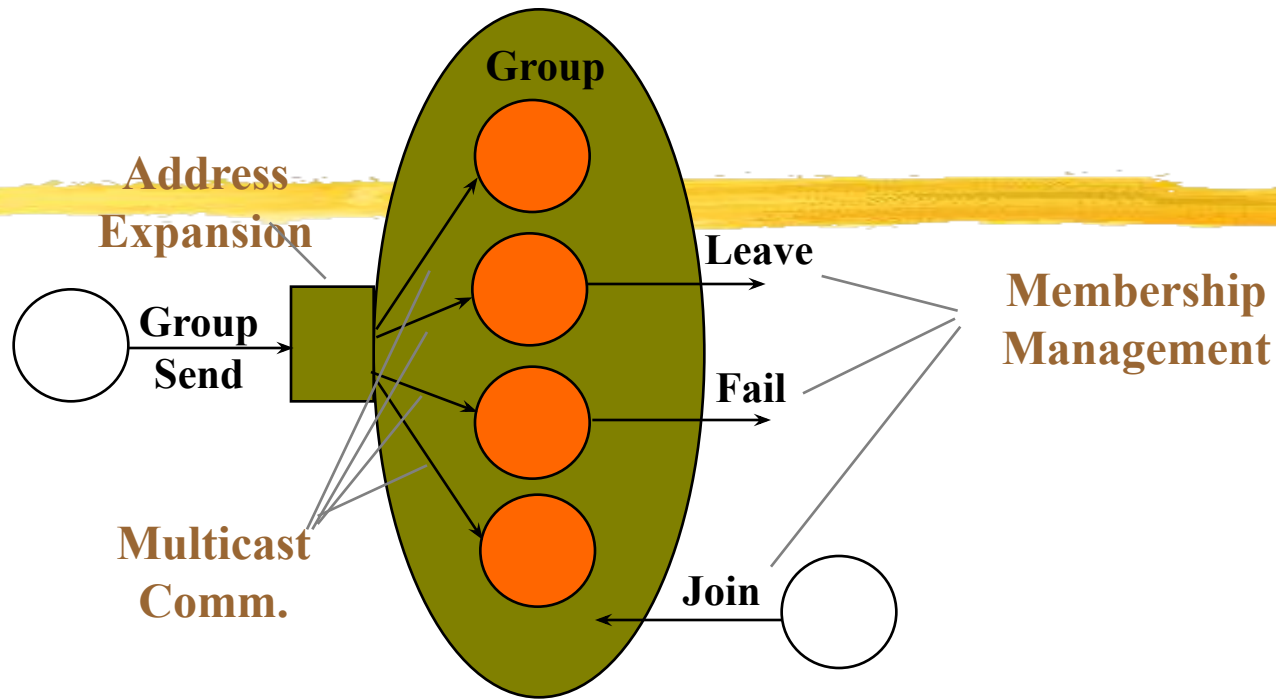
# Replication Management

- ❖ Request Communication
  - ❖ Requests can be made to a single RM or to multiple RMs
- ❖ Coordination: The RMs decide
  - ❖ whether the request is to be applied
  - ❖ the order of requests
    - ❖ FIFO ordering: If a FE issues  $r$  then  $r'$ , then any correct RM handles  $r$  and then  $r'$ .
    - ❖ Causal ordering: If the issue of  $r$  “happened before” the issue of  $r'$ , then any correct RM handles  $r$  and then  $r'$ .
    - ❖ Total ordering: If a correct RM handles  $r$  and then  $r'$ , then any correct RM handles  $r$  and then  $r'$ .
- ❖ Execution: The RMs execute the request (often they do this tentatively).

