# Distributed Solutions using Timestamps

Shankar

April 25, 2014

Event ordering in a distributed system



Let x and y be two statement executions (aka events)

#### Define x causally precedes y if

- x and y happened in that order in the same system, or
- x sent a message that y received, or
- transitive closure of above
- Causal precedence is a partial order
  - if x and y not causally related, no system can determine which happened first (without other interaction or real-time clocks)
- Timestamp mechanism extends partial order to total order for a specified set of events

### Outline

Timestamp mechanism Distributed ordering of conflicting requests Distributed lock program: algorithm level Distributed lock program: overview Cyclic timestamps

- Each system j has
  - integer "clock" clk, initially 0
- When j does an event x to be ordered:
  - increment clk, broadcast [x, clk, j]

■ *clk value*: timestamp (ts) of *x* 

[clk value, id]: extended timestamp (ets) of x

When j receives msg 
$$[y, t, k]$$
:  
clk  $\leftarrow \max(t, clk) + 1$ 

x.ets < y.ets: (x.ts < y.ts) or (x.ts = y.ts and x.id < y.id)</li>
Event x ordered before event y if x.ets < y.ets</li>

For most applications, need acks to timestamps

- Each system j has
  - integer clock *clk*, initially 0
  - $rts_k$ ,  $k \neq j$ , initially 0

## // last ts rcvd from k

■ let  $\alpha Rts$ : min([ $rts_k, k$ ] :  $k \neq j$ ) // no new ets <  $\alpha Rts$ 

- When *j* does an event *x* to be ordered:
  - increment clk, broadcast [x, clk, j]
- When j receives msg [y, t, k]:
  - $clk \leftarrow max(t, clk) + 1$
  - $rts_k \leftarrow t$ , send [ack, clk, j] to k
- When *j* receives msg [ack, *t*, *k*]:
  - $clk \leftarrow max(t, clk), rts_k \leftarrow t$

#### Properties

- Define auxiliary quantities
  - *hst*: seq of all ets's in ets-order
  - *j*.*hst*: seq of ets's seen by *j* in ets-order // initially [[0,0]]
  - $j.\alpha \overline{hst}$ : prefix of  $j.\overline{hst}$  of ets's  $\leq j.\alpha Rts$

#### Safety properties

- $Inv j.\overline{hst}$  subsequence-of  $\overline{hst}$
- $Inv j.\alpha \overline{hst}$  prefix-of  $\overline{hst}$

Progress properties (assuming no system stops rcving)

•  $j.\alpha Rts = z < \overline{hst}.last$  leads-to  $j.\alpha Rts > z$ 

•  $j.\alpha \overline{hst}.size = z < \overline{hst}.size$  leads-to  $j.\alpha \overline{hst}.size > z$ 

# // initially [[0,0]]// initially [[0,0]]

## Outline

Timestamp mechanism Distributed ordering of conflicting requests Distributed lock program: algorithm level Distributed lock program: overview Cyclic timestamps

- Collection of systems attached to a fifo channel
- Users issue conflictable requests to the systems
- Each system should serve its requests so that conflicting requests are not served simultaneously (even by different systems)
- Some special cases of the problem
  - distributed lock
    - every two requests conflict
  - distributed readers-writers lock
    - classify requests into reads and writes
    - write conflicts with every other request

- System j augments the ts mechanism as follows
- Maintain variable req: set of "ongoing" req-ets tuples
- Upon local request x: assign ts, add [x,ts,j] to req
- Upon reving [y,t,k]: process ts, add [y,t,k] to req
- Serve [x,t,j] in req when:
  - [t,j] <  $\alpha$ Rts and
  - [t,j] < [u,k] for every conflicting [y,u,k] in req</pre>
- After serving [x,t,j]: remove it from req, bcast [REL,x,t,j]
- Upon rcving [REL,y,t,k]: remove [y,t,k] from req

- System j variables
  - clk: initially 0
  - rts<sub>k</sub>: initially 0 αRts: min(rts<sub>k</sub>, k]: k ≠ j)
  - req: initially empty

// clock // highest ts rcvd from k // min ets induced by rts // set of outstanding requests-ets

- Messages
  - EREQ,x,t,k]
  - [ACK,t,k]
  - [REL,x,t,k]

// request msg
 // ack msg
 // release msg

- User isues request x
  - clk++
  - send [REQ,x,clk,j] to every system
  - add [x,c]k,j] to req
- Receive [REQ,x,t,k]
  - $clk \leftarrow max(clk, t+1)$
  - rts[k]  $\leftarrow$  t
  - send [ACK,c1k,j] to k // omit if ets > [t,k] already sent to k
  - add [x,t,k] to req
- Start serving request [x,t,j] in req when
  - [t,j]  $\leq \alpha \mathsf{Rts}$
  - for every [y,s,k] in req st x conflicts with y: [t,j]  $\leq$  [s,k]

Finish serving request [x,t,j]:

- remove [x,t,j] from req;
- send [REL,x,t,j] to every other system.
- Receive [ACK,t,k]
  - clk  $\leftarrow$  max(clk,t); rts[k]  $\leftarrow$  t
- Receive [REL,x,t,k]:
  - remove [x,t,k] from req

atomicity assumption: rules are atomic

progress assumption: weak fairness

## Outline

Timestamp mechanism Distributed ordering of conflicting requests Distributed lock program: algorithm level Distributed lock program: overview Cyclic timestamps

- Distributed program that implements a distributed lock
- Collection of systems attached to a fifo channel
- Specialize the request-ordering solution for a lock
- At most one ongoing request per system
  - so each system is thinking, hungry, or eating
  - no need for ts in release msg

Later, refine to await program implementing dist lock service

Solution: variables, functions, messages



Become hungry only if thinking

- clk++
- send [REQ,c1k,j] to every system
- ∎ req[j] ← clk

Become eating only if hungry and  $[\mathsf{req}_{\mathbf{j}},\mathbf{j}] = lpha\mathsf{Req} \leq lpha\mathsf{Rts}$ 

- Become thinking only if eating:
  - remove entry for j from req
  - send [REL, j] to every system

- Receive [REQ,t,k]
  - clk ← max(clk, t+1)
  - rts[k]  $\leftarrow$  t
  - req[k]  $\leftarrow$  t
  - send [ACK,c]k,j] to k
- Receive [ACK,t,k]:
  - clk  $\leftarrow$  max(clk,t); rts[k]  $\leftarrow$  t
- Receive [REL,k]:
  - remove entry for k from req
- atomicity assumption: rules are atomic
- progress assumption: weak fairness

Goal: Inv at most one system is eating

Inv  $A_1 - A_4$  holds, where

 $A_1$ : (([j,s] in k.req) and  $j \neq k$ )  $\Rightarrow$  ([j,s] in j.req) or ([REL,j] in transit to k)

 $A_2$ : (j eating)  $\Rightarrow$  [j.req[j], j] = j. $\alpha$ Req  $\leq$  j. $\alpha$ Rts

 $A_3$ : ((j eating) and (k eating))  $\Rightarrow$  j = k

 $A_4$ : ((j hungry) and [j.req[j], j] = j. $\alpha$ Req  $\leq$  j. $\alpha$ Rts)  $\Rightarrow$  (no one eating)

■ *Inv* A<sub>3</sub> implies desired property

Analysis: progress

Goal: (wfair, bounded eating, ongoing rx, channel progress)  $\Rightarrow$  j hungry *leads-to* j eating

- Define
  - hst: seq of all ets's in ets-order

// initially [[0,0]]

ne: # requests that have finished eating

#### Proof

- if [j,s] is in j.req, eventually [s,j]  $\leq$  j. $\alpha$ Rts holds
- after this point
  - [j,s]'s index in hst is fixed, at say n
  - entries in hst[ne+1..n] eat in order

[entry  $\overline{ne}$ 's release msg is incoming to entry  $\overline{ne+1}$ 's system. when it arrives, the latter eventually becomes eating] Timestamp mechanism Distributed ordering of conflicting requests Distributed lock program: algorithm level Distributed lock program: overview Cyclic timestamps Distributed program: implements distributed lock service

- starts a fifo channel
- starts a LockTs system at each address

LockTs: await program, refines algorithm-level system

- input functions acq and re1 // called by lock users
- output calls to tx and rx of channel access system
- one local thread to execute rx
- multiple acq calls can be ongoing but only one participates in ts mechanism

```
program LockTsDist(ADDR)
    {c<sub>j</sub>} ← start(FifoChannel(ADDR))
    for j in ADDR
        v<sub>j</sub> ← start(LockTs(ADDR, j, c<sub>j</sub>))
    return {v<sub>j</sub>}
```

Program LockTs (ADDR, j, c<sub>j</sub>) – 1

#### Main

- clk  $\leftarrow$  0
- $rts_k \leftarrow 0$ , k in ADDR-{j}}

req

- startThread(doRx())
- input mysid.acq()

```
• await (not (j in req.keys) // a1

• clk++; req_j \leftarrow clk

• for k in ADDR-{j}

c_j.tx(k, [REQ, clk, j])

• await ([req_j, j] \le \alpha Req and // a2

(ADDR.size=1 or [req_j, j] \le \alpha Rts))

• return
```

Program LockTs (ADDR, j, c<sub>i</sub>) - 2

- input mysid.rel()
  - await (true)
    - req.remove(j)
    - for k in ADDR-{j}
      cj.tx(k,[REL,j])
- function doRx()
  while true

 $//\ensuremath{\left/\right.}$  executed by a local thread

• msg  $\leftarrow$  c<sub>j</sub>.rx() ia {msg is [REQ, t, k], [ACK, t, k], or [REL,k]}

await true

do appropriate rx-msg action

- atomicity assumption {awaits}
- progress assumption {wfair threads, sfair await a1}

