

# CS 525 Advanced Distributed Systems Spring 2010

Indranil Gupta (Indy)

Failure Detectors and Membership  
Protocols

March 18, 2010

1

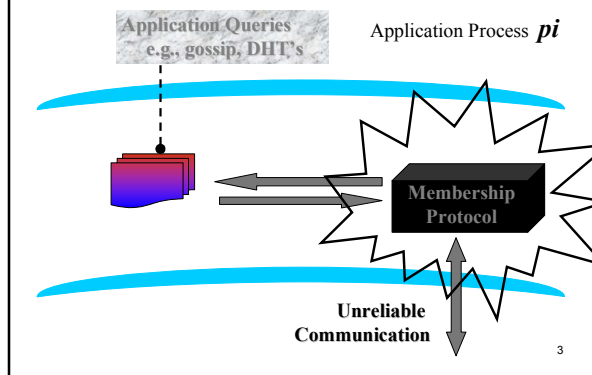
All Slides © IG

## Target Settings

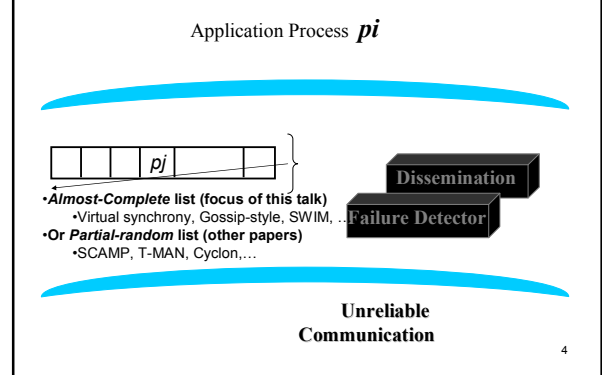
- Process 'group'-based systems
  - Clouds/Datacenters
  - Replicated servers
  - Distributed databases
- Crash-stop/Fail-stop process failures

2

## Group Membership Service

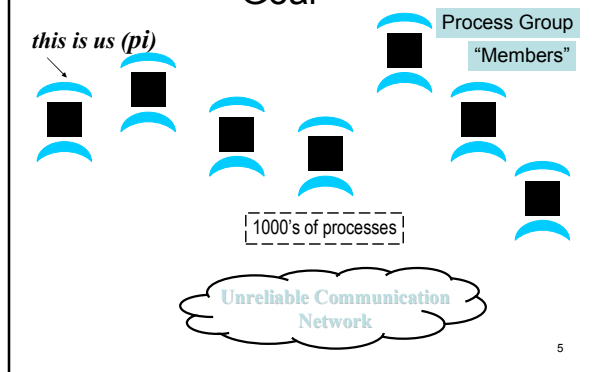


## Two sub-protocols

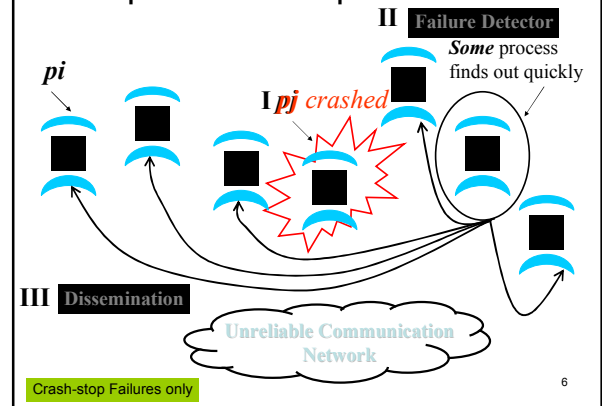


- **Almost-Complete list (focus of this talk)**
  - Virtual synchrony, Gossip-style, SWIM, ...
- **Or Partial-random list (other papers)**
  - SCAMP, T-MAN, Cyclon, ...

