# CS 425 / ECE 428
# Distributed Systems
# Fall 2014

Indranil Gupta (Indy)

*Lecture 4: Failure Detection and Membership*

# A Challenge

- You've been put in charge of a datacenter, and your manager has told you, "Oh no! We don't have any failures in our datacenter!"

- Do you believe him/her?

- What would be your first responsibility?

- Build a failure detector

- What are some things that could go wrong if you didn't do this?

# FAILURES ARE THE NORM

… not the exception, in datacenters.

Say, the rate of failure of one machine (OS/disk/motherboard/network, etc.) is once every 10 years (120 months) on average.

When you have 120 servers in the DC, the mean time to failure (MTTF) of the next machine is 1 month.

When you have 12,000 servers in the DC, the MTTF is about once every 7.2 hours!

Soft crashes and failures are even more frequent!
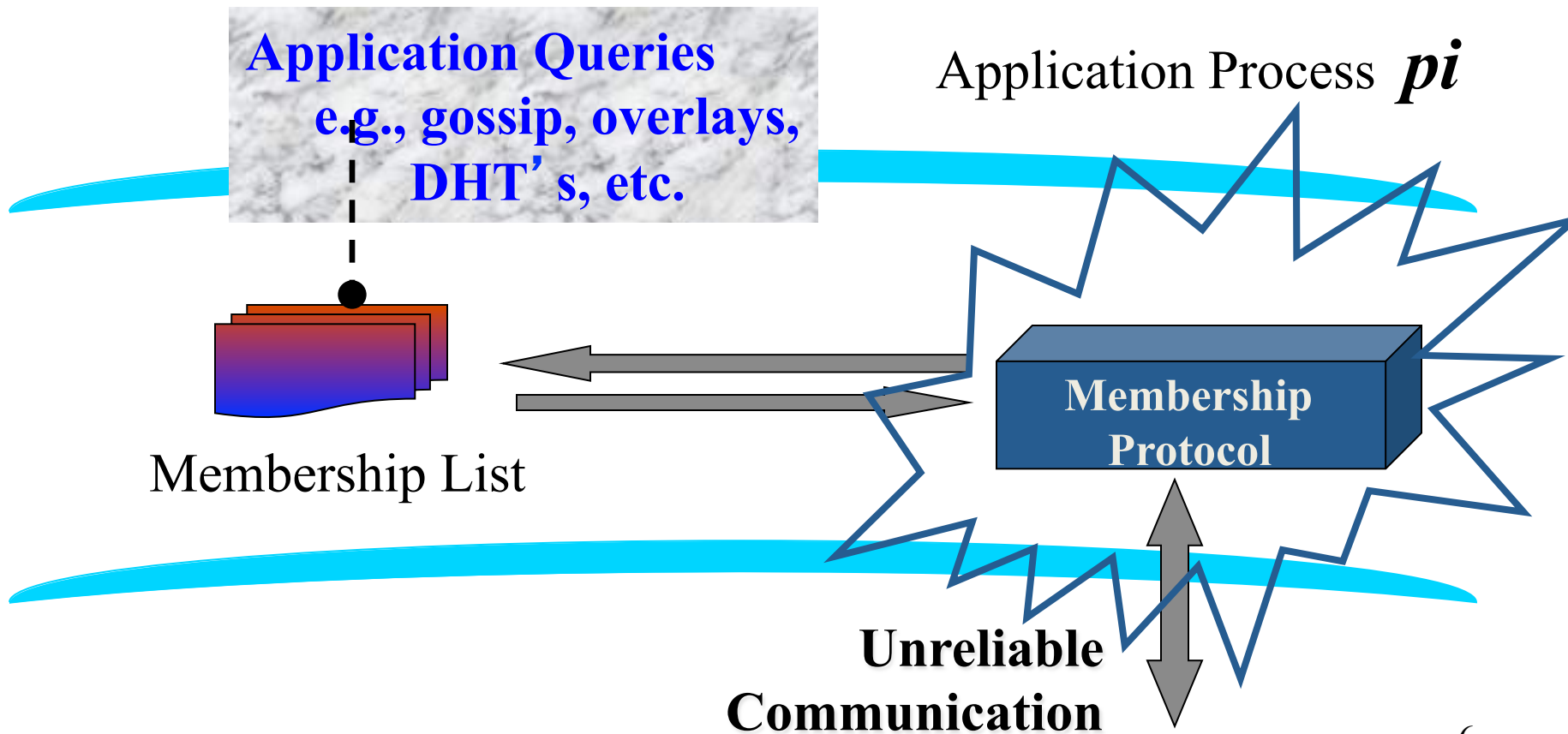
# TO BUILD A FAILURE DETECTOR

- You have a few options

  1. Hire 1000 people, each to monitor one machine in the datacenter and report to you when it fails.

  2. Write a failure detector program (distributed) that automatically detects failures and reports to your workstation.

# TARGET SETTINGS

- Process 'group'-based systems
  – Clouds/Datacenters
  – Replicated servers
  – Distributed databases


- Crash-stop/Fail-stop process failures
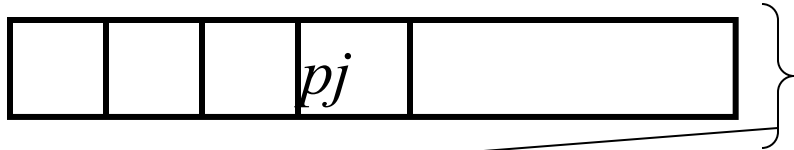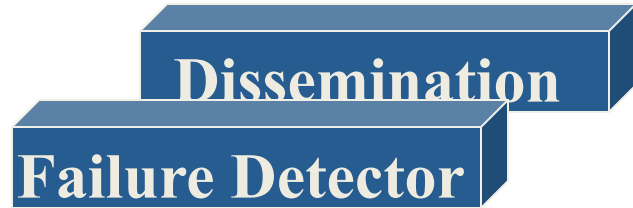
# GROUP MEMBERSHIP SERVICE



**Application Queries
e.g., gossip, overlays,
DHT's, etc.**

Application Process $pi$

Membership List

**Membership
Protocol**

**Unreliable
Communication**

# TWO SUB-PROTOCOLS

Application Process *pi*

*Group Membership List*

| | | | pj | |
|---|---|---|---|---|

- *Complete list all the time* **(Strongly consistent)**
  - Virtual synchrony
- *Almost-Complete* **list (Weakly consistent)**
  - Gossip-style, SWIM, …
- **Or** *Partial-random* **list (other systems)**
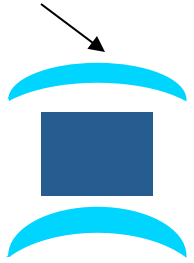  - SCAMP, T-MAN, Cyclon,…

**Focus of this series of lecture**

**Dissemination**

**Failure Detector**

**Unreliable Communication**

# LARGE GROUP: SCALABILITY A GOAL

*this is us (pi)*

Process Group "Members"

1000's of processes

**Unreliable Communication Network**

# GROUP MEMBERSHIP PROTOCOL



**II** **Failure Detector**

*Some* process finds out quickly

*pi*

**I** *pj crashed*

**III** **Dissemination**

**Unreliable Communication Network**

Crash-stop Failures only

9

# NEXT

- How do you design a group membership protocol?

# I. *pj* crashes

- Nothing we can do about it!
- A frequent occurrence
- Common case rather than exception
- Frequency goes up linearly with size of datacenter