# Replication Management

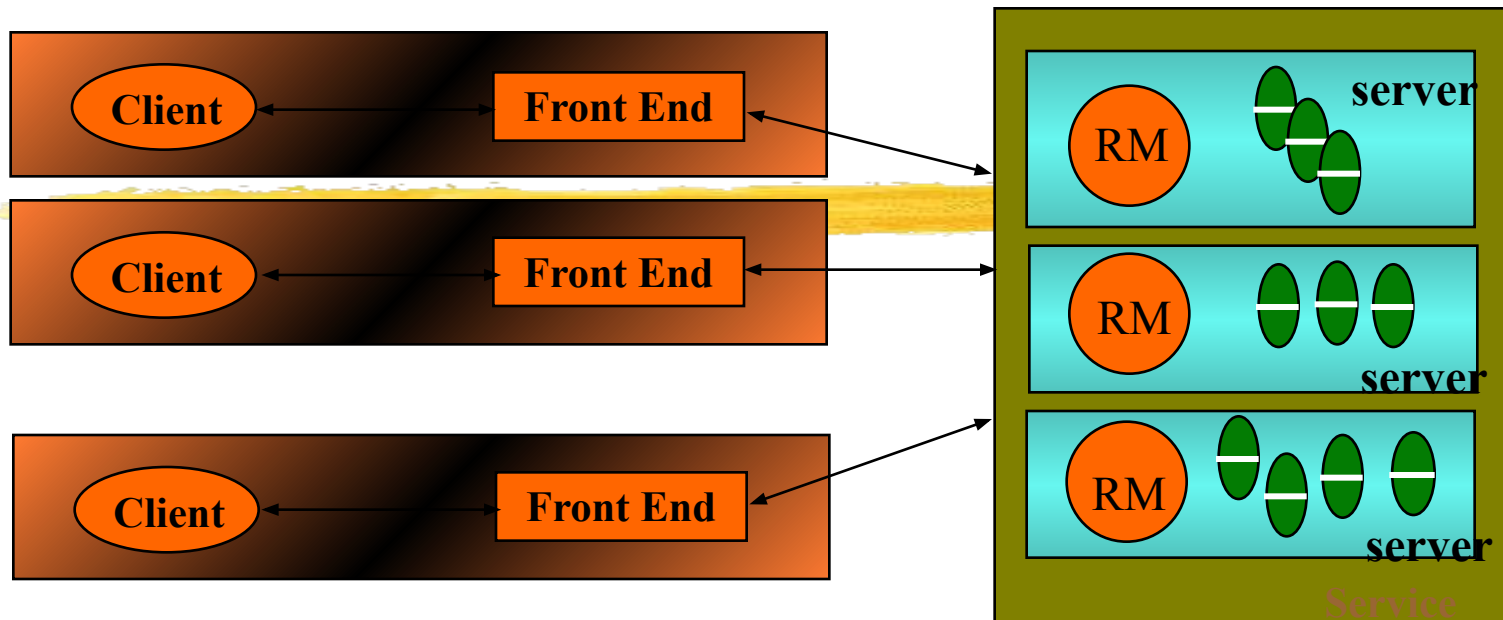
- ❖ Agreement: The RMs attempt to reach consensus on the effect of the request.
  - ❖ E.g., Two phase commit through a coordinator
  - ❖ If this succeeds, effect of request is made permanent
- ❖ Response
  - ❖ One or more RMs responds to the front end.
  - ❖ The first response to arrive is good enough because all the RMs will return the same answer.
  - ❖ Thus each RM is a **replicated state machine**
    - “Multiple copies of the same State Machine begun in the Start state, and receiving the same Inputs in the same order will arrive at the same State having generated the same Outputs.” [Wikipedia, Schneider 90]

# Group Communication: A building block



- ❖ “Member”= process (e.g., an RM)
- ❖ Static Groups: group membership is pre-defined
- ❖ Dynamic Groups: Members may join and leave, as necessary

