

Lesson 10

Distributed Mutual Exclusion

Ch 10 [BenA 06]

Distributed System
Distributed Critical Section
Ricart-Agrawala
Token Passing Ricart-Agrawala
Token Passing Neilsen-Mizuno

1.12.2009 Copyright Teemu Kerola 2009 1

(Generic) Distributed System

- Nodes have processes
- Communication channels between nodes
 - Each node connected to every other node
 - Two-way channel
 - Reliable communication channels
 - Provided by network layer below
 - Messages are not lost
 - Messages processed concurrently with other computations (e.g., critical sections)
 - Nodes do not fail
- Requirements reduced later on
 - courses on distributed systems topics

Unrealistic assumptions? Not really...

1.12.2009 Copyright Teemu Kerola 2009 2

(Generic) Distributed System

- Processes (nodes) communicate with (asymmetric) messages
 - Message arrival order is not specified
 - Transmission times are arbitrary, but finite
 - Message (header) does not include send/receiver id
 - Receiver does not know who sent the message
 - Unless sender id is in the message itself

node 5

integer k ← 20
send(request, ③ k, 30)

→

node ③

integer m, n
receive(request, m, n)

1.12.2009 Copyright Teemu Kerola 2009 3

Distributed Processes

- Sender does not block
- Receiver blocks (suspended wait) until message of the proper type is received
- Atomicity problems in each node is not considered here
 - Solved with locking, semaphores, monitors, ...
- Message receiving and subsequent actions are considered to be atomic actions
 - Atomicity within each system considered solved

1.12.2009 Copyright Teemu Kerola 2009 4

1.12.2009 Copyright Teemu Kerola 2009 5

Distributed Critical Section Problem

- Processes within one node
 - Problem solved before
- Processes in different nodes
 - More complex
- State
 - Control pointer (CP, PC, program counter)
 - Local and shared variable values
 - Messages
 - Messages, that have been sent
 - Messages, that have been received
 - Messages, that are on the way
 - Arbitrary time, but finite!

Where are these?

1.12.2009 Copyright Teemu Kerola 2009 6

Two Approaches

- Ask everybody for permission, if it is my turn now
 - Lots of questions/answers
- I'll wait until I get the token, then it is my turn
 - Pass the token to next one (which one?)
 - Wait until I get the token
 - Token (turn) goes around all the time
 - Moves only when needed?
- Both approaches have advantages/disadvantages
 - Who is "everybody"? How do I know them?
 - What if someone does not talk to me?
 - What if node/network breaks down?
 - What if token is lost?

Do not worry now about the token getting lost ...

1.12.2009

Copyright Teemu Kerola 2009

7

Ricart-Agrawala for Distributed Mutex



G. Ricart A. K. Agrawala

- Distributed Mutex, 1981 (Lamport, 1978)
- Modification of Bakery algorithm with ticket numbers
- Idea
 - Must know all other processes/nodes competing for CS
 - Choose own ticket number, "larger than previous"
 - Send it to everybody else
 - Wait until permission from everybody else
 - Exactly one will always get permission from everybody else?
 - All others will wait
 - Do your CS
 - Give CS permission to everybody else who was waiting for you

mutex, no deadlock, no starvation?

1.12.2009

Copyright Teemu Kerola 2009

8

Algorithm 10.1: Ricart-Agrawala algorithm (outline)

```

integer myNum ← 0
set of node IDs deferred ← empty set

main application process, needs distr mutex
p1: non-critical section
p2: myNum ← chooseNumber ← not trivial!
p3: for all other nodes N
p4: send(request, N, myID, myNum) ← Each one answers only when it is safe. Reply needs no content.
p5: await reply's from all other nodes
p6: critical section
p7: for all nodes N in deferred ← all those waiting for my permission
p8: remove N from deferred
p9: send(reply, N, myID)

receive server process, runs concurrently all the time
integer source, reqNum
p10: receive(request, source, reqNum) ← most recent myNum
p11: if reqNum < myNum ← make these wait by not sending reply
p12: send(reply,source,myID)
p13: else add source to deferred
    
```