# II. Distributed Failure Detectors: Desirable Properties

- Completeness = each failure is detected

- Accuracy = there is no mistaken detection

- Speed
  – Time to first detection of a failure

- Scale
  – Equal Load on each member
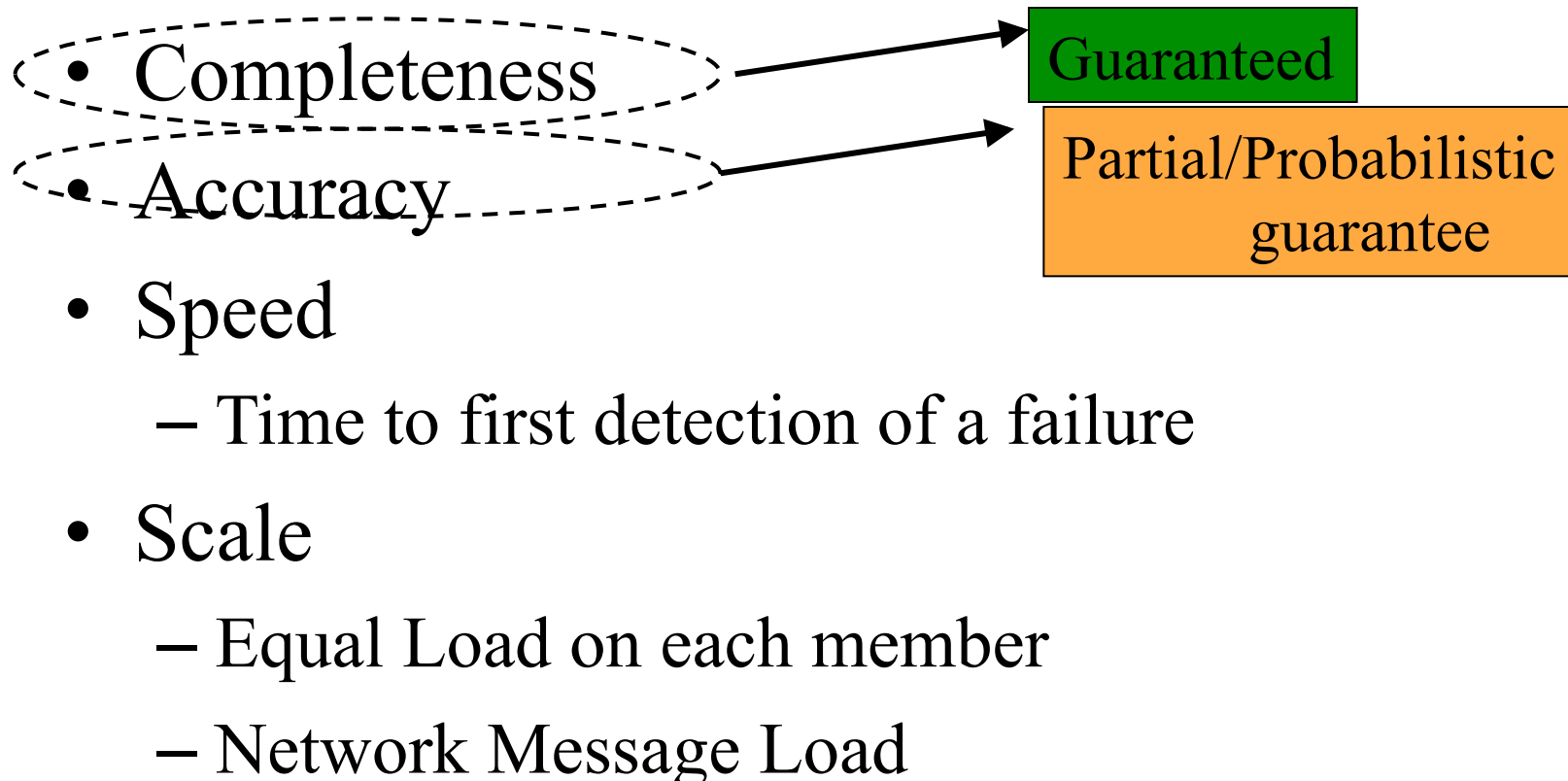  – Network Message Load

12

# DISTRIBUTED FAILURE DETECTORS: PROPERTIES

- Completeness

- Accuracy

- Speed
  - Time to first detection of a failure

- Scale
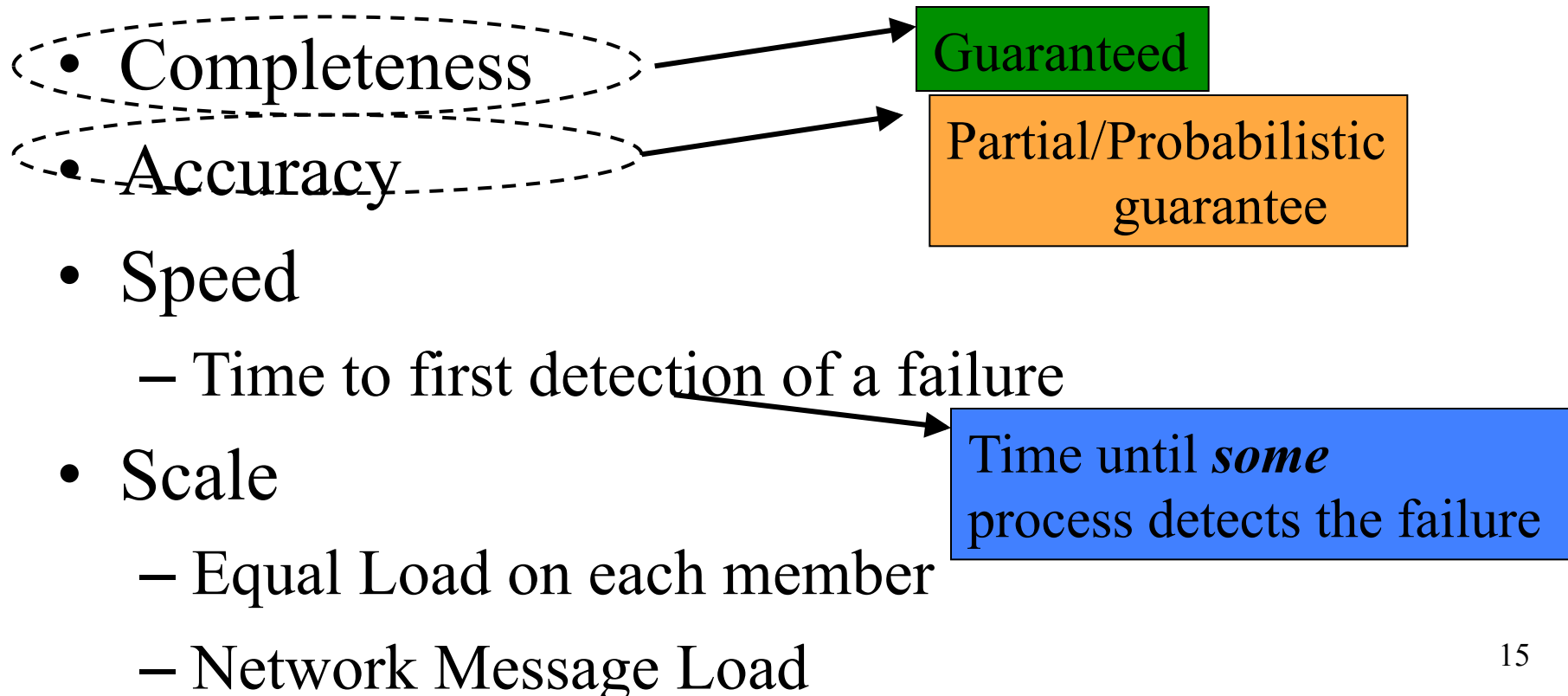  - Equal Load on each member
  - Network Message Load

Impossible together in lossy networks [Chandra and Toueg]

If possible, then can solve consensus!
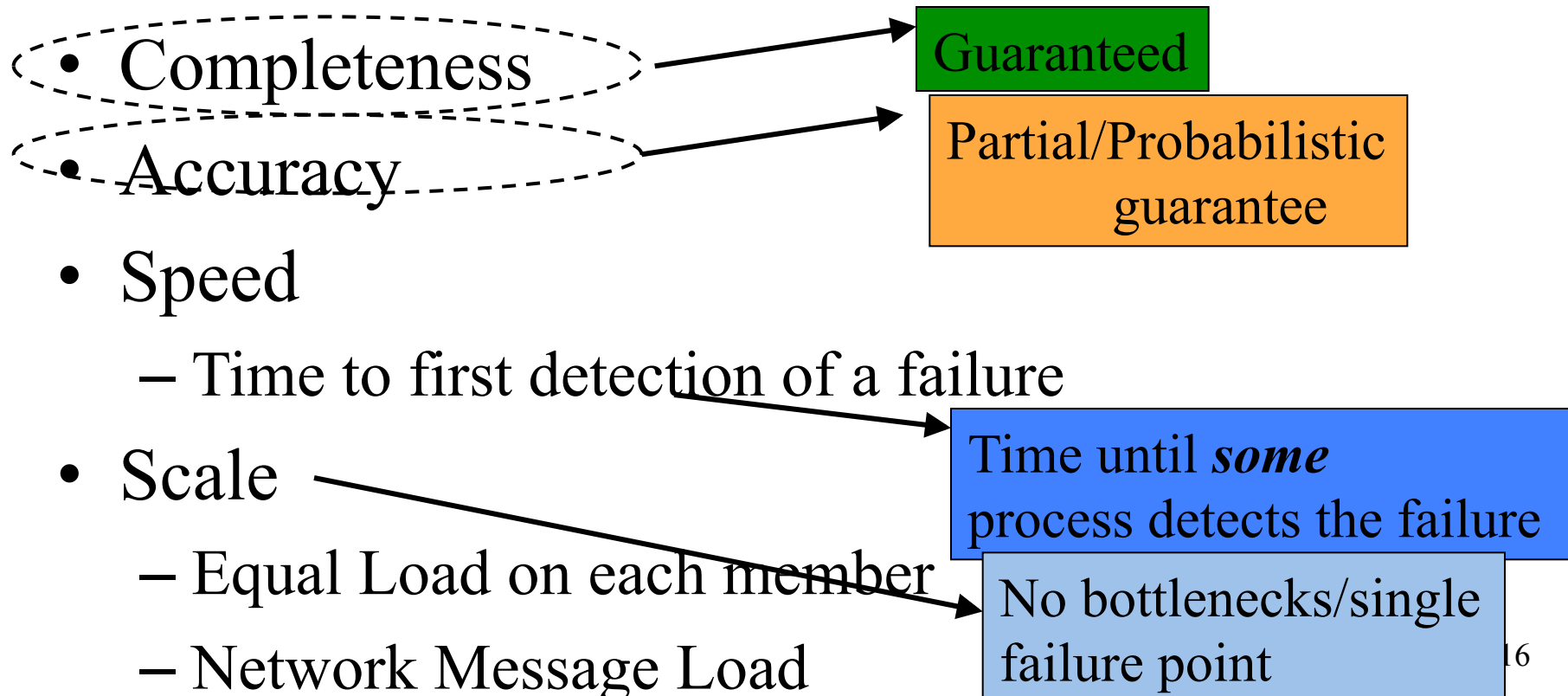
# WHAT REAL FAILURE DETECTORS PREFER

