

The Dictionary Problem by New & Bernstein

• Replicated Log to achieve mutual consistency

• The log contains: record of updates on the objects

Model

- Same as before - N nodes.

- (E, \rightarrow)

C1. events at the same node are totally ordered

C2. send^(a) and receive(e2): $e_1 \rightarrow e_2$.

The Log Problem.

type Event =

record

op: OpType

time: TimeType

node: NodeID.

event record of e eR .

$eR.op$ $op(e)$ - operation

$eR.time$ $time(e)$ -

$eR.node$ $node(e)$

$f \rightarrow e$ iff $fRGL(e)$

Dictionary

- Same as before.

P2. $x \in V(e)$ iff $e_x \rightarrow e$ \nexists x -delete event e'
s.t. $e' \rightarrow e$.

$$V(e) = \{x \mid e_x RGL(e) \wedge \nexists g RGL \text{ s.t. } gR.op = delete(x)\}$$

O1. The entire log is sent in each message.

O2. A new view is repeated computed

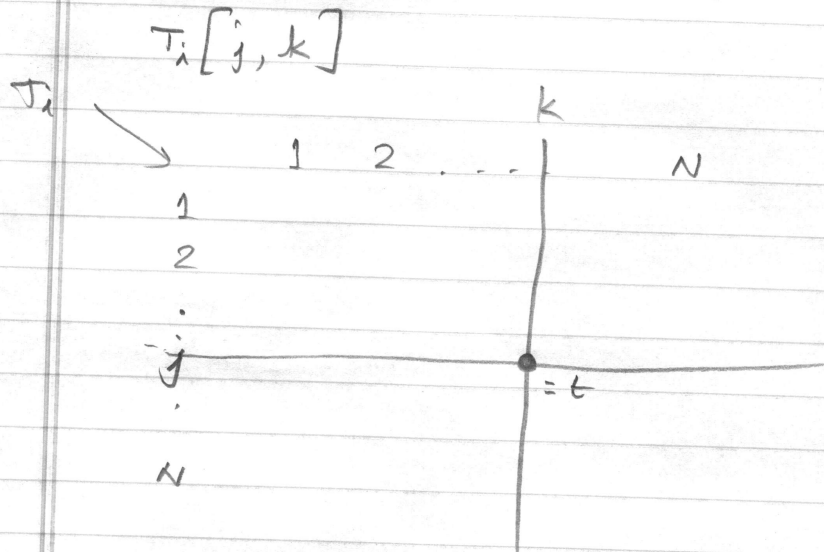
O3. The entire log is maintained

Each node N_i maintains:

1. $clock_i$

2. A two dimensional time-table T_i

$$T_i: [1..N] [1..N]$$



$T_i[j, k] = t$ - N_i knows that N_j has learned of all events up to at N_k up to time $t @ N_k$

$T_i[i, i] = N_i$'s clock.

$T_i[i, k]$ - N_i 's knowledge of events at N_k .

3. A copy of the Log.

$\text{hastecvd}(T_i, eR, k) \equiv T_i[k, eR.\text{node}] \geq eR.\text{time}$

$\Rightarrow N_i$ knows that N_k has learned of e .

Init

$$V_i = \emptyset$$

$$PL_i = \emptyset$$

$$T_i[*,*] = \emptyset$$

insert(x):

$$T_i[i,i] = \text{clock}_i++$$

$$PL_i = PL_i \cup \{ \langle \text{ins}(x), T_i[i,i], i \rangle \}$$

$$V_i = V_i \cup \{x\}$$

delete(x):

$$T_i[i,i] = \text{clock}_i++;$$

$$PL_i = PL_i \cup \{ \langle \text{del}(x), T_i[i,i], i \rangle \}$$

$$V_i = V_i \setminus \{x\}.$$

send(m) to N_k .

$$NP = \{ eR \mid eR \in PL \wedge \neg \text{hasrecvd}(T_i, eR, k) \}$$

send $\langle NP, T_i \rangle$ to N_k .

receive $\langle NP_k, T_k \rangle$ from N_k .

$$NE = \{ FR \mid FR \in NP_k \wedge \neg \text{hasrecrd}(T_i, FR, i) \}$$

$$V_i = \{ x \mid x \in V_i \text{ or } e_x R NE \}$$

$$\wedge \{ \cancel{\exists} d R NE \text{ s.t. } dR.op = del(x) \}.$$

$$\forall I \text{ do } T_i[i, I] = \max \{ T_i[i, I], T_k[k, I] \}$$

$$\forall I, J \quad T_k[I, J] = \max \{ T_i[I, J], T_k[I, J] \}$$

$$PL_i = \{ eR \mid eR \in PL_i \cup NE \} \wedge$$

$$(\exists j \text{ s.t. } \neg \text{hasrecrd}(T_i, eR, j))$$

Global Snapshots in Distributed Systems.

Source: Chandy & Lamport

ACM TOCS, 3(1):63-75, 1985.

Motivation:

- global state of a distributed computation.
- why:
 - i) deadlock in the system?
 - ii) computation has terminated?
 - iii) checkpointing the system state (why?)

Why is it difficult?

