

Distributed Systems Group
Department Software Technology
Faculty EEMCS
DELFT UNIVERSITY OF TECHNOLOGY

Distributed Algorithms (IN4150)

List of algorithms (2018-2019)

January 3, 2019

Below is a list of the algorithms in the order of treatment in the course, with the names of the original authors, if applicable.

Chapter 3: Synchronization

1. Alpha-, beta-, and gamma-synchronizer: Awerbuch
2. Causal message ordering (broadcast): Birman-Schiper-Stephenson
3. Causal message ordering (point-to-point): Schiper-Eggli-Sandoz
4. Total message ordering
5. Determining global states: Chandy-Lamport
6. Termination detection in a unidirectional ring
7. Termination detection in a general network
8. Deadlock detection for AND requests: Chandy-Misra-Haas
9. Deadlock detection for OR requests: Chandy-Misra-Haas
10. Deadlock detection for M-out-of-N requests (with/without instantaneous communication): Bracha-Toueg

Chapter 4: Coordination

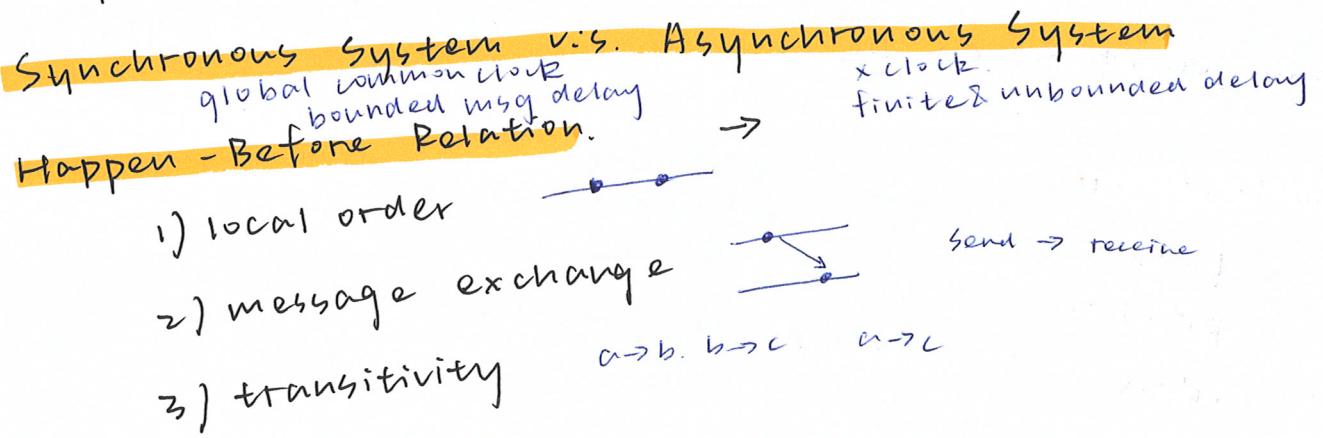
1. Assertion-based mutual exclusion: Lamport
2. Assertion-based mutual exclusion: Ricart-Agrawala
3. Assertion-based mutual exclusion: Maekawa
4. Generalized assertion-based mutual exclusion

5. Token-based mutual exclusion: Suzuki-Kasami
6. Token-based mutual exclusion: Singhal
7. Token-based mutual exclusion (tree-based): Raymond
8. Detection of loss and regeneration of a token
9. Election in a synchronous unidirectional ring (non-comparison-based)
10. Election in a bidirectional ring: Hirschberg-Sinclair
11. Election in a bidirectional ring (enhanced version)
12. Election in a unidirectional ring: Chang-Roberts
13. Election in a unidirectional ring: Peterson
14. Election in a synchronous complete network: Afek-Gafni
15. Election in an asynchronous complete network: Afek-Gafni
16. Minimum-weight spanning trees: Gallager-Humblet-Spira

Chapter 5: Fault Tolerance

1. Agreement with stopping failures
2. Byzantine agreement with oral messages: Lamport-Pease-Shostak
3. Byzantine agreement with authentication: Lamport-Pease-Shostak
4. Single-value Paxos consensus: Paxos
5. Practical Byzantine Fault Tolerance: Castro and Liskov
6. Zyzzyva: Kotla et al.
7. Randomized agreement with crash failures: Ben-Or
8. Randomized Byzantine agreement: Ben-Or
9. Stabilizing mutual exclusion: Dijkstra
10. Stabilizing stop-and-wait datalink algorithm
11. Non-stabilizing and stabilizing sliding-window datalink algorithms

chapter 3 Synchronization



Concurrency Relation

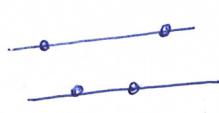
all to
neither $a \rightarrow b, b \rightarrow a$.

not transitive

Logical clock.