- Completeness → Guaranteed

- Accuracy → Partial/Probabilistic guarantee

- Speed
  - Time to first detection of a failure

- Scale
  - Equal Load on each member
  - Network Message Load

14

# WHAT REAL FAILURE DETECTORS PREFER

- Completeness ──→ Guaranteed
- Accuracy ──→ Partial/Probabilistic guarantee

- Speed
  - Time to first detection of a failure ──→ Time until *some* process detects the failure

- Scale
  - Equal Load on each member
  - Network Message Load

# What Real Failure Detectors Prefer

- • Completeness ⟶ **Guaranteed**

- • Accuracy ⟶ **Partial/Probabilistic guarantee**

- • Speed
  - – Time to first detection of a failure ⟶ **Time until *some* process detects the failure**

- • Scale
  - – Equal Load on each member ⟶ **No bottlenecks/single failure point**
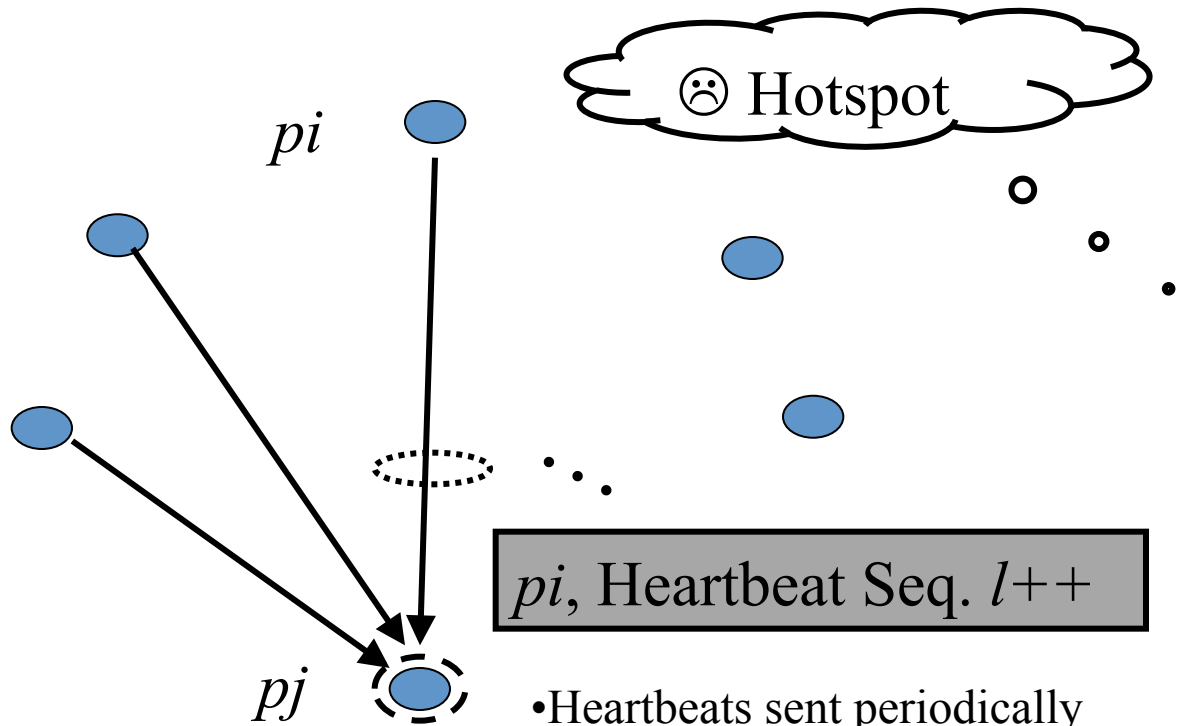  - – Network Message Load

16

# FAILURE DETECTOR PROPERTIES

- Completeness

- Accuracy

- Speed
  - Time to first detection of a failure

- Scale
  - Equal Load on each member
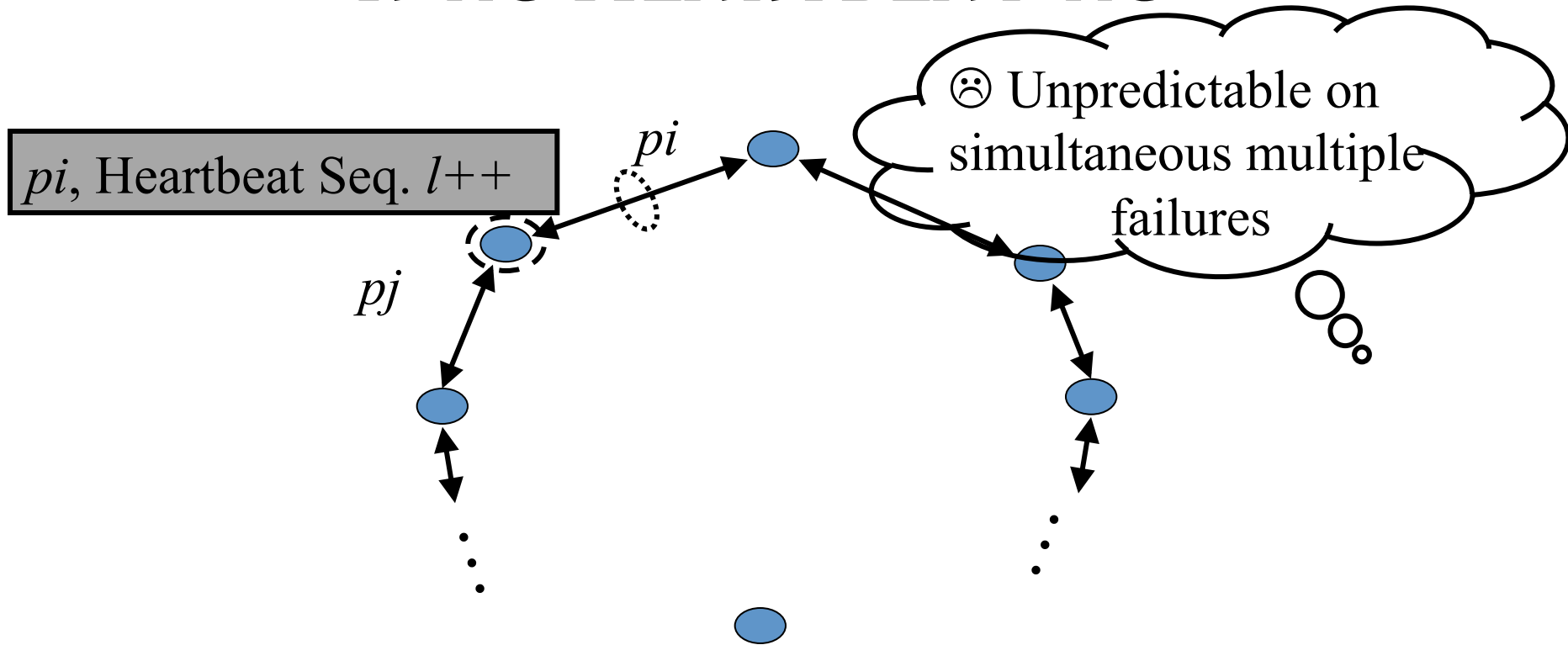  - Network Message Load