## Large Group: Scalability A Goal



## Group Membership Protocol



Crash-stop Failures only

## I. *pj* crashes

- Nothing we can do about it!
- A frequent occurrence
- Common case rather than exception

7

## II. Distributed Failure Detectors: 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

8

## 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]  
Can then solve consensus!

9

## 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

10

## Failure Detector Properties

- **Completeness** → Guaranteed
  - **Accuracy** → Partial/Probabilistic guarantee
  - **Speed**
    - Time to first detection of a failure
  - **Scale**
    - Equal Load on each member
    - Network Message Load
- Time until **some** process detects the failure

11

## Failure Detector Properties

- **Completeness** → Guaranteed
  - **Accuracy** → Partial/Probabilistic guarantee
  - **Speed**
    - Time to first detection of a failure
  - **Scale**
    - Equal Load on each member
    - Network Message Load
- Time until **some** process detects the failure
- No bottlenecks/single failure point

12

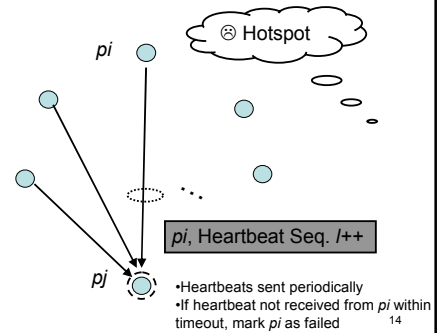
## Failure Detector Properties

- 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

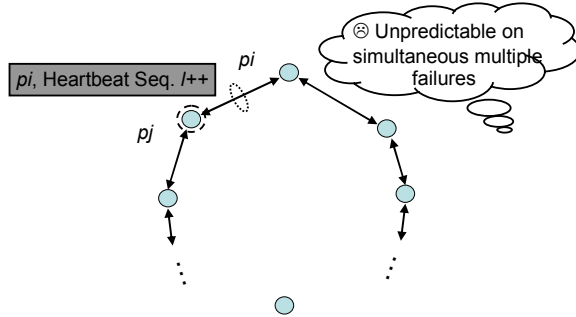
13

## Centralized Heartbeating



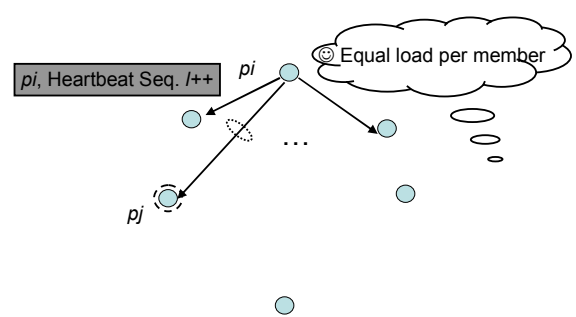
14

## Ring Heartbeating



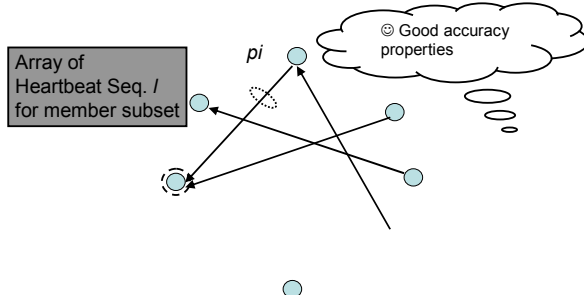
15

## All-to-All Heartbeating



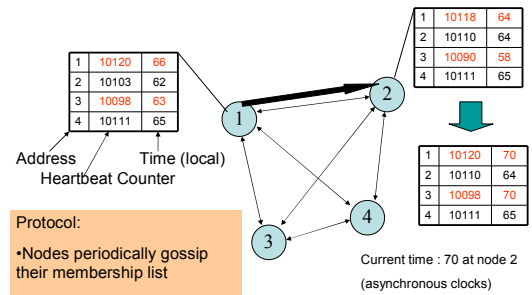
16

## Gossip-style Heartbeating



17

## Gossip-Style Failure Detection



18

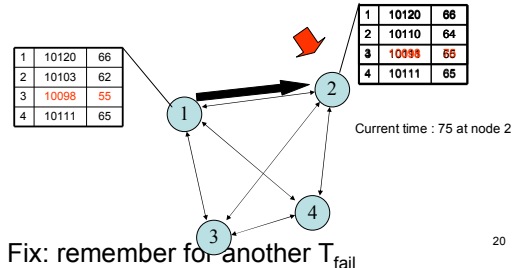
## Gossip-Style Failure Detection

- 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?