Definition

Implementation:
integer counter



assign timestamp to events $C:E \rightarrow S$.
characterize HB relation.

1) **scalar clock / Lamport clock**
Defin. & rules { Internal clock
send & receive

Problems?
concurrent processes
do not have equal clock value

2) **vector clock**
vector.

rules { internal

send & receive

$$b = c_j = \max(c_{ia} + 1, c_{j+1})$$
$$\max(v_i + e_j, v_{i+1})$$

definition

meaning.

achieve concurrent?

characterize HB relation.

time
different time granularity
achieved by msg.

Problem in Synchronization.

Synchronous Communication vs. Asynchronous communication

synchronous
occur simultaneously
a single msg transfer event
sending blocked.

asynchronous
separate
sending does not block
receiving block/interrupted.

Synchronizer

Simulation of SS on AS

Synchronizer

Purpose

problem

Idea

run synchronous algorithms
on asynchronous system

when issue clock pulse?

sender figure out all msg received.
notify receiver.

α -Synchronizer

Idea

ACK.

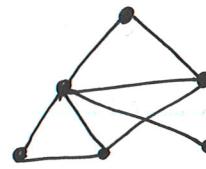
safe nodes?

when to proceed?

complexity

communication $O(N^2)$

time $O(1)$



• → ACK
← ACK
→ SAFE

β -Synchronizer

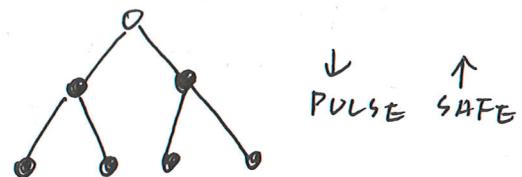
Idea

spanning tree
PULSE, SAFE

complexity

communication $O(N)$

time $O(N)$



↓ PULSE ↑
SAFE

γ -Synchronizer

Idea

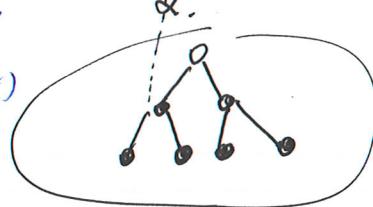
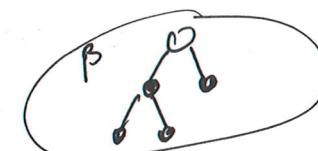
initialization: preferred-link
PULSES and SAFEs

CLUSTER-SAFEs and READYs

complexity

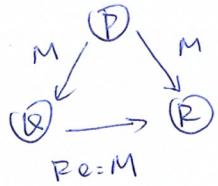
communication $O(1)$

time $O(1)$



Message Ordering

Message Order Violation



Problem

requirement on the order of msg received
e.g. see reply later than origin.
replicated database
(process update)

Message Destination types. $Dest(m)$

1) point-to-point $|Dest(m)| = 1$

2) multicast

≥ 1

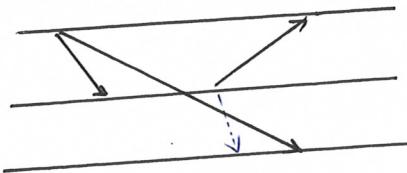
3) broadcast

processes

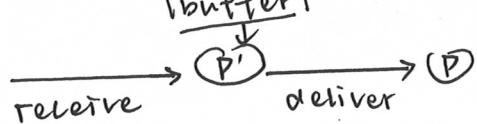
Causal message Ordering.

$$m(m_1) \rightarrow m(m_2)$$

$$d(m_1) \rightarrow d(m_2)$$



Implement of MO



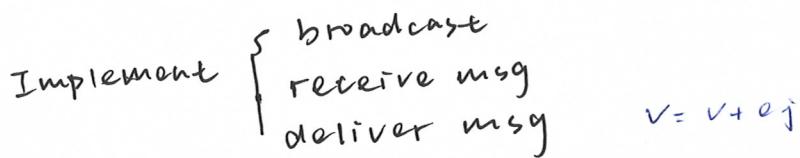
Causal Ordering: Broadcast NO FIFO Required