In spite of
arbitrary simultaneous
process failures
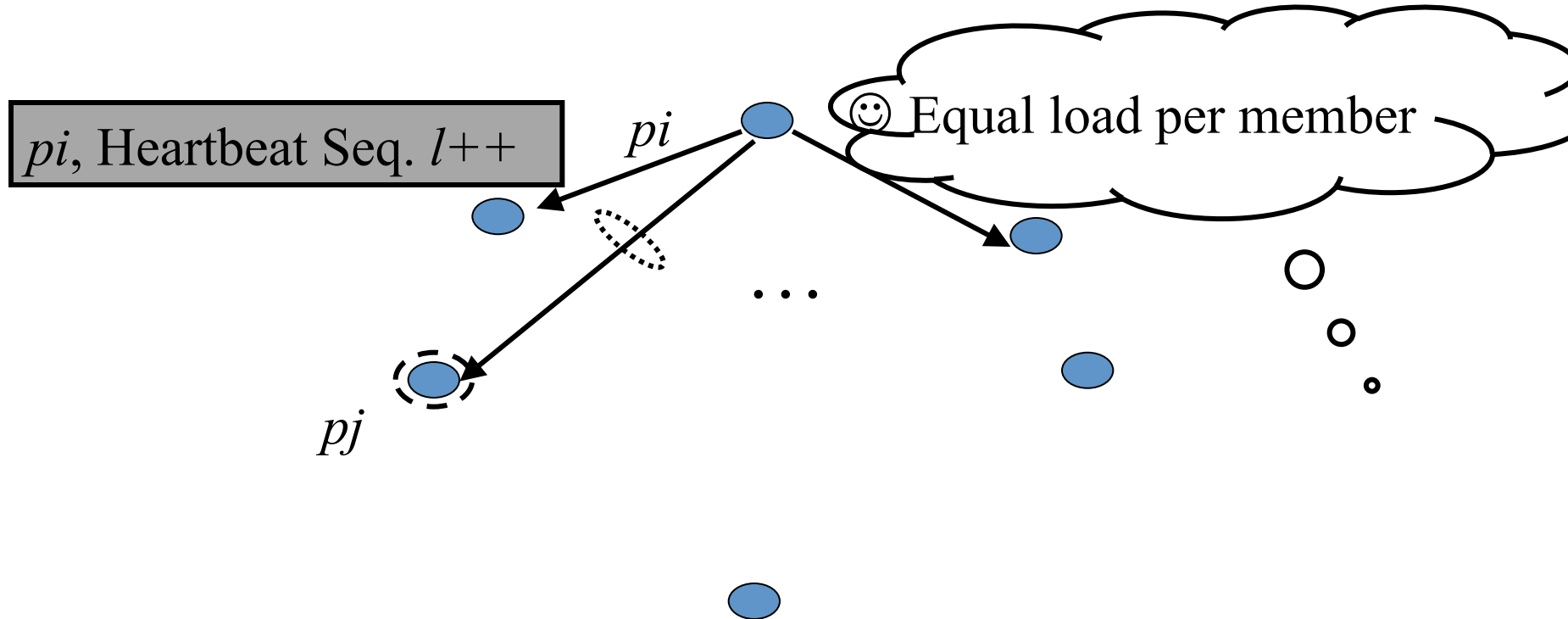
# CENTRALIZED HEARTBEATING



$pi$

☹ Hotspot

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

$pj$

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

18

# Ring Heartbeating



pi, Heartbeat Seq. $l$++

*pi*

*pj*

☹ Unpredictable on simultaneous multiple failures

# ALL-TO-ALL HEARTBEATING



*pi*, Heartbeat Seq. *l++*
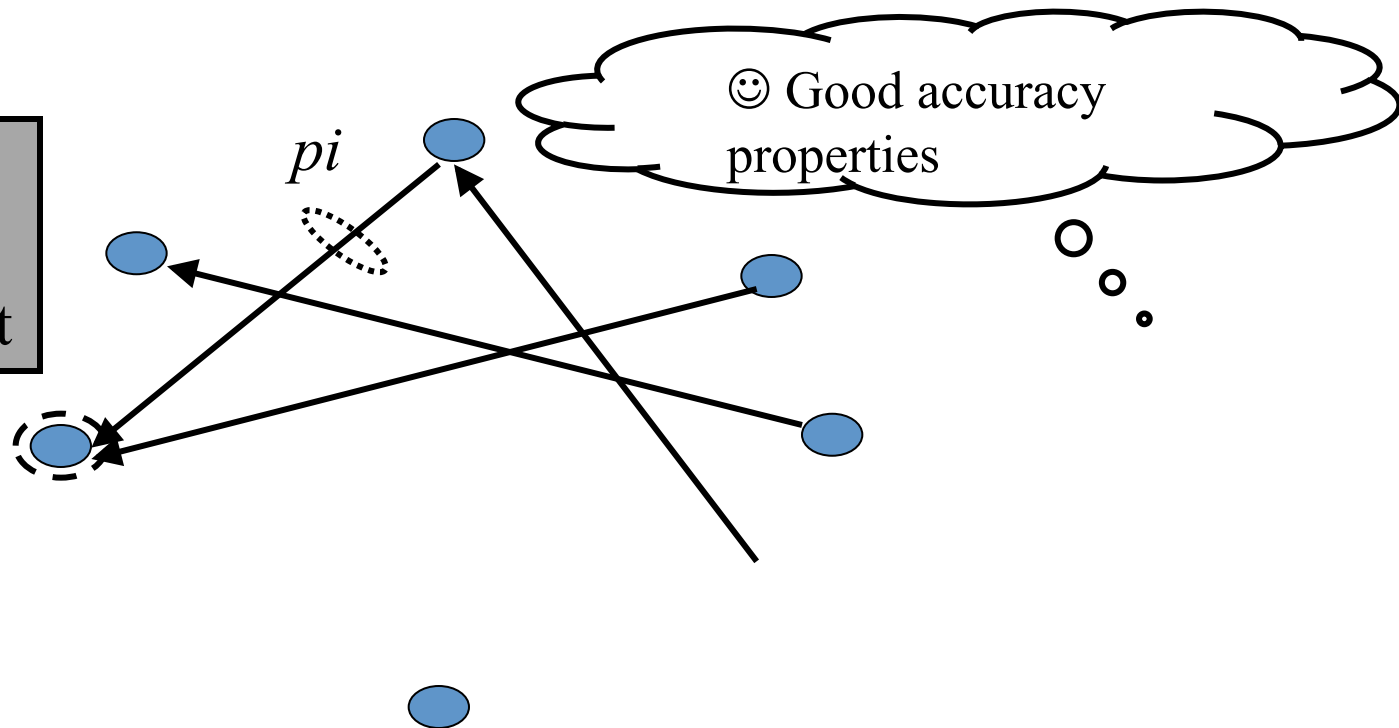
*pi*

☺ Equal load per member

*pj*

. . .

# Next

- How do we increase the robustness of all-to-all heartbeating?

# GOSSIP-STYLE HEARTBEATING

Array of
Heartbeat Seq. $l$
for member subset

$pi$

☺ Good accuracy
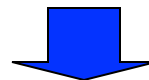properties

# GOSSIP-STYLE FAILURE DETECTION

| 1 | 10118 | 64 |
|---|---|---|
| 2 | 10110 | 64 |
| 3 | 10090 | 58 |
| 4 | 10111 | 65 |

| 1 | 10120 | 66 |
|---|---|---|
| 2 | 10103 | 62 |
| 3 | 10098 | 63 |
| 4 | 10111 | 65 |

Address        Time (local)

Heartbeat Counter

| 1 | 10120 | 70 |
|---|---|---|
| 2 | 10110 | 64 |
| 3 | 10098 | 70 |
| 4 | 10111 | 65 |

Current time : 70 at node 2

(asynchronous clocks)

**Protocol**:

• Nodes periodically gossip their membership list: pick random nodes, send it list

• On receipt, it is *merged* with local membership list

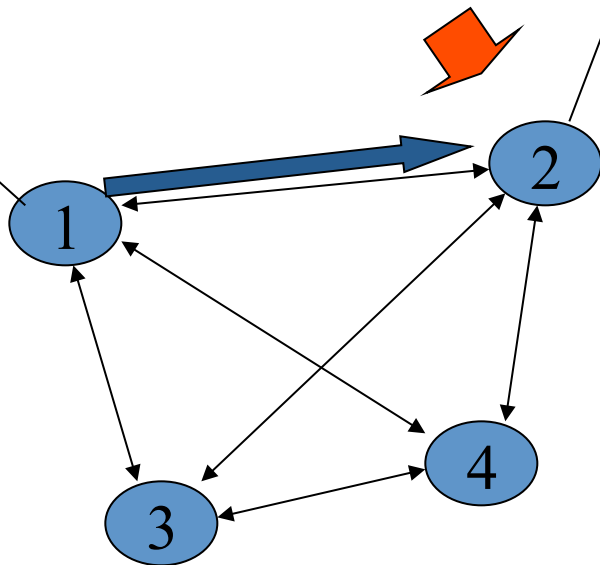• When an entry times out, member is marked as failed

23

# Gossip-Style Failure Detection

- If the heartbeat has not increased for more than $T_{fail}$ seconds,
  the member is considered failed

- And after $T_{cleanup}$ seconds, it will delete the member from the list

- Why two different timeouts?