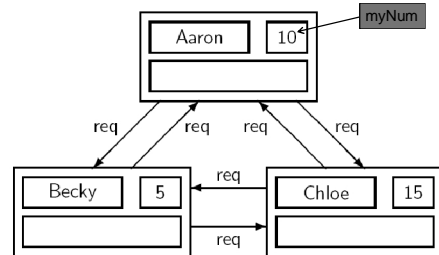
1.12.2009

Copyright Teemu Kerola 2009

9

Ricart-Agrawala Example

- 3 processes, each trying to enter CS concurrently
 - No status information needed on who had CS last



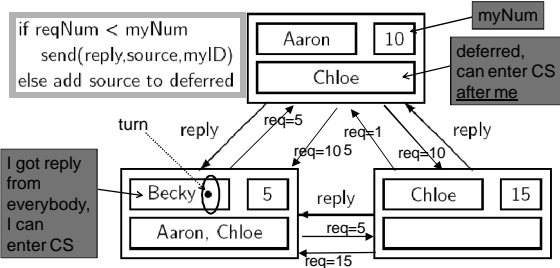
1.12.2009

Copyright Teemu Kerola 2009

10

Ricart-Agrawala Example (contd)

- Receive process runs at each node
 - What if Aaron's receive completes 1st? Last? Becky's? not yet?



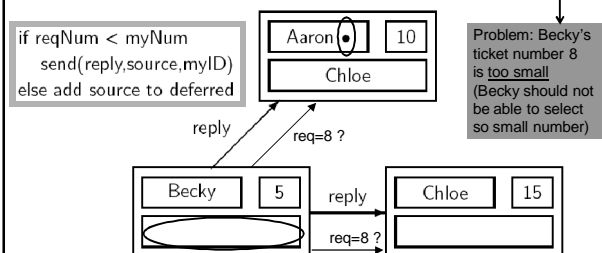
1.12.2009

Copyright Teemu Kerola 2009

11

Ricart-Agrawala Example (contd)

- Becky executes CS, and then sends deferred replies to Aaron & Chloe
- Aaron has now replies from everybody, and it can enter CS
- What if Becky now selects ticket number 8, and requests CS?
 - Aaron's and Chloe's receive will both reply immediately? Ouch!



1.12.2009

Copyright Teemu Kerola 2009

12

How to select ticket numbers

- Select always larger one than you have seen before
 - Larger than your previous *myNum*
 - Larger than any *requestedNum* that you have seen
 - They all came before you, and you should not try to get ahead of them
- What if equal ticket numbers?
 - Fixed priority, based on node/process id numbers
 - Used only with equal ticket numbers to avoid deadlock
 - Just like in Bakery algorithm

1.12.2009

Copyright Teemu Kerola 2009

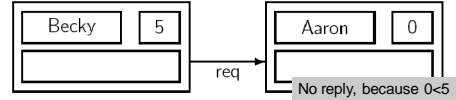
Discussion A

13

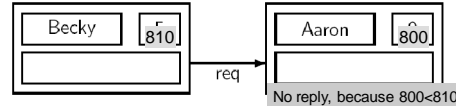
Quiescent Nodes

(hiljaiset solmut)

- Nodes that do not try to enter CS (but they could)
 - They are still listed in “all other nodes” if reqNum < myNum
 - Problem with initial value of *myNum* send(reply,source,myID)
 - Initial value zero? else add source to deferred



- Initial value $N > 0$; tickets numbers eventually will reach it



- Cure: *receive* checks for tickets numbers only if *main* wants CS

1.12.2009

Copyright Teemu Kerola 2009

14

Algorithm 10.2: Ricart-Agrawala algorithm

```
integer myNum ← 0
set of node IDs deferred ← empty set
integer highestNum ← 0
```

Main

loop forever

- p1: non-critical section
- p2: requestCS ← true
- p3: myNum ← highestNum + 1
- p4: for all other nodes N
- p5: send(request, N, myID, myNum)
- p6: await reply's from all other nodes
- p7: critical section
- p8: requestCS ← false
- p9: for all nodes N in deferred
- p10: remove N from deferred
- p11: send(reply, N, myID)

- Keep track of highest number seen
- What if one process asks for CS all the time?
- Same myNum OK?