Idea vector clock. message = content + timestamp

number broadcast msg. can't deliver. = buffer

receiving process → deliver this msg & check buffer

$$D_j(m) = V + e_j \geq V_m$$



Causal Ordering: point-to-point NO FIFO Required

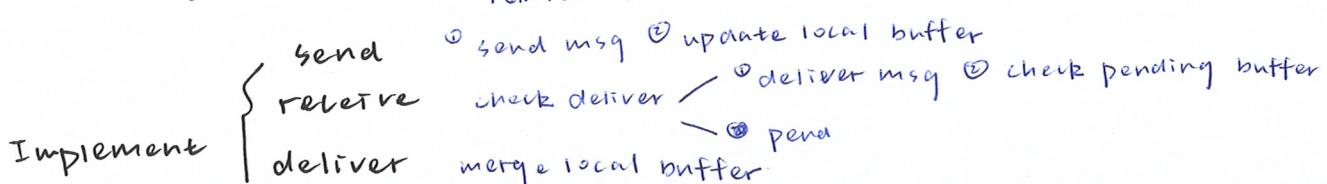
Idea vector clock.

local buffer $S = \{(P_j, V_j)\} \dots \}$

message = content + local buffer + timestamp

deliver condition $\begin{cases} \text{not in } S \\ \text{exists } i, V'_j \leq V_i \end{cases} \rightarrow \text{delivered.}$

buffer merge sm & s.



Total message ordering

all processes receive in same order

Simple solution

sequencer { number msg
broadcast
record history.

do broadcast → send to sequencer.

Sophisticate algorithm (FIFO required)

scalar clock (with process ID as tie breaker)

ordered message queue

ACK msg.

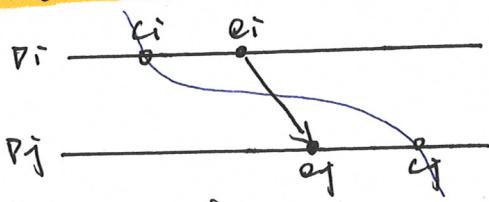
deliver condition { head of queue
all ACK received.

Global States

Problem: global state in asynchronous system

Application: ① debug ② detect e.g. bank.

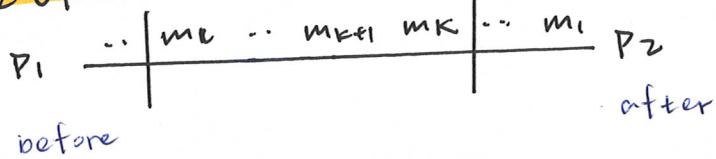
Consistent cut: internal. one \rightarrow every



vector time of C : $V(C)[i] = \max\{V(c_1)[i], \dots, V(c_n)[i]\}$

$$V(C) = (V(c_1)[1], V(c_2)[2], \dots, V(c_n)[n])$$

State of Channel



Chandy's and Lamport's Atqo (FIFO)

Idea: record state \rightarrow ① record own ② marker send

receive marker first: ① state channel ② record own

③ marker bank ④ create buffer

later: ① state channel

Termination Detection

Problem: distributed computation terminate
active \rightarrow passive
passive \rightarrow active

Termination Detection in unidirectional ring (FIFO)

Idea:

Naive Solution

Problem

Colored
solution

P₀: special

token = knowledge.

receive \ passive

token \ active.

Reactivate

introduce color in processes and token

process \rightarrow turn black: msg spots passiveness

initial white

token \rightarrow turn black = arrive black.

Final decision \rightarrow token \rightarrow black: *

\rightarrow token = white + P₀ passive

General Termination Detection.

Idea: special P
weights for processes and msgs

Initially, $w(P) = 1$ $w(P_i) = 0$.

then $w(P_i) = \frac{1}{n}$

if want to send msg: $w(P_i) = \frac{w_i}{2}$.

$$w(\text{msg}) = \frac{w_m + w_i}{2}$$

receive msg: $w(P_i) = \frac{w_i + w_m}{2}$

terminates: $w(P_i) = 0$ $w(P) =$

$$w(\text{msg}) = w(P_i).$$

Final decision: $w(P) = 1$

Deadlock Detection.