# GOSSIP-STYLE FAILURE DETECTION♪

- What if an entry pointing to a failed node is deleted right after $T_{fail}$ (=24) seconds?

| 1 | 10120 | 66 |
|---|-------|-----|
| 2 | 10110 | 64 |
| 3 | 10098 | 75 |
| 4 | 10111 | 65 |

| 1 | 10120 | 66 |
|---|-------|-----|
| 2 | 10103 | 62 |
| 3 | 10098 | 55 |
| 4 | 10111 | 65 |

Current time : 75 at node 2

25

# MULTI-LEVEL GOSSIPING

•Network topology is hierarchical

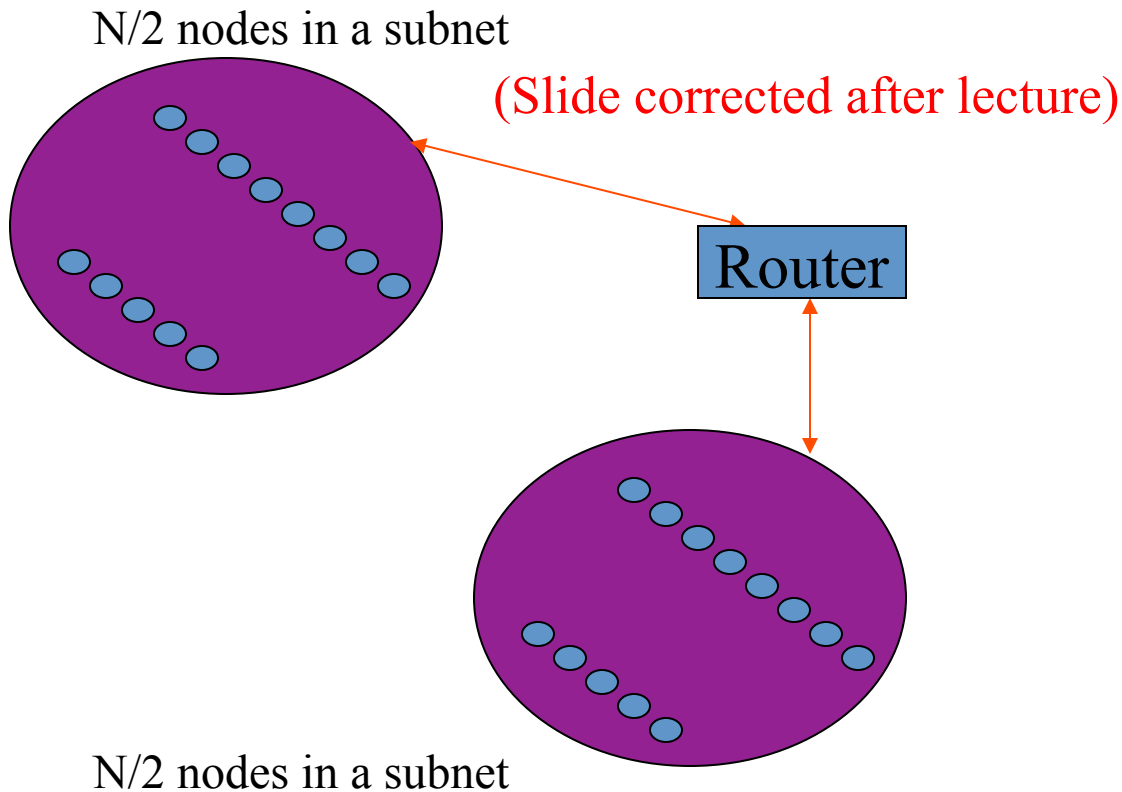•Random gossip target selection => core routers face O(N) load (Why?)

N/2 nodes in a subnet

(Slide corrected after lecture)

•Fix: In subnet $i$, which contains $n_i$ nodes, pick gossip target in your subnet with probability $(1-1/n_i)$

Router

•Router load=O(1)

•Dissemination time=O(log(N))

•What about latency for multi-level topologies?

N/2 nodes in a subnet

[Gupta et al, TPDS 06]

# ANALYSIS/DISCUSSION

- What happens if gossip period $T_{gossip}$ is decreased?
- A single heartbeat takes $O(\log(N))$ time to propagate. So: N heartbeats take:
  - $O(\log(N))$ time to propagate, if bandwidth allowed per node is allowed to be $O(N)$
  - $O(N.\log(N))$ time to propagate, if bandwidth allowed per node is only $O(1)$
  - What about $O(k)$ bandwidth?
- What happens to $P_{mistake}$ (false positive rate) as $T_{fail}$ , $T_{cleanup}$ is increased?
- Tradeoff: False positive rate vs. detection time vs. bandwidth

# Next

- So, is this the best we can do? What is the best we can do?

# Failure Detector Properties ...

- Completeness

- Accuracy

- Speed
  - Time to first detection of a failure

- Scale
  - Equal Load on each member
  - Network Message Load
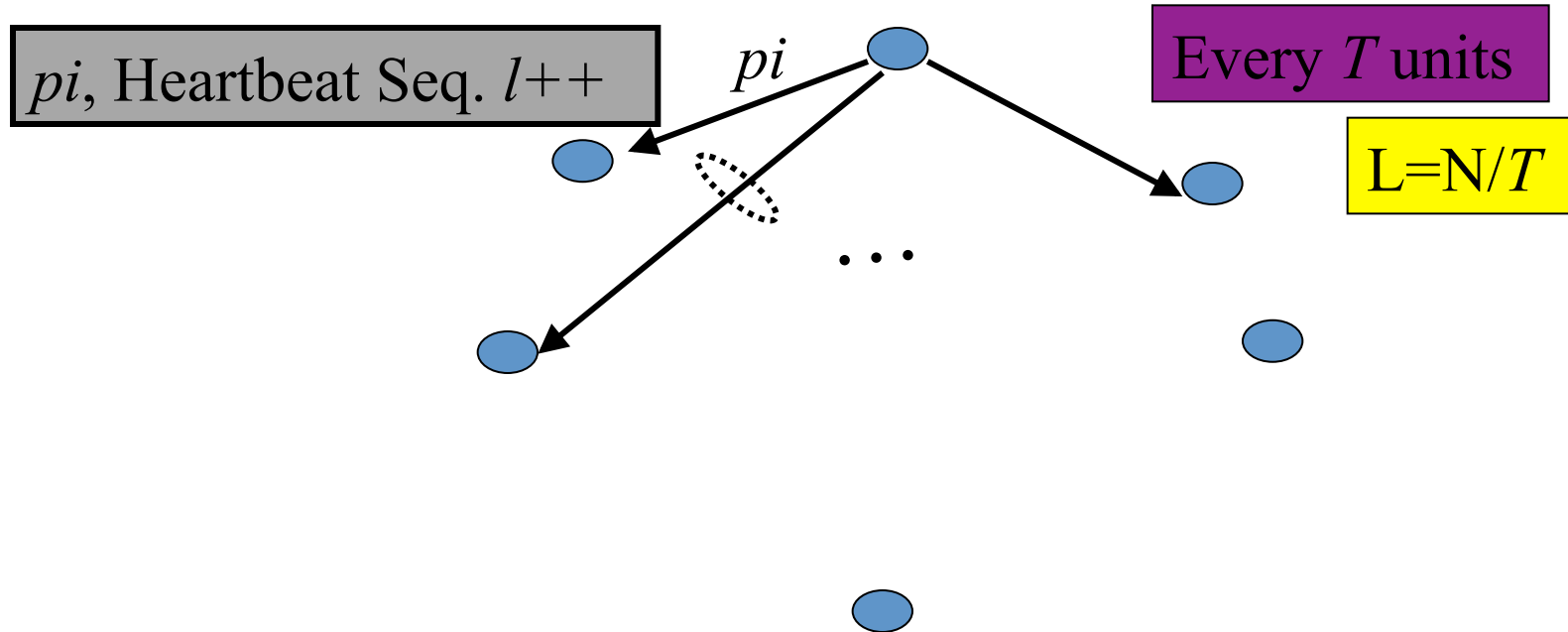
# ...Are application-defined Requirements

- Completeness → Guarantee always
- Accuracy → Probability *PM(T)*
- Speed → *T* time units
  - Time to first detection of a failure
- Scale
  - Equal Load on each member
  - Network Message Load

# ...ARE APPLICATION-DEFINED REQUIREMENTS

- Completeness → Guarantee always
- Accuracy → Probability *PM(T)*
- Speed → *T* time units
  - Time to first detection of a failure