# Replication using GC



Need consistent updates to all copies of an object

- Linearizability
- Sequential Consistency

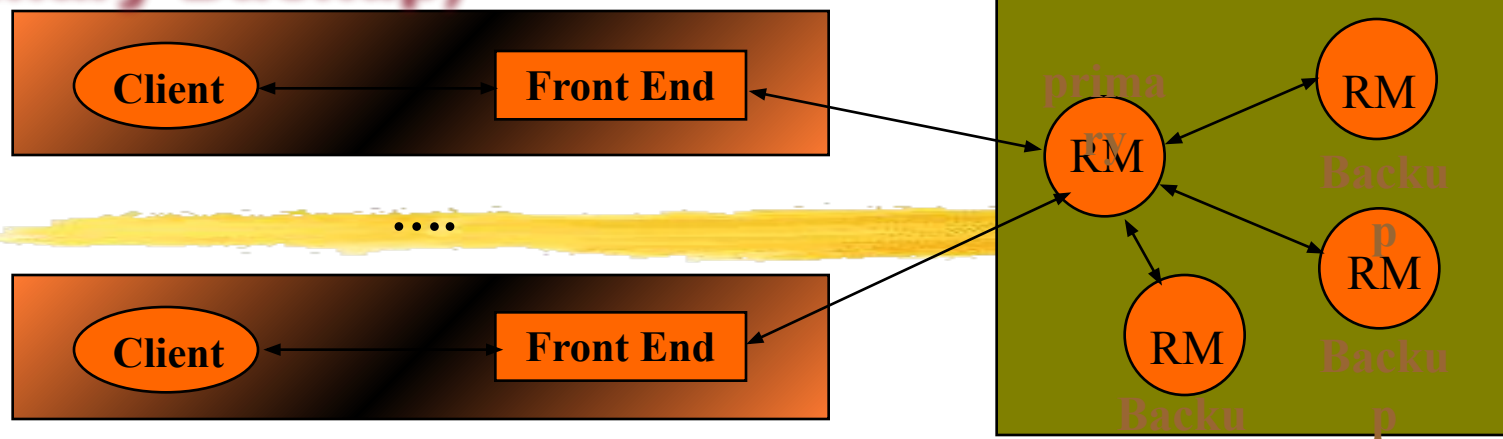
# Linearizability

- ❖ Let the sequence of read and update operations that client  $i$  performs in some execution be  $o_{i1}, o_{i2}, \dots$
- ❖ “Program order” for the client
- ❖ A replicated shared object service is **linearizable** if for any execution (real), there is some interleaving of operations (virtual) issued by all clients that:
  - ❑ meets the specification of a single correct copy of objects
  - ❑ is consistent with the real times at which each operation occurred during the execution
- ❑ Main goal: any client will see (at any point of time) a copy of the object that is correct and consistent

# Sequential Consistency

- ❖ The real-time requirement of **linearizability** is hard, if not impossible, to achieve in real systems
- ❖ A less strict criterion is **sequential consistency**: A replicated shared object service is **sequentially consistent** if for any execution (real), there is some interleaving of clients' operations (virtual) that:
  - meets the specification of a single correct copy of objects
  - is consistent with the program order in which each individual client executes those operations.
- ❖ This approach does not require absolute time or total order. Only that for each client the order in the sequence be consistent with that client's program order (~ FIFO).
- ❖ Linearizability implies sequential consistency. Not vice-versa!
- ❖ Challenge with guaranteeing seq. cons.?
  - ❖ Ensuring that all replicas of an object are consistent.

# Passive Replication (Primary-Backup)



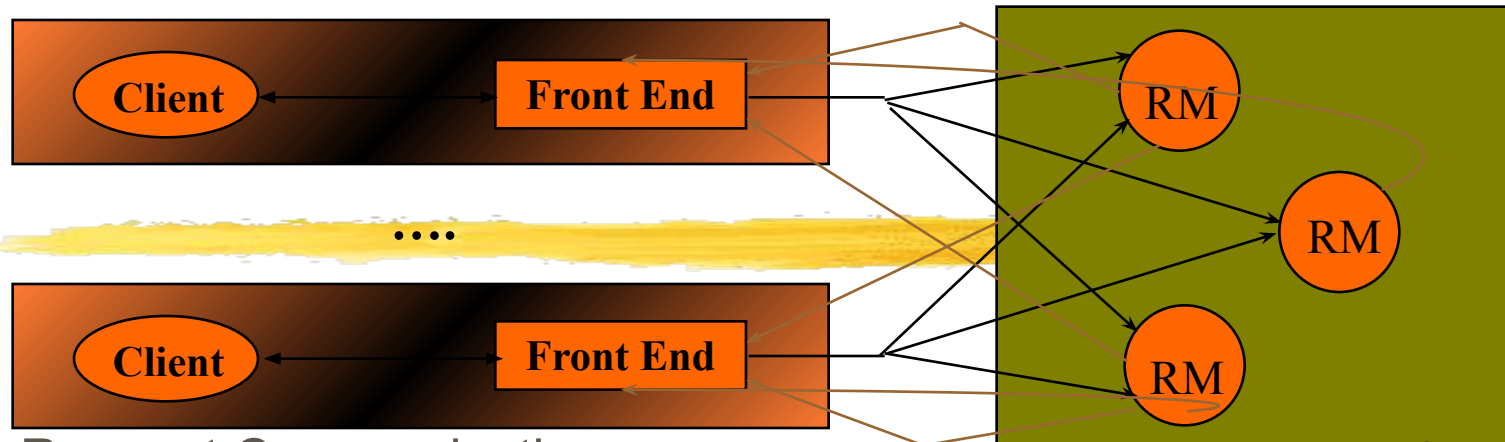
- ❖ **Request Communication:** the request is issued to the primary RM and carries a unique request id.
- ❖ **Coordination:** Primary takes requests atomically, in order, checks id (resends response if not new id.)
- ❖ **Execution:** Primary executes & stores the response
- ❖ **Agreement:** If update, primary sends updated state/result, req-id and response to all backup RMs (1-phase commit enough).
- ❖ **Response:** primary sends result to the front end