(Receive on next slide)

1.12.2009

Copyright Teemu Kerola 2009

15

Algorithm 10.2: Ricart-Agrawala algorithm (continued)

Receive

- ```
integer source, requestedNum
loop forever
p1: receive(request, source, requestedNum)
p2: highestNum ← max(highestNum, requestedNum)
p3: if not requestCS or requestedNum < myNum
p4: send(reply, source, myID)
p5: else add source to deferred
```

original article

<http://www.cs.gatech.edu/class/s/AY2002/cs6210/fallpapers/MutualExclusion.pdf>

- Mutex between main & receive?
  - Exact mutex boundaries?
- What to do when myNum overflows?
  - Restart everybody? When? How?
  - Fairness is not the problem, mutex is
- Correctness proofs
  - Mutex? No deadlock? No starvation?

1.12.2009

Copyright Teemu Kerola 2009

Discussion B

16

### Token Based Algorithms

- Problems with permission based algorithms
  - Need permission from everybody (very many?)
  - Inactive participants (those not wanting in CS) slow you down
    - Need reply from all of them!
    - Lots of synchronization even if only one tries to get into CS
    - →→→ Lots of communication (many messages)
- Token based algorithms
  - Have token, that is enough
    - No synchronization with everybody else needed
  - Get token, send token is simple
    - Communicate only with a few (fewer) nodes
    - Scalable?
  - Mutex is trivial, how about deadlock and starvation?

1.12.2009

Copyright Teemu Kerola 2009

17

1.12.2009

Copyright Teemu Kerola 2009

18

### Ricart-Agrawala ideas

- Send token to next one only when I know that someone wants it
  - o/w keep token until needed
- Keep local *requested* array for best knowledge for the most recent CS request times
  - Update this based on received CS request messages
- Keep *granted* array, that has precise knowledge when each node actually was last granted CS
  - Update it only when CS granted
  - Pass it with token to next node
    - Only this *granted* array (with token) is exactly correct!
    - Other nodes have (slightly) old *granted* array

1.12.2009

Copyright Teemu Kerola 2009

19

### Algorithm 10.3: Ricart-Agrawala token-passing algorithm

```

boolean haveToken ← true in node 0, false in others
integer array[NODES] requested ← [0,...,0] ← local data in node
integer array[NODES] granted ← [0,...,0] ← distributed global data
integer myNum ← 0
boolean inCS ← false

sendToken
 if exists N such that requested[N] > granted[N]
 for some such N
 send(token, N, granted)
 haveToken ← false
 If no one else wants token, I will keep it

Receive
 server process, runs all the time
integer source, reqNum
loop forever
 receive(request, source, reqNum)
 requested[source] ← max(requested[source], reqNum)
 if haveToken and not inCS
 sendToken ← Give also most recent granted[]

```

1.12.2009

Copyright Teemu Kerola 2009

20

### Algorithm 10.3: Ricart-Agrawala token-passing algorithm (continued)

```

Main application process, needs distr mutex
loop forever
 non-critical section
 if not haveToken
 myNum ← myNum + 1
 for all other nodes N
 send(request, N, myID, myNum)
 receive(token, granted)
 haveToken ← true
 inCS ← true
 critical section
 granted[myID] ← myNum
 inCS ← false
 sendToken

```

1.12.2009

Copyright Teemu Kerola 2009

Discussion C 21

### Algorithm 10.3: Ricart-Agrawala token-passing algorithm (continued)

```

Main application process
loop forever
 non-critical section
 if not haveToken
 myNum ← myNum + 1
 for all other nodes N
 send(request, N, myID, myNum)
 receive(token, granted)
 haveToken ← true
 inCS ← true
 critical section
 granted[myID] ← myNum
 inCS ← false
 sendToken

```

1.12.2009

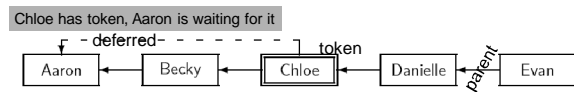
Copyright Teemu Kerola 2009

22

### Neilsen-Mizuno Token Based Algorithm



- Ricart-Agrawala: token carries queue of waiting processes
  - Token can be very large, which may be problematic
- Neilsen-Mizuno: virtual tree structure within the nodes implements the queue
  - Algorithm utilizes *virtual spanning tree* of nodes
    - Spanning tree: all nodes linked as a tree, no cycles
  - Simple *token* indicates "turn" for critical section
  - *Parent* link points to the direction of last in line for CS
    - Parent == 0: node may have token and is last in line for CS
  - *Deferred* link points to next in line for CS



