



CS 142: Lecture 5.2 Mutual Exclusion

Richard M. Murray 1 November 2019

Goals:

- Review Lamport's algorithm (from Mon) and prove correctness
- Describe token-based algorithms for mutual exclusion (ring + tree)
- If time: performance comparison

Reading:

- P. Sivilotti, Introduction to Distributed Algorithms, Chapter 7
- [SS94] M. Singhal and N. G. Shivaratri. *Advanced Concepts in Operating Systems*. McGraw-Hill, 1994. (Chapter 6: Distributed Mutual Exclusion)

Lamport's Mutual Exclusion Algorithm





Idea: treat request queue as a distributed atomic variable

- reqQ: queue of timestamps requests for CS (sorted in increasing order)
- knownT: list of last "known times" for other processes

UNITY program: list of actions that can be executed by each agent (in any order)

- SendReq: mode = NC \rightarrow mode := TRY || ($\forall j$:: send(i, j, $\langle req_i, t_i \rangle$))
- EnterCS: mode = TRY ^ recQ[head] = $\langle req_i, t_i \rangle$ ^ ($\forall j :: knownT[j] > ti) \rightarrow mode := CS;$
- ReleaseCS: mode = CS \rightarrow mode := NC || reqQ.pop($\langle req_i, t_i \rangle || (\forall j :: send(i, j, \langle rel_i, t_i \rangle))$
- RecvReq: $(\exists j :: recv(i, j) = \langle req_j, t_j \rangle \rightarrow recQ.push/sort(\langle req_j, t_j \rangle) || send(i, j, \langle ack_i, t_i \rangle))$
- RecvAck: $(\exists j :: recv(i, j) = \langle ack_j, t_j \rangle \rightarrow knownT[j] := t_j)$
- RecvReI: $(\exists j :: recv(i, j) = \langle rel_j, t_j \rangle \rightarrow reqQ.pop(\langle rel_j, t_j \rangle)$

Sample Execution

 $\begin{array}{ll} {SendReq:} \ mode = NC \rightarrow mode = TRY \mid (\forall j :: send(i, j, \langle req_i, t_i \rangle)) \\ [] {RecvReq:} (\exists j :: recv(i, j) = \langle req_j, t_j \rangle \rightarrow recQ.push/sort(\langle req_j, t_j \rangle) \mid send(i, j, \langle ack_i, t_i \rangle)) \\ [] {RecvAck:} (\exists j :: recv(i, j) = \langle ack_j, t_j \rangle \rightarrow knownT[j] := t_j) \\ [] {EnterCS:} \ mode = TRY \ recQ[head] = \langle req_i, t_i \rangle \ (\forall j :: knownT[j] > ti) \rightarrow mode = CS; \\ [] {ReleaseCS:} \ mode = CS \ \rightarrow mode = NC \mid reqQ.pop(\langle rel \backslash q, t_i \rangle \mid (\forall j :: send(i, j, \langle rel_i, t_i \rangle)) \\ [] {RecvRel:} (\exists j :: recv(i, j) = \langle rel_j, t_j \rangle \rightarrow reqQ.pop(\langle ack_j, t_j \rangle) \\ recQ: {\langle reqj, tj \rangle, ...} \\ knwT1: [log, kT2, kT3] \end{array}$



Proof of Correctness

Safety: need to show that no two processes are in CS at the same time

- \bullet Assume the converse: U_i and U_j are both in CS
- Both U_i and U_j must have their own requests at head of queue
- Head of U_i: <req_i, t_i>
- Head of U_j: <req_j, t_j>
- Assume WLOG $t_i < t_j$ (if not, switch the argument)
- Since U_j is in its CS, then we must have t_j < U_j.knownT[i] ⇒ <req_i, t_i> must be in U_j.reqQ (since messages are FIFO)
- $t_i < t_j \implies req_j$ can't be at the head of U_j.reqQ
- $\rightarrow \leftarrow$ (contradiction)

Progress: need to show that eventually every request is eventually processed

- Approach: find a metric that is guaranteed to decrease (or increase)
- One metric: number of entries in Ui.knownT that are less than its request time (ti)
 - Represents number of agents who might not have received our request
- Is this a good metric?
 - Bounded below by zero and if at zero then we eventually enter our critical section
 - Must always decrease as other processes enter their critical section (and someone will execute their CS at some point in time)

Ui reqQ

 $\langle reqi, ti \rangle$

⟨reqj, tj⟩

 $- t_i < t_j < U_j.knownT[i] ⇒ U_i$

has acknowledged reqi

Uj reqQ

(reak, 1)