19

## Gossip-Style Failure Detection

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



20

## Multi-level Gossiping

• Network topology is hierarchical

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

• Fix: Select gossip target in subnet I, which contains  $n_i$  nodes, with probability  $1/n_i$

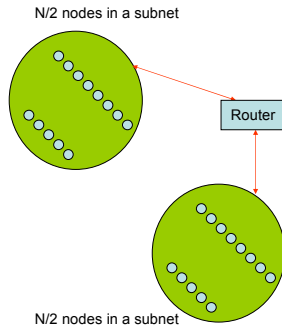
• Router load =  $O(1)$

• Dissemination time =  $O(\log(N))$

• Why?

• What about latency for multi-level topologies?

[Gupta et al, TPDS 06]



21

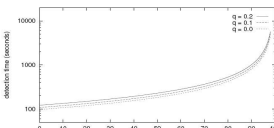
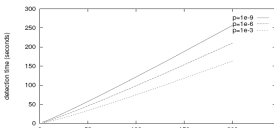
## Analysis/Discussion

- 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 are allowed to be  $O(N)$
  - $O(N \cdot \log(N))$  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

22

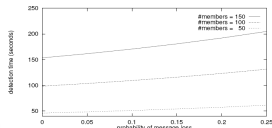
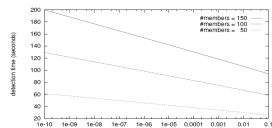
## Simulations

- As # members increases, the detection time increases



- As # failed members increases, the detection time increases significantly

- As requirement is loosened, the detection time decreases



- The algorithm is resilient to message loss

23

## Failure Detector Properties ...

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

24

### ...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

25

### ...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

**N\*L: Compare this across protocols**

26

### All-to-All Heartbeating

$p_i$ , Heartbeat Seq.  $l++$

Every  $T$  units

$L=N/T$

27

### Gossip-style Heartbeating

Array of Heartbeat Seq.  $l$  for member subset

$p_i$

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

$T=\log N * t_g$

$L=N/t_g=N*\log N/T$

28

### What's the Best/Optimal we can do?

- **Worst case load  $L^*$** 
  - as a function of  $T, PM(T), N$
  - Independent Message Loss probability  $p_m$
- $L^* = \frac{\log(PM(T))}{\log(p_m)} \cdot \frac{1}{T}$  (proof in PODC 01 paper)

29

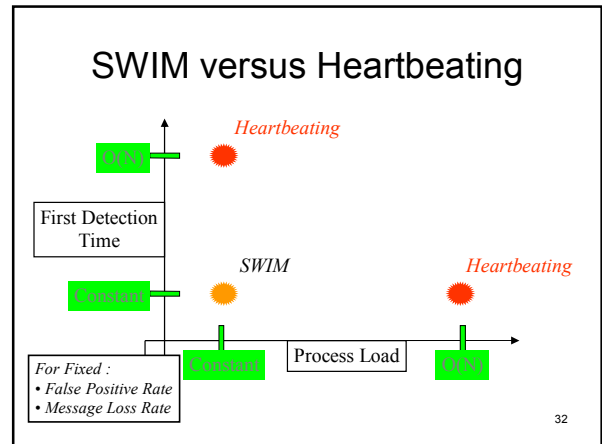
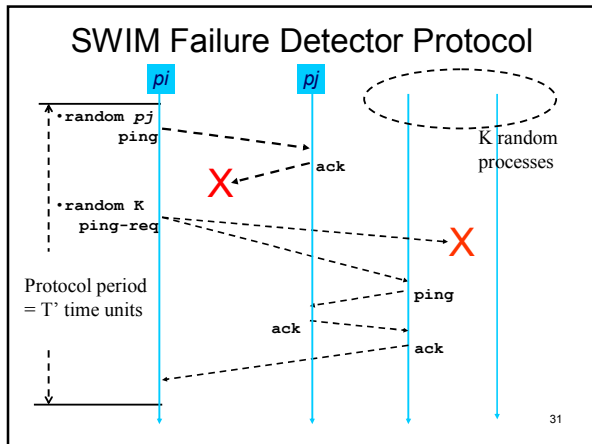
### Heartbeating