Problem: each request, wait

Wait-for-Graph: nodes
(WFG)
links

Resource Model: Resource RA

process request

resource granted

release resource

Condition: cycle

msg

block and wait, from dependent set

→ active, receive single msg.

Condition: Rule { a path

{ no edge in → out. }

Communication

Model:

Deadlock detection for AND request

Idea: suspect deadlock → send PROBE
propagate.

Condition: receive own PROBE

Δ_i dependent. P_j : path.

Condition: depth(i) = true

Deadlock detection for OR request

Idea: suspect deadlock → send QUERY
→ propagate. → construct tree
dependent set

receive QUERY → 1st engaged.

↓ later send REPLY, back

all REPLY received

Condition: every QUERY → REPLY

Deadlock detection for N-out-of-M Requests

Processors
— active
— blocked.

msg: REQUEST, REPLY, RELINQUISH

With instantaneous communication

Idea:

$$\text{OUT_IN} \rightarrow \text{DUT_IN}$$

2 phase {
 ① notify ② notify() \rightarrow spanning tree
 simulate. if no request pending \rightarrow free = true.
 grant()

end of receiving DONE

② GRANT simulate REPLY

sufficient GRANT \rightarrow active.

condition: initiator free = false

With messages in transit

Idea: color represent link.

{
 grey
 black
 white
 translucent

(optimizer) grey, white, translucent disappear \rightarrow apply only o.

grey, white edges resource granted

active: $n \in \# \text{greywhite}$

know color of edge \rightarrow send coloring with to IN & OUT set

Chapter 4 Coordination

Mutual Exclusion

Problem: grant exclusive privilege

Algo: { token-based: $\text{process} = \text{execute CS}$
assertion-based: request permission from all/part
reply → continue only one

Lamport's mutual-exclusion algorithms: Assertion-based (FIFO)

Idea: msg = scalar clock + PID. ordered request queue.

want CS: → REQUEST to all (+itself).

receive REQUEST → ① enter queue.
② send REPLY

enter CS if ① all REPLY ② head of queue.

leave CS: send RELEASE to all
↑ remove request.

Complexity $3(n-1)$.

Optimization: suppress unnecessary REPLY.
x needed if older REQUEST.

Assertion based: Ricart's and Agrawala's

Idea: combine REPLY and RELEASE RELEASE not needed.
REPLY { immediately. request higher > own.
deferr

enter CS: all REPLY.

leave CS: send deferred REPLY.

Complexity $2(n-1)$.

Optimization reduce request set.

Maekawa's Assertion-based algo:

Requirements

1) any two non-empty intersection (must)

Properties

2) every process in own

3) request size same size

4) every process in same number request set

Lemma

minimum size of request set is $O(\ln n)$

Idea: request list.

want CS? → send REQUEST to R.

enter CS = receive GRANT from R.

leave CS = send RELEASE to R.

Handling
deadlock

granted + receive REQUEST → compare timestamp

if new lat older = ① enter queue ② send POSTPONE

if old lat older: send INQUIRE to old.

receive INQUIRE: wait until

→ receive GRANT from all = complete CS ③ send RELEASE

→ receive one POSTPONE = ④ send RELINQUISH

receive RELINQUISH: ⑤ enqueue request

⑥ send GRANT to oldest

A generalized mutual-exclusion algorithm

Token based: Suzuki's and Kasami's

① Request for token.

Idea: token permission to CS.

want CS? ① number request.

② send REQUEST to all (itself).

Process maintain array H: for every p last request knows

TH: ... last request granted.

Token

compare N & TH. = forward token.

Avoid starvation: check own index.

Leave CS = update TH.

Send token to first $H[i] \rightarrow TH[i]$

else wait

Variation = In token maintain queue = pid should receive token.

(no starvation) leave CS = ① add to queue

② remove itself

③ send token to head.

Optimization reduce REQUEST

Token based: Singhal's

Optimization send REQUEST to may have token.