Proof of Correctness (Progress)

```
 \begin{array}{l} {\rm SendReq:} \ {\rm mode} = {\rm TRY} \mid (\forall j :: {\rm send}(i, j, \langle {\rm req}_i, t_i \rangle)) \\ {\rm RecvReq:} \ (\exists j :: {\rm recv}(i, j) = \langle {\rm req}_j, t_j \rangle \rightarrow {\rm recQ.push/sort}(\langle {\rm req}_j, t_j \rangle) \mid {\rm send}(i, j, \langle {\rm ack}_i, t_i \rangle)) \\ {\rm RecvAck:} \ (\exists j :: {\rm recv}(i, j) = \langle {\rm ack}_j, t_j \rangle \rightarrow {\rm knownT[j]} := t_j) \\ {\rm ReterCS:} \ {\rm mode} = {\rm TRY} \ {\rm recQ[head]} = \langle {\rm req}_i, t_i \rangle \ {\rm (\forall j :: {\rm knownT[j]} > ti)} \rightarrow {\rm mode} = {\rm CS}; \\ {\rm ReleaseCS:} \ {\rm mode} = {\rm CS} \ \rightarrow {\rm mode} = {\rm NC} \mid {\rm reqQ.pop}(\langle {\rm rel} \backslash q, t_i \rangle \mid (\forall j :: {\rm send}(i, j, \langle {\rm rel}_i, t_i \rangle)) \\ {\rm RecvRel:} \ (\exists j :: {\rm recv}(i, j) = \langle {\rm rel}_j, t_j \rangle \rightarrow {\rm reqQ.pop}(\langle {\rm ack}_j, t_j \rangle) \end{array}
```

Proof steps

- Metric: number of entries in Ui.knownT that are less than its request time (ti)
- Need to show that this is guaranteed to decrease \Rightarrow eventually U_i can enter CS
- Consider an agent U_j with with an entry less than U_i's request time: U_i.knownT[j] < t_i (where <req_i, t_i) is the request from U_i)
- Agent U_j's logical time is guaranteed to increase when $\langle req_i,\,t_i\rangle$ is received by U_j
- Agent U_j will send U_i an acknowledgement with t_j > t_i (logical clock property) => metric will decrease by 1

Additional steps for complete proof:

Optimizations on Lamport's Algorithm

Optimization #1 - if request sent with later timestamp than received, no ack needed



Optimization #2 (Ricart-Agrawala): merge release messages with replies

- Receive req: add to reqQ || send <ack_i, t_i> to U_j only in certain conditions:
 - if not requesting access to CS nor in CS
 - if requesting CS and U_j timestamp is smaller than Ui request
- When U_i exits CS, send release (clears all deferred acknowledgements)
- Idea: eventually, the pending request will be granted and we will send a <release> message, which will serve as the acknowledgement



Optimizations on Lamport's Algorithm

Optimization #1 - if request sent with later timestamp than received, no ack needed

Optimization #2 (Ricart-Agrawala): merge release messages with replies

Optimization #3 (Maekawa): request permission from a subset of agents [SS94]

- U_i only sends requests to a subset of sites R_i, chosen such that $(\forall i, j : R_i \cap R_j \neq \{\})$
- Rough idea: every pair of agent has an agent that "mediates" between the pair
- Each agent sends only one ack at time; wait until a release is received => mediating agents can only grant permission if they haven't granted permission to another agent
- Proof is more complicated, but many fewer messages required

Timing analysis [SS94]

- Response time = amount of time to execute critical section (if no queue)
- Sync delay: U_i exits CS to U_j enters CS
- T = transport time: E = execution time



NON-TOKEN	Resp. time (ll)	Sync Delay	Messages (ll)	Messages (hl)	ll = low load condition
Lamport	2T+E	T	3(N-1)	3(N-1)	(one req at a time)
Ricart-Agrawala	2T+E	T	2(N-1)	2(N-1)	$n_l = nign load cond.$
Maekawa	2T+E	2T	$3\sqrt{N}$	$5\sqrt{N}$	(U/ seldom in NC)
Maekawa	2T+E	2T	$3\sqrt{N}$	$5\sqrt{N}$	(Ui seldom in f

Token-Based Approaches to Mutual Exclusion

Simple token ring: send token clockwise (CW) around ring

• Problem: potentially long sync delays (especially in low load)