- 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

30



### SWIM Failure Detector

Parameter	SWIM
First Detection Time	<ul style="list-style-type: none"> <li>Expected <math>\left[ \frac{e}{e-1} \right]</math> periods</li> <li>Constant (independent of group size)</li> </ul>
Process Load	<ul style="list-style-type: none"> <li>Constant per period</li> <li>&lt; <math>8 L^*</math> for 15% loss</li> </ul>
False Positive Rate	<ul style="list-style-type: none"> <li>Tunable</li> <li>Falls exponentially as load is scaled</li> </ul>
Completeness	<ul style="list-style-type: none"> <li>Deterministic time-bounded</li> <li>Within <math>O(\log(N))</math> periods w.h.p.</li> </ul>

33

### Accuracy, Load

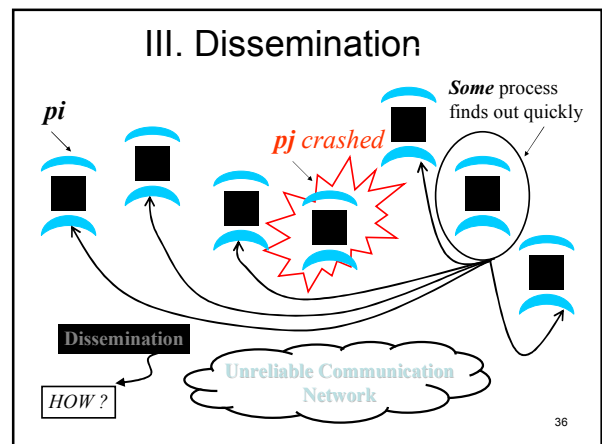
- $PM(T)$  is exponential in  $K$ . Also depends on  $pml$  (and  $pf$ )
  - See paper
- $\frac{L}{L^*} < 28$       $\frac{E[L]}{L^*} < 8$      for up to 15 % loss rates

34

### Detection Time

- Prob. of being pinged in  $T^* = 1 - (1 - \frac{1}{N})^{N-1} = 1 - e^{-1}$
- $E[T] = T^* \frac{e}{e-1}$
- Completeness: Any alive member detects failure
  - Within worst case  $O(N)$  protocol periods

35

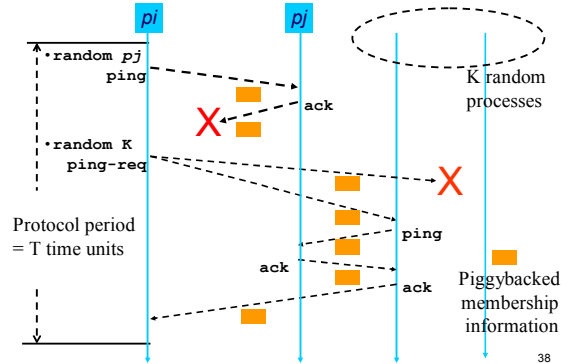


## Dissemination Options

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

37

## Infection-style Dissemination



38

## Infection-style Dissemination

- Epidemic style dissemination
  - After  $\lambda \cdot \log(N)$  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  $\lambda \cdot \log(N)$  protocol periods; this defines weak consistency

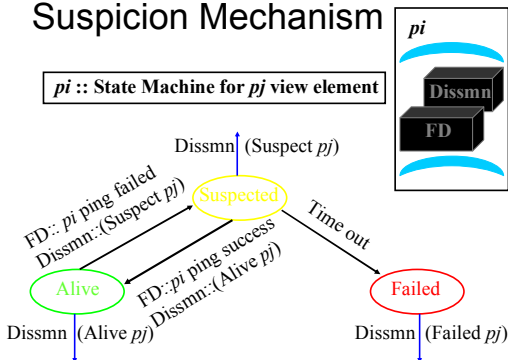
39

## Suspicion Mechanism

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

40

## Suspicion Mechanism



41

## 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
- Precedence rules for (Alive, inc #), (Suspect inc #), (Failed, inc #)
  - See paper

42

## Time-bounded Completeness

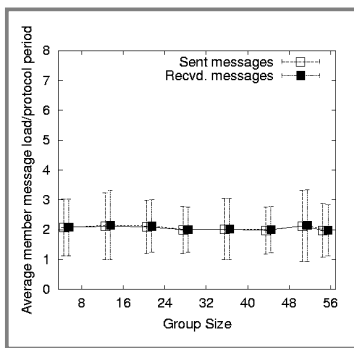
- Key: select each membership element once as a ping target in a traversal
  - Round-robin ping
  - Random permutation of list after each traversal
- Each failure is detected in worst case  $2N-1$  (local) protocol periods
- Preserves FD properties