Idea Process state $\begin{cases} R & \text{requesting} \\ E & \text{executing CS} \\ H & \text{holding token} \\ O & \text{other} \end{cases}$

Process maintain array N = request number

S = States

Token

TM

S Initialization:

TS

$P_0 = S(0) = H \quad S(1) = O \quad \dots$

$P \in S \mapsto R. 1 \mapsto S(j) = O. 1 \mapsto 1$

want CS? ① set own requesting. ② send REQUEST to R.

Receive REQUEST = ① update request number.

② if local state

→ E.O. remote requesting

→ R. \leftarrow send REQUEST

→ H send token.

Receive token: update TH, N, S, TS.

Token based: Raymond let a single msg represent all requests. LIFO

Idea assume spanning tree is created.
root = head token. change v.

every process maintains ~~holder~~ owner
neighbor on the path to root

request queue owner
neighbor id.

want CS? ① append owner id in queue ② send REQUEST in direction root.

receive REQUEST: ① append neighbor id

② send REQUEST to parent

/ if is root: ① send token to head. ② dequeue

③ reassign parent

receive token: if is head, = enter CS

else: ① propagate to head ② send REQUEST

dequeue

assign parent

if queue not empty

Election

Problem: single process gets privilege
elect with largest id

Anonymous Network \rightarrow no id

Election in a synchronous unidirectional ring (non-comparison based)

Idea: elect minimum

ring size = n, known to all.

In round 1, inspect $id=1 \rightarrow$ send id to neighbor.

relay.

n, no msg \rightarrow doesn't exist!

= 2.

In round 2, inspect $id=2 \rightarrow$ send id to neighbor.

K.

(K-1)n+1

Complexity: $O(n)$

Election in bidirectional ring - Solution 1

Idea: check if largest id in every larger segment ($size=2^k+1$)

In round K, check 2^{K-1} both sides.

PROBE msg = id + hop counter

receive PROBE: if own larger: discard.

else : forward, hop--.

If hop count = 0, send OK.

initiate next round if receive 2. OK.

Elect Condition: if not known size = 2^{OK} from wrong side.
1 processor 2 PROBE.

Election in bidirectional ring - Solution 2

Idea: neighbour bounded by active processes.

In round 1, exchange msgs.

If own = largest \rightarrow active.

else \rightarrow passive

Subsequent phase: execute in virtual ring.

Elect Condition: receive own id

Election in unidirectional ring: Chang Robert

Idea: send msg toward to neighbour.

compare $\swarrow =$ elected.
 \searrow own larger. discard

\searrow smaller. relay

Election in unidirectional ring: Peterson

Idea simulation of bit slot.

send id to right. maxId, received from left \rightarrow right

if $id = \text{largest}$ \rightarrow active

else \rightarrow passive

Elect condition: receive own id.

$$\begin{array}{l} \circ \text{id} \\ \circ \text{max(id, id)} \end{array}$$

Election in complete network: Afek and Gafni. (synchrony)

Idea send id to ever larger set of nodes and wait for ACK.

Candidate process

if receive larger: \circ send ACK \circ adopt id.

if receive smaller msg: killed.

size is twice

Elect condition: receive ACK from all

node of candidate process = level

ordinary process: \circ receive order. \circ send ACK.

msg = level + id

elected = first start.

Election in asynchronous network: Afek and Gafni.

Idea no make sense = wait for all ACK.
react if receive single msg

msg = level + id level = captured process number.

$\xrightarrow{\text{capture}}$ IP (owned by R). R also needed - captured.

(P) $\xrightarrow{(1)}$ $\xrightarrow{(2)}$ node keep links { father
potential father }

In asynchronous system, previous owner might be killed.

Minimum Spanning Tree

Application: broadcast in minimum lost

Lemma: A weighted connected graph in which all weights are different has a unique MST.



some edges in T_1 but not T_2 .

e : lowest weight.

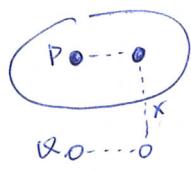
add e to $T_2 \rightarrow$ cycle.

($\therefore e'$ in T_2 not in T_1)

$w(e') > w(e)$

throw e' \rightarrow lower weight.

Lemma: F is fragment and e is MCE, F with e and node incident with e



Suppose not

final MST must have a $P \rightarrow Q$.

\therefore at least one edge outgoing edge

\because MCE $w(e) < w(x)$

fragment. subtree

outgoing edge

MCE (minimum-weight outgoing edges)

Idea: build from fragment, start with single node.