Token ring with requests: add request message type

• Minimize sync delay by waiting for requests to arrive (CCW)

```
SendReq (Try \rightarrow CS)
                                  RecvReq
if holder \neq self
                                      if holder = self \land \neg using
                                           send token on (CW)
      hungry := true
                                       else
     if \neg asked
                                                                          hungry = Ui wants resource
                                            pending\_requests := true
          send request (CCW)
                                                                          asked = token requested
                                           if holder \neq self \land \neg asked
           asked := true
                                                                          holder = token is with Ui
                                                send request (CCW)
                                                                          using = Ui in critical section
     wait until using
                                                asked := true
else
                                  RecvToken
      using := true
[use the resource]
                                       asked := false
using := false
                                      if hungry
if pending_requests
                                            using := true
     send token on (CW)
                                            hungry := false
                                       else
      pending\_requests := false
                                                                               requests
                                            send token on (CW)

  Both approaches require

                                            pending\_requests := false
```

existence of ring topology

Token Tree Algorithm (Raymond)

Basic idea: pass request toward token & token toward request U1 Maintain a tree with root being agent with token Each agent gets requests and passes them toward token • Agent with token either enters CS (if earliest request) U2 U3 passes token to node w/ earlier request Passing the token also updates the direction of links (so that root stays with the token) • Invariant: graph is (rooted) tree, root at token U4 U5 U6 U7 Question: how do we prove this works?

- Safety: no two nodes are in CS at same time (easy)
- Progress: need to show that all requesting nodes eventually get access
 - Trick (as usual): figure out a metric that captures this



Richard M. Murray, Caltech CDS

Remarks on Token Tree Algorithm

Algorithm

- If node wants CS, doesn't hold token and requestQ is empty, send request toward the root (= token location) and adds request to its requestQ
- 2. If agent receives request, place on requestQ and pass message toward the token
- 3. When root receives request, send token toward request and update parent link
- 4. When agent receives the token, remove top entry in requestQ, send token to that entry, and update parent link. If requestQ is non-empty, send request to (new) parent
 - Necessary step since we need token back to send to next entry in queue
- 5. Agent enters critical section when it has the token and its entry is at top of requestQ
- 6. After site has finished execution in critical session, got to step 4

Basic idea behind the proof [SS94]

- Agents execute requests in first come first serve (FCFS) order
- Let Ui = agent requesting access, Uh = agent holding the token
 - Always exists path Ui, Ui₁, Ui₂, ... Ui_{k-1}, Ui_k, Uh
 - When Uh gets request from Uik, there are two possibilities
 - Ui_{k-1} is at the top of Ui_k 's requisite => send token there => chain gets shorter
 - Some other agent Uj is at top queue => send token there first
 - Can show Ui_k will eventually get token back => will eventually get to Ui_{k-1} (\rightarrow)



Comparisons between different algorithms

Timing analysis [SS94]

- Response time = amount of time to execute critical section (if no queue)
- Sync delay: Ui exits CS to Uj enters CS
- T = transport time: E = execution time
- II = low load condition (one req at a time)
- hl = high load condition (Ui seldom in NC)

CS Request arrives Its request sent out The site enters The site exits the CS the CS

TABLE 6.1

A comparison of performance (ll = light load, hl = heavy load)

NON-TOKEN	Resp. time (ll)	Sync Delay	Messages (ll)	Messages (hl)
Lamport	2T+E	T	3(N-1)	3(N-1)
Ricart-Agrawala	2T+E	T	2(N-1)	2(N-1)
Maekawa	2T+E	2T	$3\sqrt{N}$	$5\sqrt{N}$
TOKEN	Resp. time (ll)	Sync Delay	Messages (ll)	Messages (hl)
TOKEN Suzuki and Kasami	Resp. time (<i>ll</i>) 2 <i>T</i> + <i>E</i>	Sync Delay	Messages (ll)	Messages (hl)
TOKEN Suzuki and Kasami Singhal's Heuristic	Resp. time (<i>ll</i>) 2 <i>T</i> + <i>E</i> 2 <i>T</i> + <i>E</i>	Sync Delay T T	Messages (ll) N N/2	Messages (hl) N N

Summary: Mutual Exclusion

Key ideas:

- Distributed protocol for allow access to a shared resource ("critical section")
- Two approaches: distributed atomic variables (Lamport + variants) or token-based
- User process specifications:

NC next $NC \lor TRY$ stable.TRY

CS next $CS \lor NC$

transient.CS

- System specifications:
 - Safety: no two users (Ui) are in critical section (CS) at the same time
 - Progress: all agents will get a chance (as long as they keep requesting)

Friday: problem solving session (board talk)

Next week: specification refinement, conflict resolution (dining philosophers)