Map alg-level state to await-program state

alg-level		await-program
j hungry	$\leftrightarrow$	thread in j.acq.a2
j eating	$\leftrightarrow$	[j,.] in j.req, no thread in j.acq.a2
j thinking	$\leftrightarrow$	no[j,.] in j.req

■ Show that alg-level properties are preserved (★)

- Prove: LockTsDist(ADDR) implements DistLockService(ADDR)
  - define program of implementation and service inverse
  - identify effective atomicity breakpoints
  - obtain assertions
  - prove program satisfies assertions

// easy given ★

## Outline

Timestamp mechanism Distributed ordering of conflicting requests Distributed lock program: algorithm level Distributed lock program: overview Cyclic timestamps

## Using cyclic timestamps

- Goal: cyclic timestamps in the distributed lock solution
- Easily achieved by modifying solution slightly
- Existing solution: request [t,j] eats when
   1. [t,j] = j.αReq
   2. [t,j] ≤ j.αRts
- Impose additional requirement:
  - 3. j eats only after rcving ack from every system
- Resulting simplification
  - ack's ts always higher than request's ts
    - so no need for ack's ts
    - no need for {rtsk}
  - sufficient to track # acks rcvd
    - no need for ack's sender id

### Solution: variables, functions, messages

cyclic timestamps

- System j variables
  - ∎ clk, <del>{rts<sub>k</sub>}</del>, req
  - na // # acks due
- System j functions
  - $\alpha Rts$ ,  $\alpha Req$
- Messages
  - [REQ,t,k], [ACK,t,k], [REL,k]

## Become hungry only if thinking

- clk++
- req[j]  $\leftarrow$  clk
- send [REQ,c1k,j] to every system

```
∎ na ← 0
```

```
Become eating only if hungry and

[req_j, j] = \alpha Req and na = ADDR.size - 1
```

- Become thinking only if eating:
  - remove entry for j from req
  - send [REL, j] to every system

- Receive [REQ,t,k]
  - $clk \leftarrow max(clk, t+1)$
  - req[k]  $\leftarrow$  t
  - send [ACK] to k
- Receive [ACK]:
- Receive [REL,k]:
  - remove entry for k from req
- atomicity assumption: rules are atomic
- progress assumption: weak fairness

## Analysis: conventions

#### Abbreviations for readability:

- hstj.ts to mean hst[j][0]
- N to mean ADDR.size

#### Define

- hst: seq of all ets's, initially [[0,0]]
- ne: # releases globally, initially 0
- j.ne: # releases seen by j
- tse: ts of last request to release
- j.tse: ts of last request released at j

// as before
// as before
// may lag ne
// hstne.ts
// hstj.ne.ts

- Prove: any ts in transit is in tse.. tse + 2N
- Prove: j.tse is in  $\overline{\text{tse}} 2N..\overline{\text{tse}}$
- Hence: any ts in transit is in j.tse..j.tse + 4N
- Hence can use modulo-M ts, for M  $\geq$  4N
- Modify system j to use cyclic ts
  - add variable tse, initially 0
  - when a request msg is sent, set its to mod(clk, M)
  - when a request msg [REQ, ct, j] is rcvd, treat ct as unbounded ts tse + mod(ct-tse, M)
  - when [k,t] is removed from req, set tse to t

Following are invariant

 $C_1$ : ([REQ,t,j] rcvable)  $\Rightarrow \overline{\text{hst}}_{\overline{\text{ne}}}$ .ts  $\leq t \leq \overline{\text{hst}}$ .last.ts

 $C_2$ : forsome(x in hst: x.ts  $\leq$  i.clk  $\leq$  x.ts+1)

$$C_3$$
:  $\overline{hstp.ts} \le \overline{hst}_{p+1}.ts \le \overline{hstp.ts}+2$ 

$$\begin{array}{ll} C_4: ([REQ,t,j] \mbox{ rcvable}) & \Rightarrow \\ & \overline{\mbox{hst}}_{\overline{\mbox{ne}}}.\mbox{ts} \leq \mbox{t} \leq & \overline{\mbox{hst}}_{\overline{\mbox{ne}}}.\mbox{ts} + 2 \times \mbox{ADDR.size} \end{array}$$

- Inv C<sub>1</sub>: [REQ,t,j] in transit implies req [t,j] hungry
- Inv C<sub>2</sub>: C<sub>2</sub> satisfies invariance rule
- Inv  $C_3$ :  $C_3$  satisfies invariance rule assuming  $Inv C_2$
- Inv  $C_4$ : follows from Inv  $C_1$ ,  $C_3$

Following are invariant

- $C_5$ : (# REL msgs incoming to j) < ADDR.size
- $C_6$ :  $\overline{\text{ne}}$  j.ne = (# REL msgs incoming to j)
- Inv  $C_5$ : [t,k] is acked only after k's previous REL msgs are rcvd
- Inv  $C_6$ :  $C_6$  satisfies invariance rule
- Inv  $C_7$ :  $C_7$  implied by  $C_6$ ,  $C_5$ ,  $C_3$
- Inv  $C_8$ :  $C_8$  implied by  $C_7$ ,  $C_4$