Fragment has level. 1 node = level 0.

nodes cooperate find MCE. fragment connect among MCE

merge: new fragment.

Absorb



same level, more

$v > v'$. postpone.

any start spontaneously. receive \rightarrow start

Byzantine Agreement with authentication

Idea: msg carry signature. ⚡ forged
loyal could verify

msg = initial value + signature. ($V_0 = v_0, \dots, s_{ik}$)

tenant wait until length = $f+1$.

maintain. set $V =$ received order.

execute same choice function

Randomized Byzantine Agreement with crash failure

Idea: every start with v . $f < \frac{n}{2}$

rounds consist 3 phases

{ notification phase = ⚡ broadcast NOTIFICATION. ⚡ wait for ($n-f$) H.

proposal: ⚡ if enough notification for one. broadcast it
 $\geq \frac{n}{2}$ else broadcast random?

decision if decided = stop

else wait $(n-f)$ P.

if > 1 P. adopt value

$\geq f$ P. decide that

else new round with random v .

Randomized Byzantine Agreement

Idea: every start with v . $n \geq 5f$.

notification phase ⚡ broadcast. ⚡ wait $(n-f)$ H.

proposal if enough support $\geq \frac{nf}{2}$ = propose that.

else

if decided = stop

else wait $(n-f)$ P.

if $> f$. Propose = adopt. value

$\geq 3f$ = decide

decision phase else. new round with random

Chapter 5 Fault Tolerance

Consensus: agree on same value

Application: 1) commit in distributed DB.

2) modify record in DB with replica

Fault classification

{
fail-stop
omission
performance

Byzantine

Byzantine Agreement

Problem

Goal: non-faulty processes agree on same

source/commander start

Lieutenant decide

Agreement = agree on same

Validity = source is non-faulty, decide initial

Termination = decide with finite time

Agreement with stopping failure

Idea:

$f =$

Every process starts with v .

Maintain set w , $w = \{v\}$ initially.

Each round: ① broadcast ② receive w_j .
 $w = w \cup w_j$.

Decide. w single element. \rightarrow decide

\rightarrow default

Byzantine Agreement with Oral Messages

Idea: $f+1$ levels.

Bottom case: $OM(0), f=0$.

① commander broadcast.

② decide that value.

$OM(f), f > 0$.

① command broadcast

② each as commander. execute $OM(f-1)$

Lieutenant decide on majority {commander, lieutenant, commander}

Stabilization

→ must implement stabilize function for each process

Problem:

- transparent fault
- transient
- bring incorrect \rightarrow correct state

Stabilizing Mutual-Exclusion in unidirectional ring.

Legal state $\{ \text{fin-}c_i=1 \} \subseteq$

Idea In a shared memory model: suspect own & predecessor
composite atomicity = R&W single step
Pi maintain single integer modulo K.
Ordinary protocol:

- ① compare
- ② unequal \rightarrow copy

Special protocol:

- ① compare
- ② equal \rightarrow increment by 1

legal: $K, K \dots K, K-1, K-1 \dots$

Stabilizing stop-and-wait datalink algorithm

Idea Sender & receiver maintain counter
Send msg \rightarrow receive ACK then send next.
timeout: resend.

A non-stabilizing datalink algorithm

Idea Sender maintain $\{ w, n_s, n_a \}$

Receiver maintain nr.

msg = msg number.

Receiver msg:

- ① nr++
- ② send ACK.

Receive ACK:

- ① if ACK higher = set na.

timeout: resend unacked

Predicate $(na \leq nr) \wedge (nr \leq ns) \wedge (ns \leq na + w)$

⇒

$(msg, i) \quad i \leq ns$

⇒

$(ack, i) \quad i \leq nr$

A stabilizing datalink algorithm

Idea: msg carry 2 integer. 2nd = na. \rightarrow give catch up with sender.

predicate

then \exists $i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

then $\exists i \in \text{enr. } \exists n \in \mathbb{N}$ $\forall j \in \text{enr. } i \leq j \leq i + n \Rightarrow \text{msg}[j] = \text{msg}[i]$

State Machine Replication

Problem execute client request in same order

Single Value Paxos

Idea Goal: decide a value

proposer, acceptor, learner.

Alg proceed in round. int.

round num assign to proposer.

proposers ignore msg with smaller round.