# Fault Tolerance in Passive Replication

- ❖ The system implements linearizability, since the primary sequences operations in order.
- ❖ If the primary fails, a backup becomes primary by leader election, and the replica managers that survive agree on which operations had been performed at the point when the new primary takes over.
  - ❖ The above requirement can be met if the replica managers (primary and backups) are organized as a group and if the primary uses view-synchronous group communication to send updates to backups.
- ❖ Thus the system remains linearizable in spite of crashes

# Active Replication



- ❖ Request Communication: The request contains a unique identifier and is multicast to all by a reliable totally-ordered multicast.
- ❖ Coordination: Group communication ensures that requests are delivered to each RM in the same order (but may be at different physical times!).
- ❖ Execution: Each replica executes the request. (Correct replicas return same result since they are running the same program, i.e., they are *replicated protocols* or *replicated state machines*)
- ❖ Agreement: No agreement phase is needed, because of multicast delivery semantics of requests
- ❖ Response: Each replica sends response directly to FE

# FT via Active Replication

- ❖ RMs work as replicated state machines, playing equivalent roles. That is, each responds to a given series of requests in the same way. One way of achieving this is by running the same program code at all RMs (but only one way – why?).
- ❖ If any RM crashes, state is maintained by other correct RMs.
- ❖ This system implements sequential consistency
  - ❖ The total order ensures that all correct replica managers process the same set of requests in the same order.
  - ❖ Each front end's requests are served in FIFO order (because the front end awaits a response before making the next request).
- ❖ So, requests are FIFO-total ordered.
- ❖ Caveat (Out of band): If clients are multi-threaded and communicate with one another while waiting for responses from the service, we may need to incorporate causal-total ordering.

# **EXTRA SLIDES**



# **Proof of FLP '83 (Impossibility of Consensus in an Asynchronous System)**

**(adapted from UIUC CS425 - Indranil Gupta)**

# Recall

**Asynchronous system:** All message delays and processing delays can be arbitrarily long or short.

**Consensus:**

- Each process  $p$  has a **state**
  - program counter, registers, stack, local variables
  - input register  $x_p$  : initially either 0 or 1
  - output register  $y_p$  : initially  $b$  (undecided)
- Consensus Problem: design a protocol so that either
  - all processes set their output variables to 0 (all-0's)
  - Or all processes set their output variables to 1 (all-1's)
  - Non-triviality: at least one initial system state leads to each of the above two outcomes

