Distributed Computing

José Orlando Pereira

Grupo de Sistemas Distribuídos Departamento de Informática Universidade do Minho

2010/2011



Goal 2: Why is it correct?

- With synchronous rounds, local state easily reflects global state
- What about in an asynchronous system?





- Remote invocation
- All processes request and reply to invocations
- A mutex is held while invoking remotely or handling remote invocations
- Distributed deadlock possible when multiple processes invoke each other



Distributed Computing

Global Predicates

Example: Distributed deadlock

Deadlock-free run:



© 2007-2010 José Orlando Pereira GSD/DI/U.Minho



Distributed Computing

Global Predicates

Example: Distributed deadlock

Distributed deadlock:





Instant observation is impossible:



* 🔿

Deadlock detection with a "wait for" graph:



A more complex deadlock-free run:





Distributed Computing

Global Predicates

Example: Distributed deadlock

• A deadlock-free WFG:



※ 🗘

• A WFG with a ghost deadlock:



Global Property Evaluation

- All these problems are instances of the Global Property Evaluation (GPE) problem
- Can it be solved in an asynchronous system?
- Methods that can be used? Relative cost?



Passive monitor process

Report all events to monitor:





First try: Synchronous system

- Global clock, δ upper bound on message delay
- Tag events with real time
- Consider events only up to t- δ
 - With synchronous rounds, this means using messages from the previous round!



Distributed Computing

Global Predicates

First try: Synchronous system





Clock properties

- What properties of a real-time clock make this approach correct?
- RC(i) the time at which i happened



Definition: Causality

- Events i and j are <u>causally related</u> $(i \rightarrow j)$ iff:
 - i precedes j in some process p
 - for some m, i=send(m) and j=receive(m)
 - for some k, $i \rightarrow k$ and $k \rightarrow j$ (transitivity)
- Events i and j are concurrent (i||j) iff neither i→j or j→i

Causality





Clock properties

- If $i \rightarrow j$ then RC(i)<RC(j)
- For some event j:
 - When we are sure that there is no unknown i such that RC(i)<RC(j)
 - Then there is no i such that $i \rightarrow j$
- Can we build a logical clock with the same property?

Second try: Logical clock

- Tag events as follows:
 - Local events: increment counter
 - Send events: increment and then tag with counter
 - Receive events: update local counter to maximum and then increment
- Use FIFO channels
- Consider events only up to the minimum of maximum tags



Distributed Computing

Global Predicates

Second try: Logical clock





Scalar clocks

Synchronous system (RC):

- Delay δ to consistency
- Asynchronous system (LC):
 - Possible unbounded delay to consistency
 - Blocks if some process stops sending messages



Third try: Vector clock

- Tag events with a vector as follows:
 - Local event at i: increment counter i
 - Send event at i: increment counter i and tag with vector
 - Receive event at i: update each counter to maximum and increment counter i



Third try: Vector clock





Causal delivery

The monitor delivers events as follows:

- With local vector I[...]
- For some r[...] from i
- Wait until:
 - I[i]+1=r[i]
 - I For all j≠i: r[i]≤l[i]
- The monitor is always in a consistent cut
- Blocking can be avoided by forwarding past messages



No reporting to monitor process

- Reporting all events to a monitor causes a large overhead
- Can a query be issued at some point in time?



Fourth try: No reporting, synchronous

- Monitor broadcasts tss in the future
- At tss, each process:
 - Records state
 - Sends messages to all others
 - Starts recording messages until receiving a message with RC > tss
- After stopping, sends all data to monitor



Distributed Computing

Global Predicates

Fourth try: No reporting, synchronous



© 2007-2010 José Orlando Pereira

Distributed Computing

Global Predicates

Fifth try: No reporting, logical clock





Chandy and Lamport

- Send a "Snapshot" message to some process
- Upon receiving for the first time:
 - Records state
 - Relays "Snapshot" to all others
 - Starts recording on each channel until receiving "Snapshot"
- Send all data to monitor



Chandy and Lamport





Global Property Evaluation

- GPE requires no gaps in observed history, regarding causality
- What properties can be evaluated?



Cuts and consistency

- A <u>cut</u> is the union of prefixes of process history
- A <u>consistent cut</u> includes all causal predecessors of all events in the cut
- Intuitive methods:
 - If a cut is an instant, there are no messages from the future
 - In the diagram, no arrows enter the cut
 - All events in the frontier are concurrent



Consistent cuts





Consistent global states





Consistent global states

- Includes the true sequence of states in the system
- An observer within the system cannot deny any of the possible paths





Stable predicates

- Once true, always true
- Examples:
 - Deadlock detection
 - Termination
 - Loss of token
 - Garbage collection
- Can be evaluated periodically on snapshots



Stable predicates







Non-stable predicates

- True in a subset of observable states
- Some are <u>possibly true</u>: an observer in the system cannot deny having been true
- The predicate does not hold on some paths



Non-stable predicates

- True in a subset of observable states
- Some are <u>definitely true</u>: an observer in the system is sure of having been true
- The predicate holds on all possible paths





Non-stable predicates

Examples:

- Total size of queues in the system
- Number of messages in transit
- Amount of memory used
- Can be detected by full monitoring of all (relevant) events



Conclusion

Second goal achieved:

- Causality
- Global predicate evaluation



* 🔿