N*L: Compare this across protocols

- Scale
  - Equal Load on each member
  - Network Message Load

# ALL-TO-ALL HEARTBEATING

*pi*, Heartbeat Seq. *l*++

*pi*

Every *T* units

L=N/*T*

...
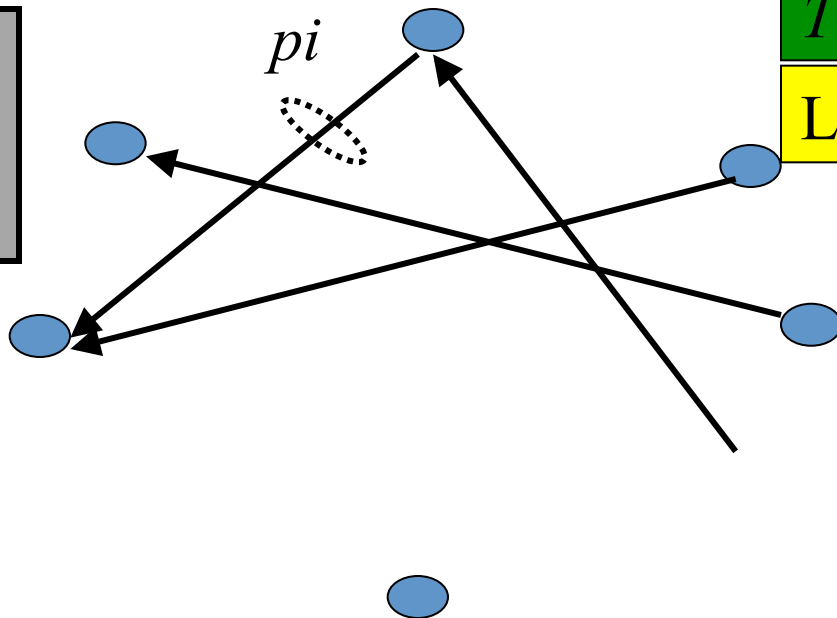
# GOSSIP-STYLE HEARTBEATING

Array of Heartbeat Seq. $l$ for member subset

Every tg units =gossip period, send O(N) gossip message

$pi$

$T$=logN * tg

L=N/tg=N*logN/$T$

# WHAT'S THE BEST/OPTIMAL WE CAN DO? <span style="color:red">Slide changed after lecture</span>

- *Worst case* load L* <span style="color:red">per member</span> in the group (messages per second)
  - as a function of *T, PM(T)*, N
  - Independent Message Loss probability $p_{ml}$

- $$L^* = \frac{\log(PM(T))}{\log(p_{ml})} \cdot \frac{1}{T}$$

# Heartbeating

- Optimal L is independent of N (!)
- All-to-all and gossip-based: sub-optimal
  - L=O(N/T)
  - try to achieve simultaneous detection at **all** processes
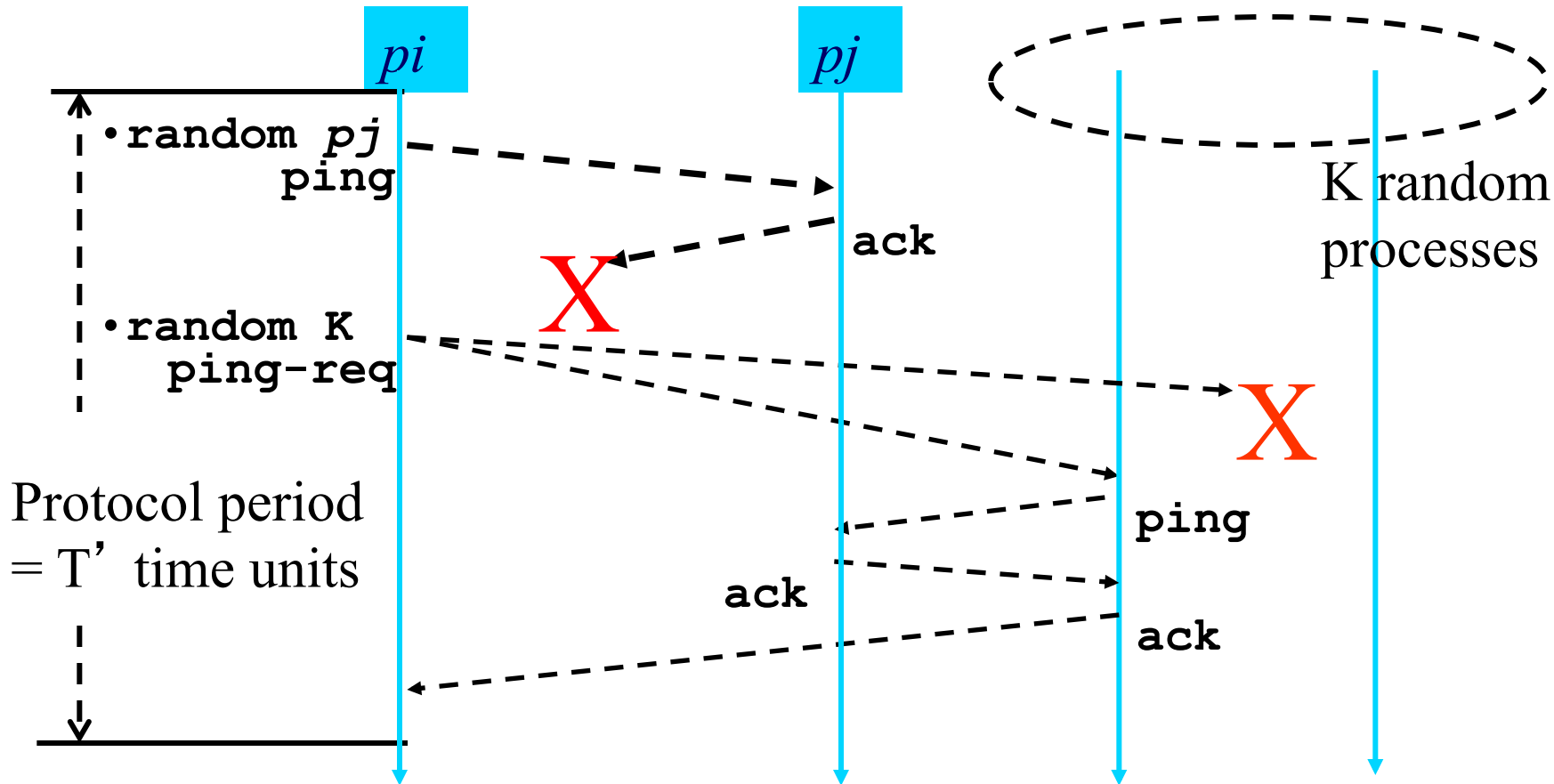  - fail to distinguish *Failure Detection* and *Dissemination* components

➲Key:
  ▫ Separate the two components
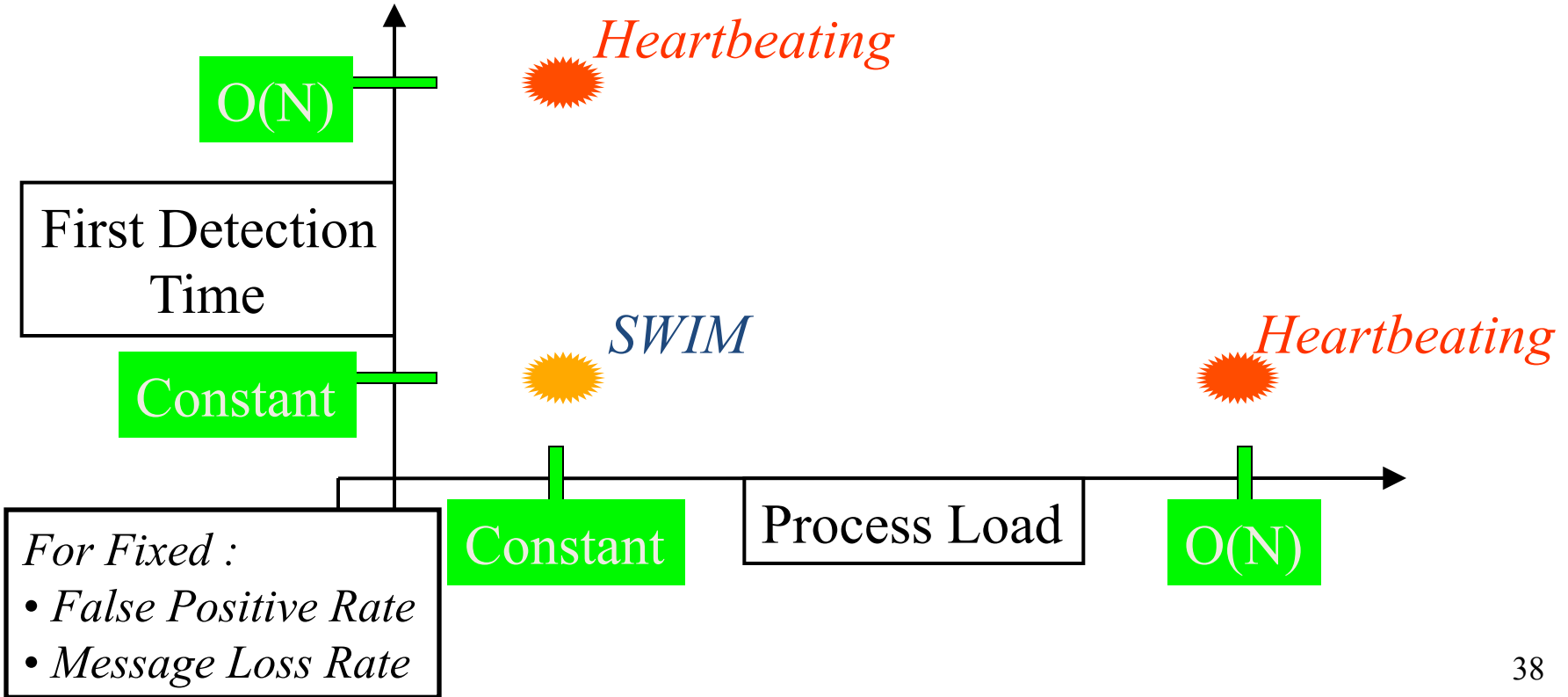  ▫ Use a non heartbeat-based Failure Detection Component

# Next

- Is there a better failure detector?