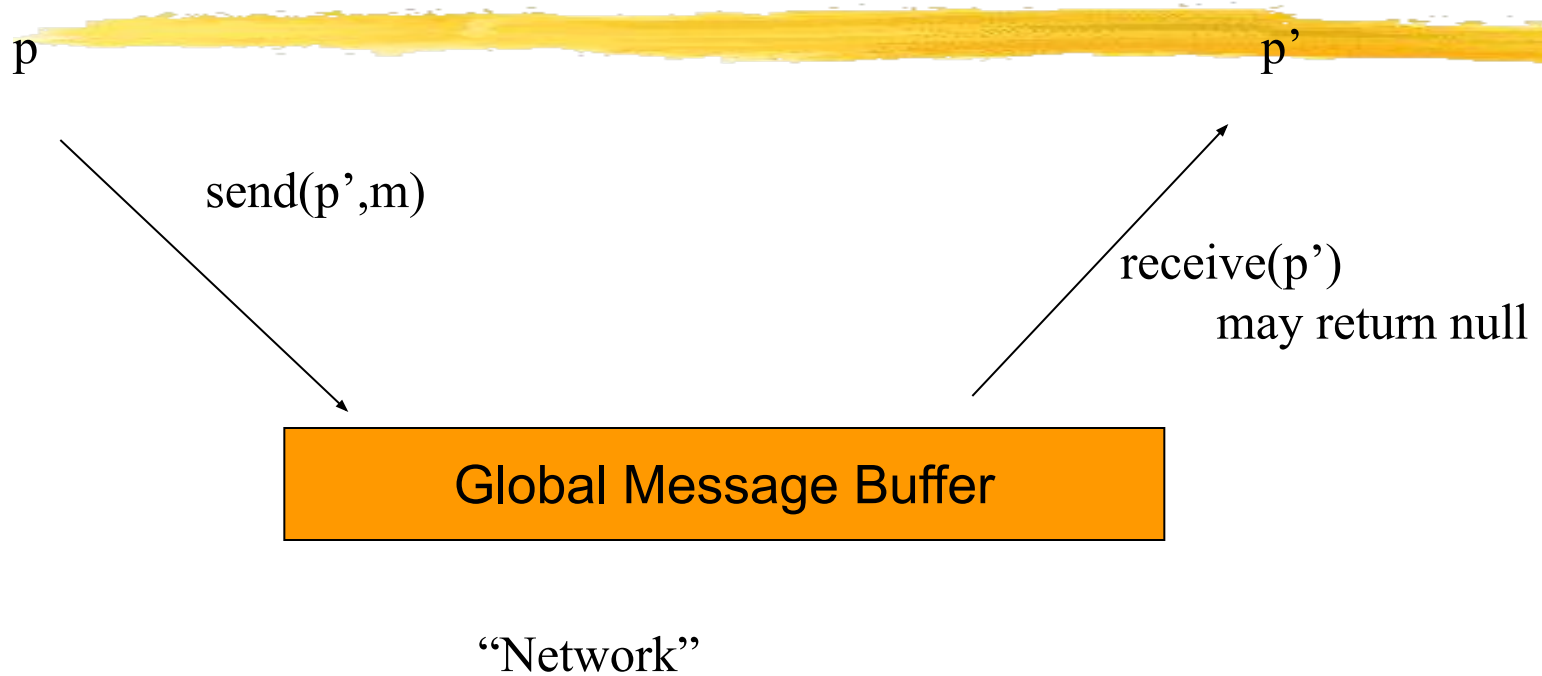


# Proof Setup

- For impossibility proof, OK to consider
  1. more restrictive system model, and
  2. easier problem
    - Why is this is ok?



# Network

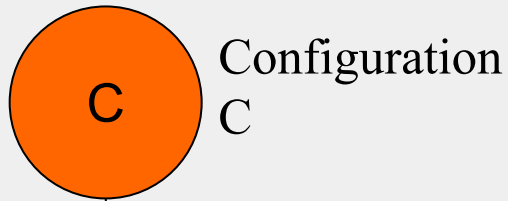




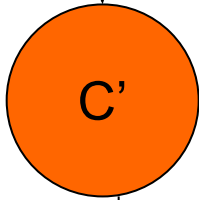
# States

- State of a process
- **Configuration=global state.** Collection of states, one for each process; alongside state of the global buffer.
- Each **Event** (different from Lamport events) is atomic and consists of three steps
  - receipt of a message by a process (say p)
  - processing of message (may change recipient's state)
  - sending out of all necessary messages by p
- **Schedule:** sequence of events

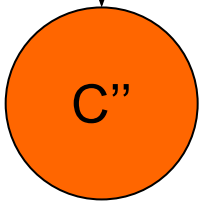




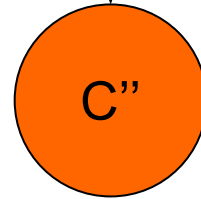
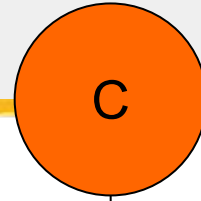
Event  
 $e'=(p',m')$