43

## Results from an Implementation

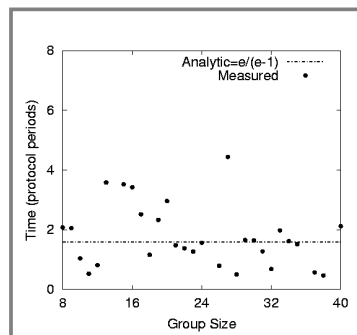
- Current implementation
  - Win2K, uses Winsock 2
  - Uses only UDP messaging
  - 900 semicolons of code (including testing)
- Experimental platform
  - Galaxy cluster: diverse collection of commodity PCs
  - 100 Mbps Ethernet
- Default protocol settings
  - Protocol period=2 s; K=1; G.C. and Suspicion timeouts= $3 \cdot \text{ceil}[\log(N+1)]$
- No partial membership lists observed in experiments

44



Per-process Send and Receive Loads are independent of group size

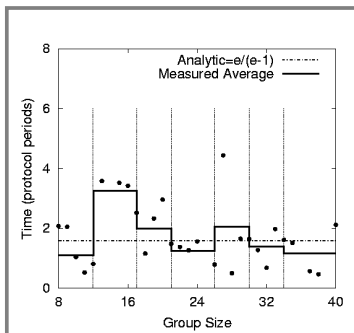
45



Time to First Detection of a process failure

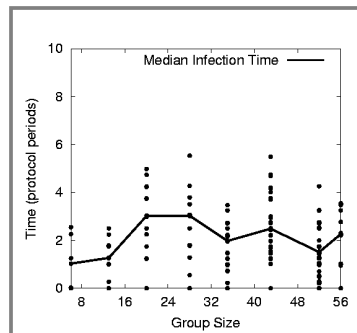
T1  
T1+T2+T3

46



Time to First Detection of a process failure apparently uncorrelated to group size

47



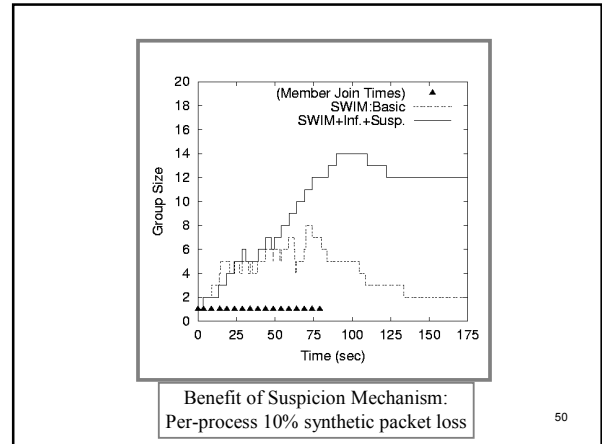
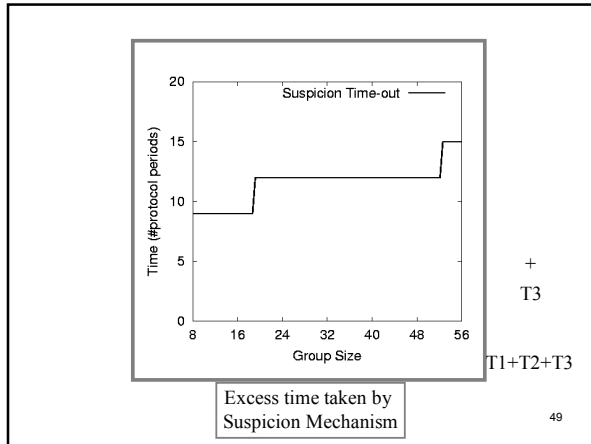
Membership Update Dissemination Time is low at high group sizes

+  
T2

T1+T2+T3

48





- ## More discussion points
- It turns out that with a partial list that is *uniformly random*, gossiping retains same properties as with complete lists
    - Why?
    - Partial membership protocols
      - SCAMP, Cyclon, TMAN, ...
  - Gossip-style failure detection underlies
    - Astrolabe
    - Amazon EC2/S3 (rumored!)
  - SWIM used in
    - CoralCDN/Oasis anycast service: <http://oasis.coralcdn.org>
    - Mike Freedman used suspicion mechanism to blackmark frequently-failing nodes
- 51

## Questions

52