# SWIM FAILURE DETECTOR PROTOCOL

*pi*    *pj*

K random processes

- **random *pj* ping**

**ack**

X

- **random K ping-req**

X

**ping**

Protocol period = T' time units

**ack**

**ack**

# SWIM versus Heartbeating

# SWIM FAILURE DETECTOR

| Parameter | SWIM |
|---|---|
| First Detection Time | • Expected $\left\lceil \dfrac{e}{e-1} \right\rceil$ periods<br>• Constant (independent of group size) |
| Process Load | • Constant per period<br>• < 8 L* for 15% loss |
| False Positive Rate | • Tunable (via K)<br>• Falls exponentially as load is scaled |
| Completeness | • Deterministic time-bounded<br>• Within O(log(N)) periods w.h.p. |

# Accuracy, Load

- *PM(T)* is exponential in -*K*. Also depends on *pml* (and *pf* )
  - See paper

- $\dfrac{L}{L*} < 28$   $\dfrac{E[L]}{L*} < 8$   for up to 15 % loss rates

# DETECTION TIME

- Prob. of being pinged in T'= $1-(1-\dfrac{1}{N})^{N-1}=1-e^{-1}$

- $\mathrm{E}[T]=\mathrm{T'}.\dfrac{e}{e-1}$

- Completeness: *Any* alive member detects failure
  - Eventually
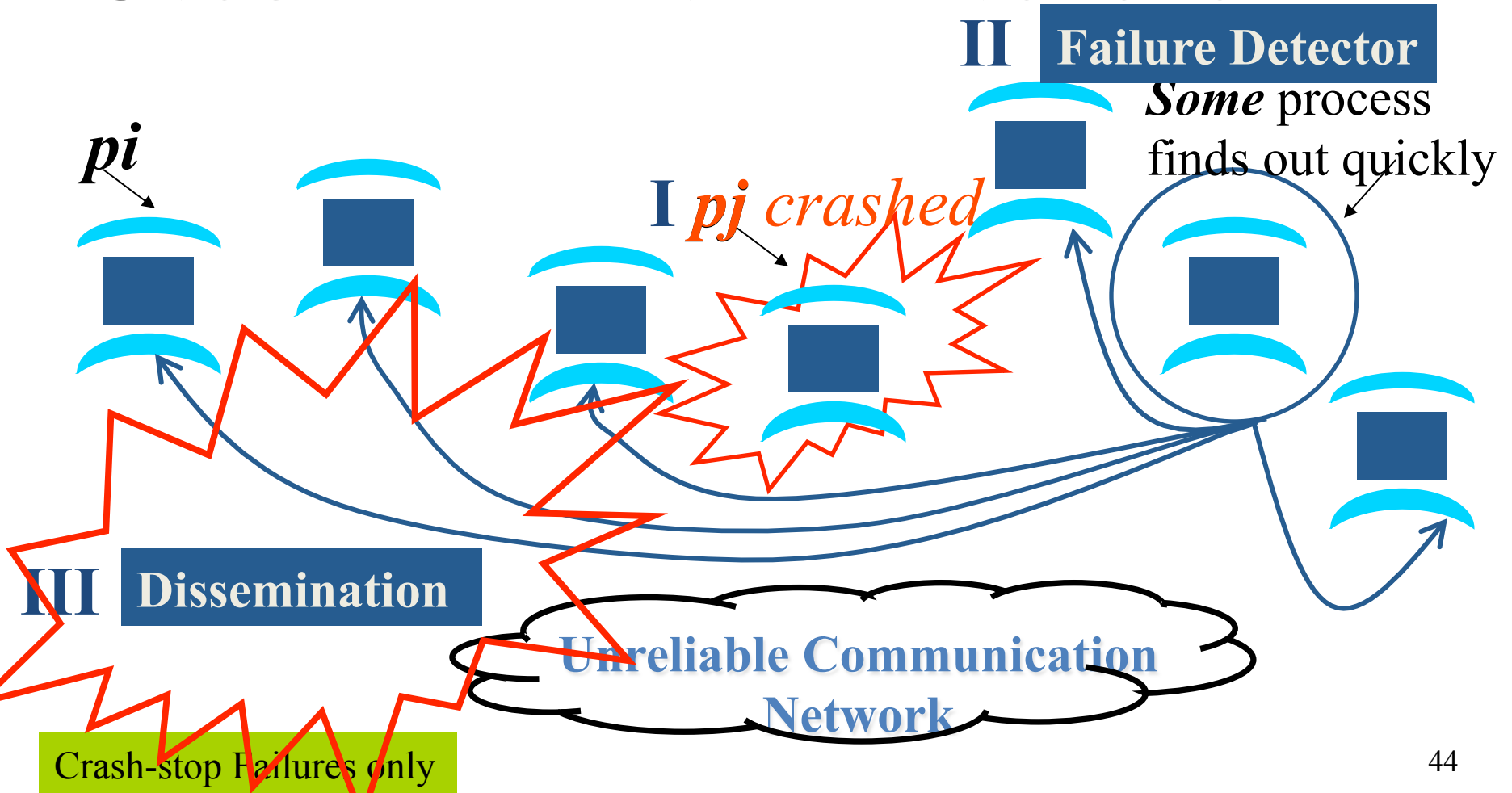  - By using a trick: within worst case *O(N)* protocol periods

41

# TIME-BOUNDED COMPLETENESS

- Key: select each membership element once as a ping target in a traversal
  - Round-robin pinging
  - Random permutation of list after each traversal
- Each failure is detected in worst case 2N-1 (local) protocol periods
- Preserves FD properties

42

# NEXT

- How do failure detectors fit into the big picture of a group membership protocol?
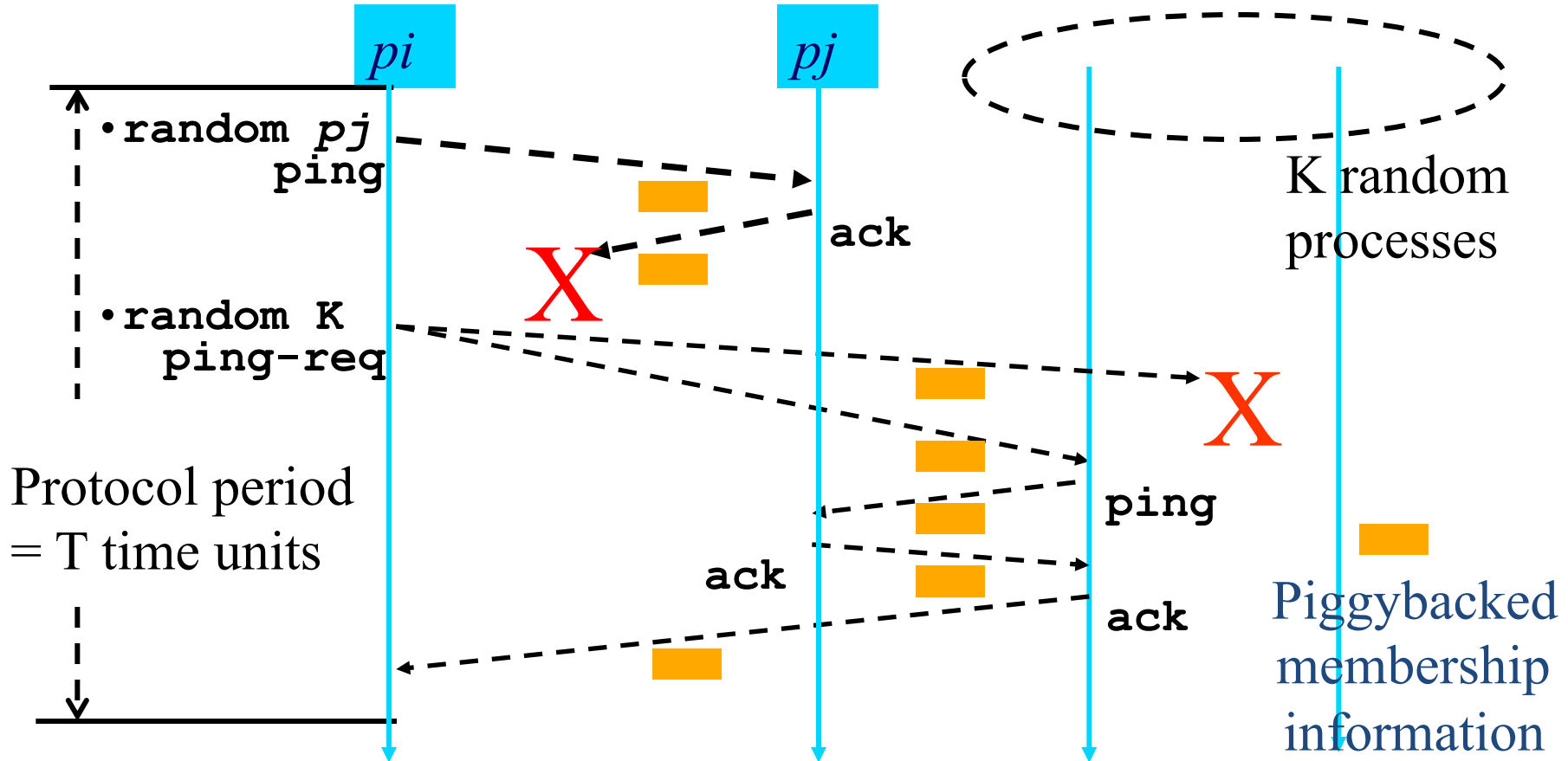
- What are the missing blocks?