Event  $e''=(p'',m'')$

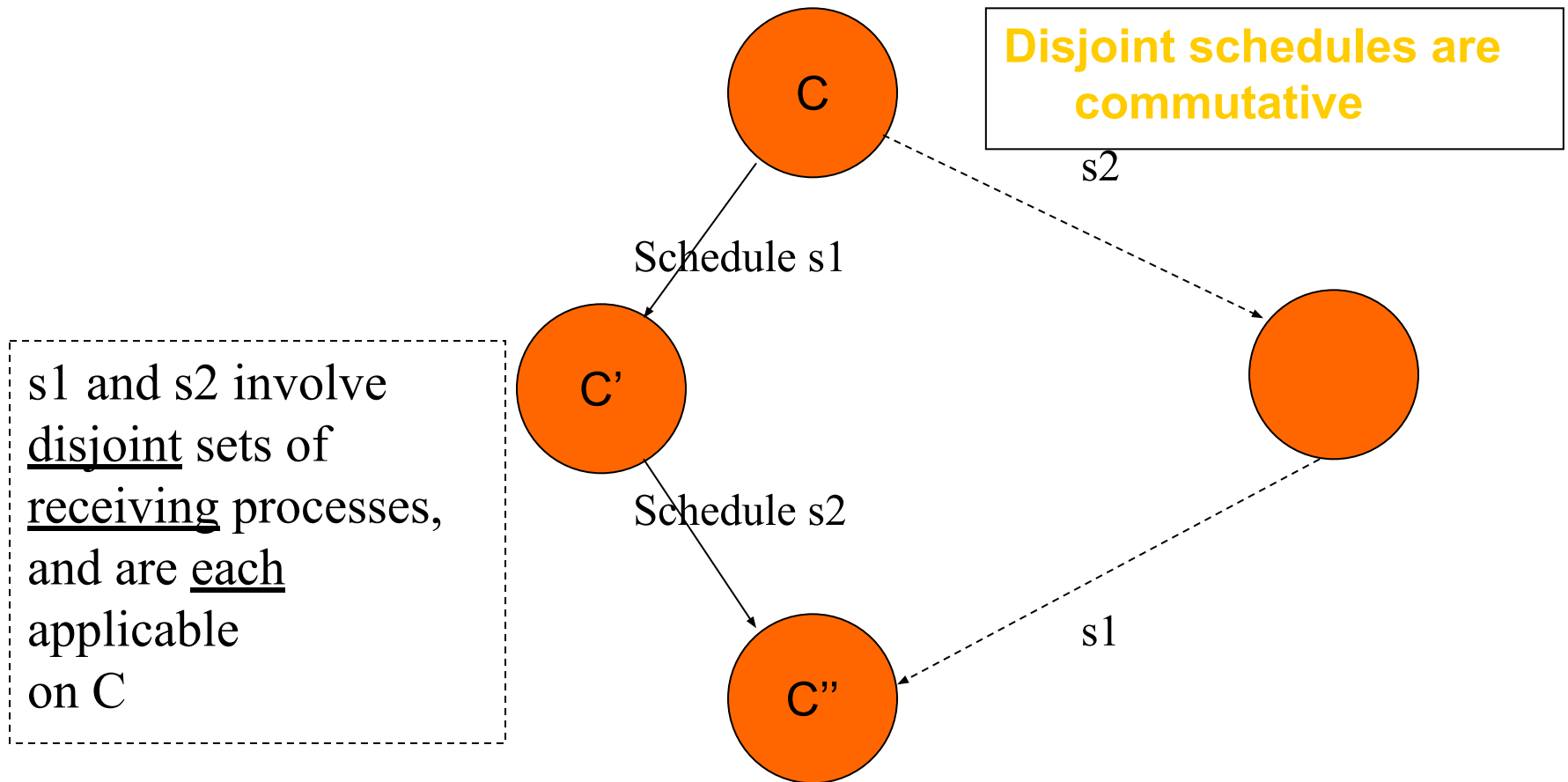


Schedule  
 $s=(e',e'')$



Equivalence

# Lemma 1



# Easier Consensus Problem

Easier Consensus Problem: **some** process eventually sets  $y_p$  to be 0 or 1

**Only one process crashes** – we're free to choose which one



# Easier Consensus Problem

- Let config.  $C$  have a set of decision values  $V$  reachable from it
  - If  $|V| = 2$ , config.  $C$  is bivalent
  - If  $|V| = 1$ , config.  $C$  is 0-valent or 1-valent, as is the case
- **Bivalent** means **outcome is unpredictable**



