CS 525 Advanced Distributed Systems Spring 2011

Indranil Gupta (Indy)

Membership Protocols (and Failure Detectors) March 31, 2011

All Slides © IG

Target Settings

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

Crash-stop/Fail-stop process failures

Group Membership Service



Two sub-protocols

Application Process *pi*



Unreliable Communication





I. pj crashes

- Nothing we can do about it!
- A frequent occurrence
- Common case rather than exception
- Frequency goes up at least 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

Distributed Failure Detectors: Properties



- 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



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

Failure Detector Properties



Network Message Load

Failure Detector Properties



Failure Detector Properties

- Completeness
- Accuracy
- Speed

In spite of arbitrary simultaneous process failures

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

Centralized Heartbeating







Gossip-style Heartbeating



Gossip-Style Failure Detection



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} seconds?



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]

N/2 nodes in a subnet



Analysis/Discussion

- What happens if gossip period $\mathsf{T}_{\mathsf{gossip}}$ 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

Simulations

As # members increases, the detection time increases

As requirement is loosened, the detection time decreases



Failure Detector Properties ...

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



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

All-to-All Heartbeating

Gossip-style Heartbeating

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_{ml}

•
$$L^* = \frac{\log(PM(T))}{\log(p_m)} \cdot \frac{1}{T}$$

(proof in PODC 01 paper)

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

SWIM Failure Detector Protocol

SWIM versus Heartbeating

SWIM Failure Detector $1-(1-\frac{1}{N})^{N-1} = 1-e^{-1}$

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

Accuracy, Load

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

for up to 15 % loss rates

Detection Time

- Prob. of being pinged in T'= $1 (1 \frac{1}{N})^{N-1} = 1 e^{-1}$
- $E[T] = T' \cdot \frac{e}{e-1}$
- Completeness: Any alive member detects failure
 - Eventually
 - By using a trick: within worst case O(N) protocol periods

III. Dissemination

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

Infection-style Dissemination

- Epidemic style dissemination
 - After $\lambda \log(pr)$ otocol periods, 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; 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
 - e.g., correlated message losses near pinged host
- Key: suspect a process before declaring it as failed in the group

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 #)
- Nothing overrides a (Failed, inc #)
 - See paper

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

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*ceil[log(N+1)]
- No partial membership lists observed in experiments

More discussion points

- It turns out that with a partial membership list that is *uniformly random*, gossiping retains same properties as with complete membership lists
 - Why? (Think of the equation)
 - 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

Reminder – Due this Sunday April 3rd at 11.59 PM

- Project Midterm Report due, 11.59 pm [12pt font, single-sided, 8 + 1 page Business Plan max]
- Wiki Term Paper Second Draft Due (Individual)
- Reviews you only have to submit reviews for 15 sessions (any 15 sessions) from 2/10 to 4/28. Keep track of your count! Take a breather!

Questions