# GROUP MEMBERSHIP PROTOCOL

**II** **Failure Detector**

*Some* process finds out quickly

*pi*

**I** *pj crashed*

**III** **Dissemination**

**Unreliable Communication Network**

Crash-stop Failures only

44

# DISSEMINATION OPTIONS

- Multicast (Hardware / IP)
  - unreliable
  - multiple simultaneous multicasts
- Point-to-point (TCP / UDP)
  - expensive
- Zero extra messages: Piggyback on Failure Detector messages
  - Infection-style Dissemination
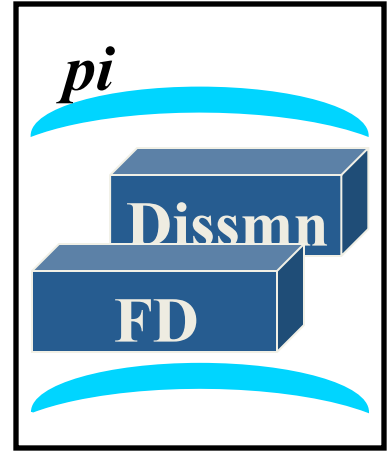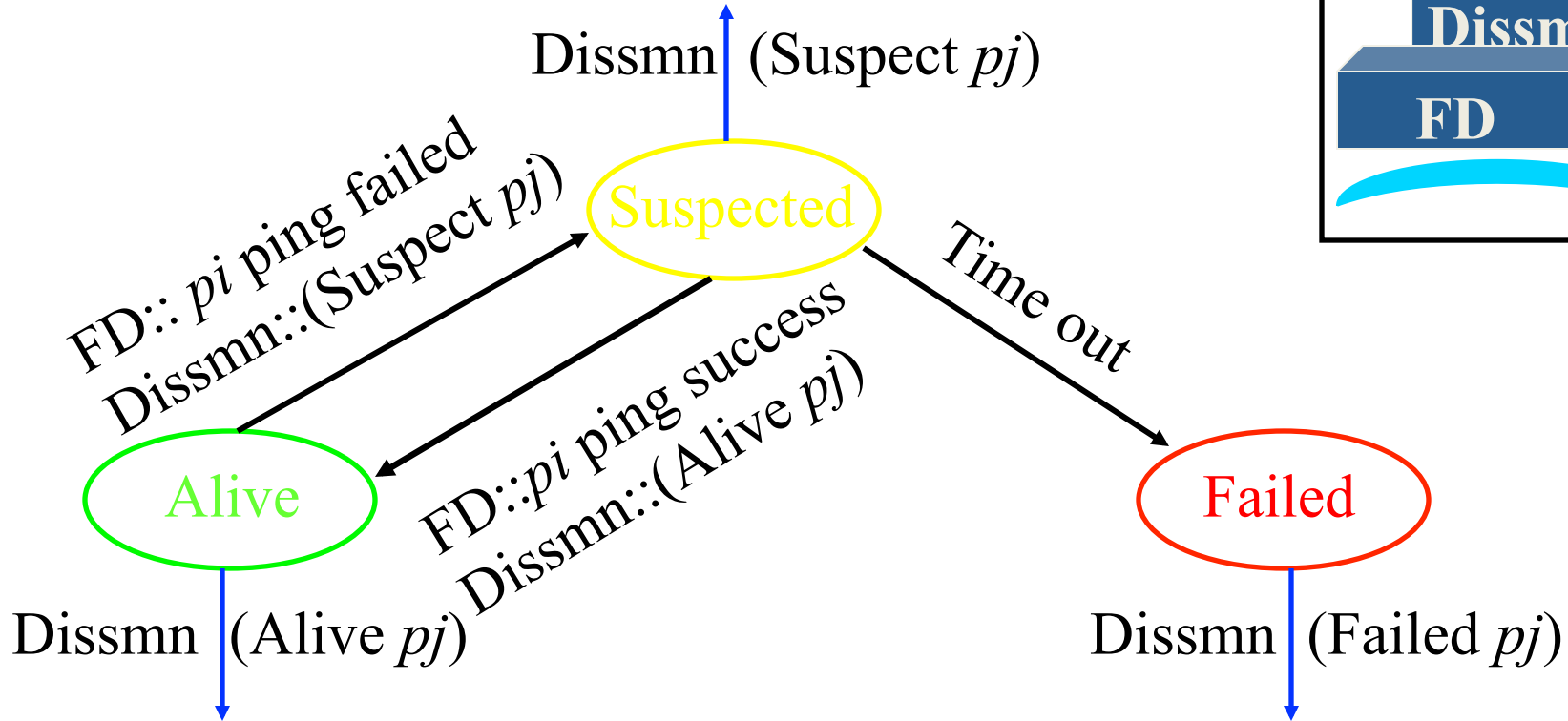
# INFECTION-STYLE DISSEMINATION

*pi*

*pj*

- **random *pj***
  **ping**

**ack**

X

- **random K**
  **ping-req**

X

K random
processes

Protocol period
= T time units

**ping**

**ack**

**ack**

Piggybacked
membership
information

# INFECTION-STYLE DISSEMINATION

- Epidemic/Gossip style dissemination
  - After $\lambda.\log(N)$ protocol periods, $N^{-(2\lambda-2)}$ processes would not have heard about an update

- Maintain a buffer of recently joined/evicted processes
  - Piggyback from this buffer
  - Prefer recent updates

- Buffer elements are garbage collected after a while
  - After $\lambda.\log(N)$ protocol periods, i.e., once they've propagated through the system; this defines weak consistency

# Suspicion Mechanism

- False detections, due to
  - Perturbed processes
  - Packet losses, e.g., from congestion
- Indirect pinging may not solve the problem
- Key: *suspect* a process before *declaring* it as failed in the group

# SUSPICION MECHANISM

# SUSPICION MECHANISM

- Distinguish multiple suspicions of a process
  - Per-process *incarnation number*
  - *Inc* # for *pi* can be incremented only by *pi*
    - e.g., when it receives a (Suspect, *pi*) message
  - Somewhat similar to DSDV
- Higher inc# notifications over-ride lower inc#'s
- Within an inc#: (Suspect inc #) > (Alive, inc #)
- (Failed, inc #) overrides everything else

# WRAP UP

- Failures the norm, not the exception in datacenters
- Every distributed system uses a failure detector
- Many distributed systems use a membership service

- Ring failure detection underlies
    - IBM SP2 and many other similar clusters/machines

- Gossip-style failure detection underlies
    - Amazon EC2/S3 (rumored!)

# IMPORTANT ANNOUNCEMENT

- Next week Tue and Thu: We'll have a **flipped classroom**! (like Khan Academy)
- Homework **before** Next week
    - Please see video lectures for two topics
        - Timestamps and Ordering before Tue
        - Global Snapshots before Thu
- When you come to class on Sep 9th (Tue) and Sep 11th (Thu) the TAs will be helping you do **exercises in class** (not HW problems, but other exercise problems we will give you)
- We will not replay videos in class, i.e., there will be no lecturing.
- If you don't see the videos **before** class, you will flounder in class. So make sure you see them before class.
- Exercises *may* count for grades.
- Please bring a pen/pencil and paper to both classes.