# What the FLP proof shows

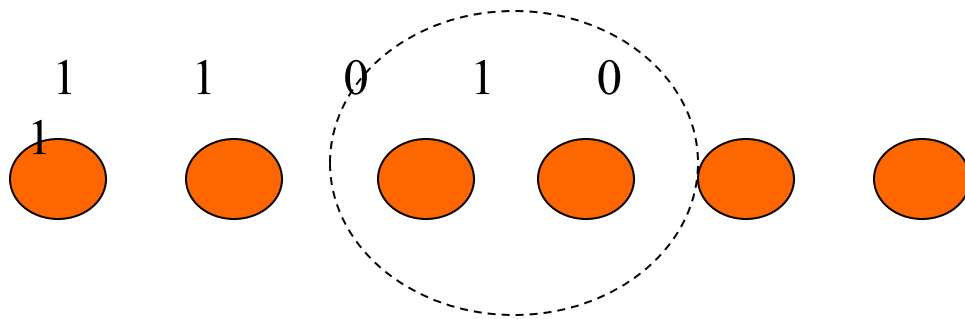
1. There exists an initial configuration that is bivalent
2. Starting from a bivalent config., there is always another bivalent config. that is reachable



# Lemma 2

Some initial configuration is bivalent

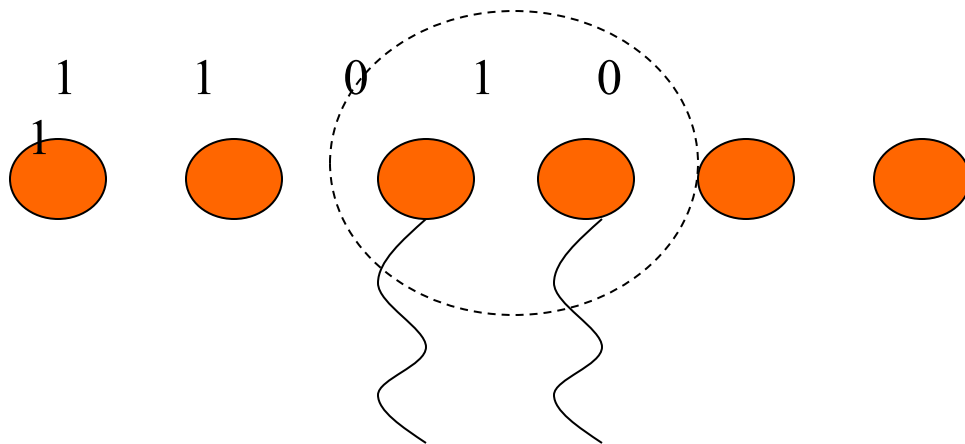
- Suppose all initial configurations were either 0-valent or 1-valent.
- If there are  $N$  processes, there are  $2^N$  possible initial configurations
- Place all configurations side-by-side (in a lattice), where adjacent configurations differ in initial xp value for exactly one process.



- There has to be **some** adjacent pair of 1-valent and 0-valent configs.



- There has to be **some** adjacent pair of 1-valent and 0-valent configs.
- Let the process  $p$ , that has a different state across these two configs., be the process that has crashed (i.e., is silent throughout)



Both initial configs. will lead to the same config. for the same sequence of events

Therefore, both these initial configs. are bivalent when there is such a failure



# What we'll show

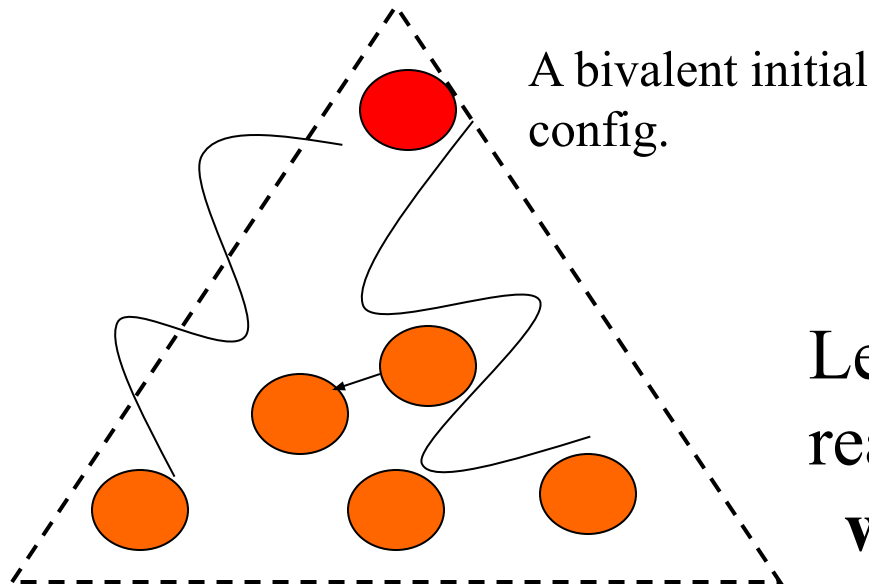
1. There exists an initial configuration that is bivalent
2. Starting from a bivalent config., there is always another bivalent config. that is reachable