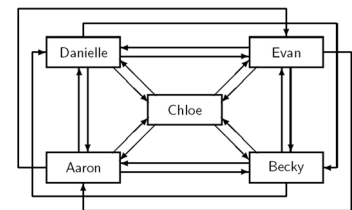
1.12.2009

Copyright Teemu Kerola 2009

23

### Neilsen-Mizuno Example

- Fully connected nodes
- Chloe is in CS
- No one waits for CS



1.12.2009

Copyright Teemu Kerola 2009

24

### Neilsen-Mizuno Example (contd)

- Chloe has token, nobody waits for it
- Aaron requests CS
  - Sends msg=(req, Aaron, Aaron) on parent link
  - Removes himself from parent spanning tree
- Becky receives msg, and forwards the request "upward"
  - Sends msg=(req, Becky, Aaron) to Chloe
  - Moves to new parent spanning tree, points to Aaron
  - Aaron is now last to request CS

1.12.2009 Copyright Teemu Kerola 2009 25

### Neilsen-Mizuno Example (contd)

- Chloe receives msg (req, Becky, Aaron)
  - Chloe in CS, sets deferred field to Aaron and sets parent field to Becky
  - Chloe was (also) last in line for CS
- When Chloe completes CS, she will pass token to Aaron
- Token transferred directly to the next process in line for critical section (if any)
  - Just token is passed, no big array with it

1.12.2009 Copyright Teemu Kerola 2009 26

### Neilsen-Mizuno Example (contd)

- Chloe still has CS, Evan wants CS
  - Sends (req, Evan, Evan) to Danielle
  - Danielle sends (req, Danielle, Evan) to Chloe
  - Chloe sends (req, Chloe, Evan) to Becky
  - Becky sends (req, Becky, Evan) to Aaron
  - Aaron makes a deferred link to Evan

1.12.2009 Copyright Teemu Kerola 2009 27

### Neilsen-Mizuno Example (contd)

- Chloe completes CS, passes token to Aaron
- Aaron completes CS, passes token to Evan
- Evan completes CS, keeps token

1.12.2009 Copyright Teemu Kerola 2009 28

### Algorithm 10.4: Neilsen-Mizuno token-passing algorithm

```

integer parent ← (initialized to form a tree)
integer deferred ← 0
boolean holding ← true in the root, false in others

Main
loop forever
 p1: non-critical section
 p2: if not holding
 p3: send(request, parent, myID, myID)
 p4: parent ← 0
 p5: receive(token)
 p6: holding ← false
 p7: critical section
 p8: if deferred ≠ 0
 p9: send(token, deferred)
 p10: deferred ← 0
 p11: else holding ← true

```

1.12.2009 Copyright Teemu Kerola 2009 29

### Algorithm 10.4: Neilsen-Mizuno token-passing algorithm

```

Receive (runs concurrently with main, mutex problems solved...)
integer source, originator
loop forever
 p12: receive(request, source, originator)
 p13: if parent = 0
 p14: if holding
 p15: send(token, originator)
 p16: holding ← false
 p17: else deferred ← originator
 p18: else send(request, parent, myID, originator)
 p19: parent ← source

```

1.12.2009 Copyright Teemu Kerola 2009 Discussion D 30

### Ricart-Agrawala vs. Neilsen-Mizuno

- Number of messages needed
- Size of messages
- Size of data structures in each node
- Behaviour with heavy load
  - Many need CS at the same time
- Behaviour with light load
  - Requests for CS do not come often
  - Usually only one process requests CS at a time

1.12.2009

Copyright Teemu Kerola 2009

31

### Other Distributed Mutex Algorithms

- Other token-based algorithms
  - Token ring: token moves all the time
  - Lots of token traffic even when no CS requests
- Centralized server
  - Simple, not very many messages
  - Not scalable, may become bottleneck
- Give up unrealistic assumptions
  - Nodes may fail
  - Messages may get lost, token may get lost
- See other courses



Courses on distributed systems topics (hajautetut järjestelmät)

1.12.2009

Copyright Teemu Kerola 2009

32

### Summary

- Distributed critical section is hard, avoid it
  - Use centralized solutions if possible?
- Permission based solutions
  - Ricart-Agrawala – ask everyone
- Token based solutions
  - Ricart-Agrawala – centralized state in granted[]
  - Neilsen-Mizuno – queue kept in spanning tree
- There are other algorithms
- How do they scale up?

1.12.2009

Copyright Teemu Kerola 2009

33