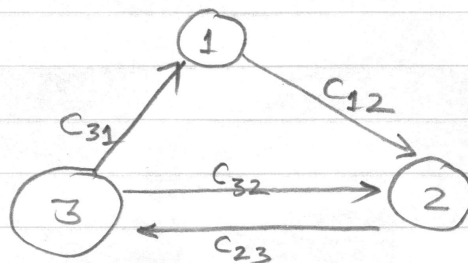
- process can record its own state + the messages it sends + receives.
- it can record nothing else.
- processes must cooperate with others
- all processes cannot record their local states at the same time.
- no common clock / no shared memory.
- the problem is compute global snapshot:
local states + comm. channel states.

Finally, should not impede / interfere with the underlying computation.

System Model.

A distributed system:

- a finite set of processes, and
- a finite set of channels

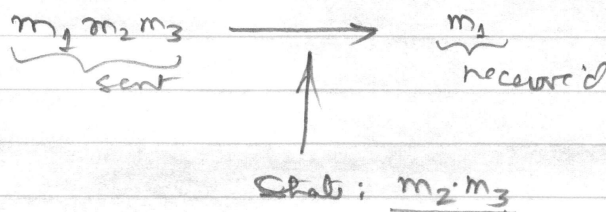


Channels:

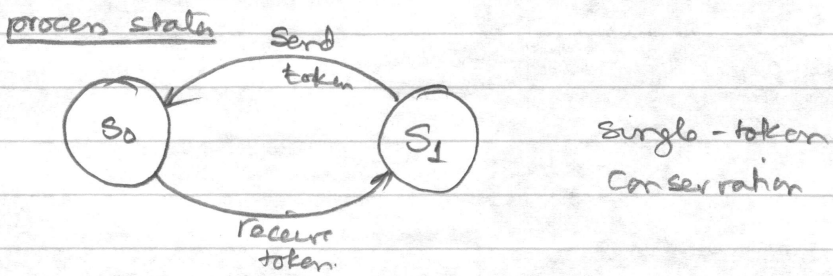
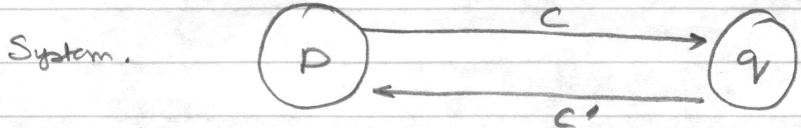
- infinite capacity
- error-free
- deliver messages in the order sent.

The state of a channel:

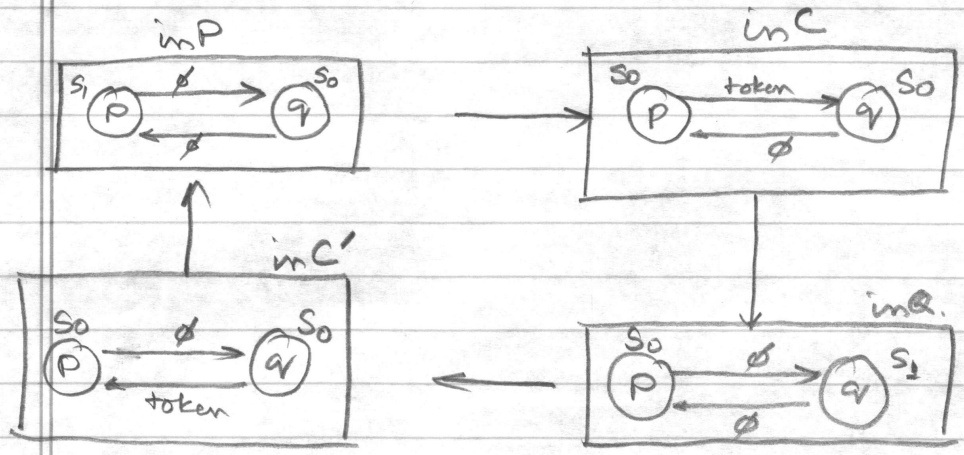
Sequence of messages sent excluding the messages received



Example 2.1



P in S1 ; q in S0.



The Algorithm.

- use an example to motivate the steps of the algorithm.
- ^{assume that} can record the state of the channel instantaneously.

Ex.

1 record state of P in in-P.

↓

in-C

↓

record: state of P
state of C
state of C'

⇒ 2 tokens in-P + in-C.

Inconsistency: state of P before P sent a message
along C

+

State of C after P sent the message.

Let n be the # of messages sent along C before P's state is
Let n' be the # " " " C before C's state

Global state inconsistent if
 $n < n'$.

Alternate scenario:

state(c) recorded in-P,

state(P), state(c), state(c') recorded in-C,

⇒ No token.

Inconsistent if c records before P sends a message along c and state of P after P send that is

$$n > n'$$

⇒ Consistent global state
 $n = n'$.

Similarly if $m = \#$ of messages rec'd along C before
 $m' = \#$ of messages rec'd along C ~~before~~ before
C' state is recorded.

$m = m'$ for consistency.

$$\underbrace{n'}_{\text{sent}} \geq \underbrace{m'}_{\text{rec'd}}$$

$$m \geq m'$$

Alg.

Marker-Sending Rule for a Process p .

For each output channel C of p :

p sends a "marker" along C after recording its state and before sending any further messages.

Marker-receiving Rule for a Process q :

q receiving a marker along a channel C :

if q has not recorded its state then:

q records its state

q records the state C as the empty seq.

q sends marker on all outgoing C'

else

q records the state of C as the seq. of messages rec'd. along C after q 's state was recorded and before q rec'd. the marker along C .