**Starting from a bivalent config., there is always  
another bivalent config. that is reachable**

# Lemma 3



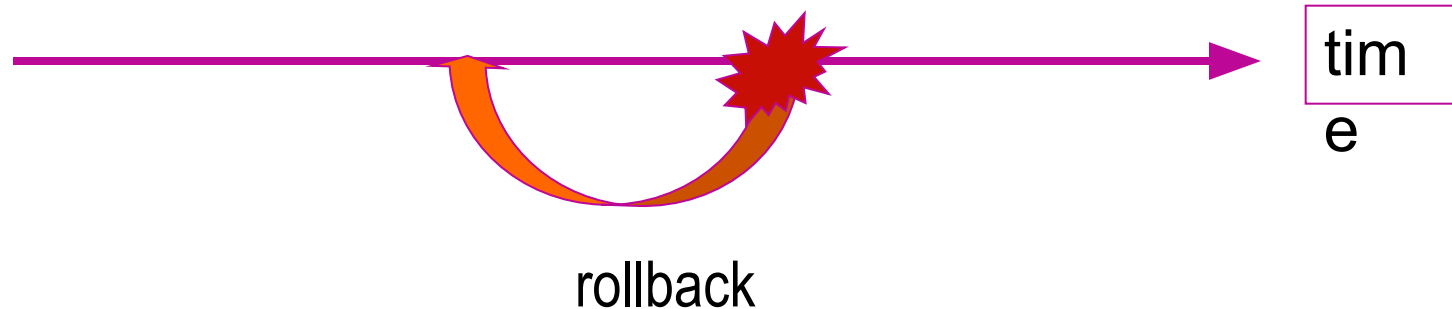
let  $e=(p,m)$  be some event  
applicable to the initial  
config.

Let  $C$  be the set of configs.  
reachable  
**without** applying  $e$

# Backward vs. forward error recovery

## *Backward error recovery*

When safety property is violated, the computation **rolls back** and resumes from a previous correct state.



## *Forward error recovery*

Computation does not care about getting the history right, but moves on, as long as eventually the safety property is restored. True for **self-stabilizing systems**.

# Message Logging



- Tolerate crash failures
- Each process periodically records its local state and log messages received after
  - Once a crashed process recovers, its state must be consistent with the states of other processes
  - Orphan processes
    - surviving processes whose states are inconsistent with the recovered state of a crashed process
  - Message Logging protocols guarantee that upon recovery no processes are orphan processes

# Message logging protocols



- Pessimistic Message Logging
  - avoid creation of orphans during execution
  - no process  $p$  sends a message  $m$  until it knows that all messages delivered before sending  $m$  are logged; quick recovery
  - Can block a process for each message it receives - slows down throughput
  - allows processes to communicate only from recoverable states; synchronously log to stable storage any information that may be needed for recovery before allowing process to communicate

# Message Logging



- Optimistic Message Logging
  - take appropriate actions during recovery to eliminate all orphans
  - Better performance during failure-free runs
  - allows processes to communicate from non-recoverable states; failures may cause these states to be permanently unrecoverable, forcing rollback of any process that depends on such states

# Causal Message Logging



- Causal Message Logging
  - no orphans when failures happen and do not block processes when failures do not occur.
  - Weaken condition imposed by pessimistic protocols
  - Allow possibility that the state from which a process communicates is unrecoverable because of a failure, but only if it does not affect consistency.
  - Append to all communication information needed to recover state from which communication originates - this is replicated in memory of processes that causally depend on the originating state.