Phase 1: build majority. acceptor + propose.

receive: if newest, reply latest vote & round number

Phase 2: propose latest vote to majority

(if not voted, random)

acceptor vote positive / abstain to learner

Practical Byzantine Fault Tolerance

View = primary backup. view num

client send request to primary

primary assign number, multicast to backup.

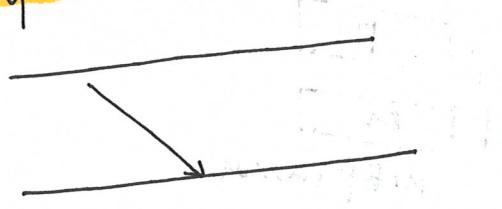
replica execute & reply to client

client wait for $F+1$ same result

if no, client send all req to replica.

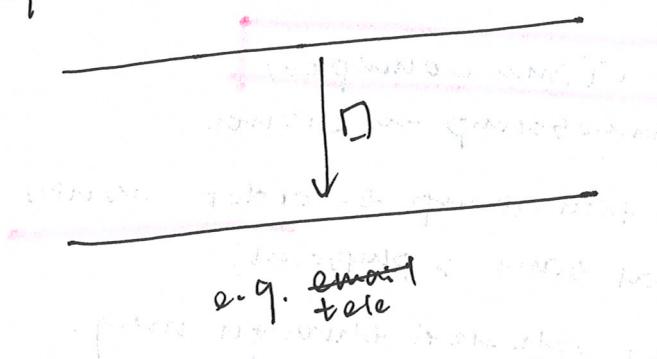
17/11/15 02 Modeling distributed system

Asynchronous communication (AC) / Non-blocking course



just send, doesn't check whether receive it / wait

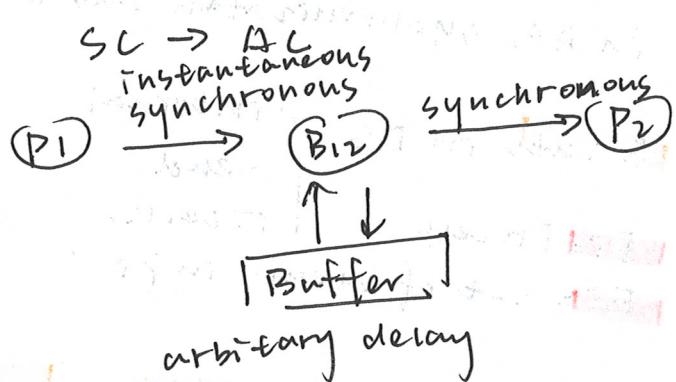
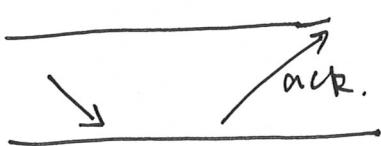
Synchronous communication (SC)



send message and wait to receive it.
i.e. waiting process
receiving block
no message in transit

simulate

AC \rightarrow SC



Link Properties

FIFO. receive message in order they sent

Buffer (unbounded.

store all mess in transit

Mess delay finite

finite & unbounded

Mess delay bounded if process expects

msg, after time expire
it will sure that it
won't receive it in future

receiver don't receive
for a long time, he is never
be sure whether mes will
arrive

Msg may not lost

synchronous

Msg may not be damaged.

Asynchronous System

no common clock

finite & unbounded msg delay

event-driven. (internal/reception a event)

upon receipt of msg

when condition

Synchronous System

a common clock

bounded msg delay

time - driven

round { receive
do
send

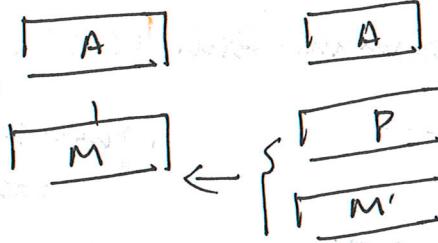
(in reality x exist)

Simulation

① Local simulation

(M', P) looks like M

~ from outside differ



② Global simulation.

to outside observer

there is no difference

Chapter 3 Synchronization (Time Concepts)

Time is needed to assign timestamp to events

compare timestamp to order events
(logical time \rightarrow physical)

in AS, synchronization has to be achieved through msg.

Events in DS { internal
send
receive

P_i = Process

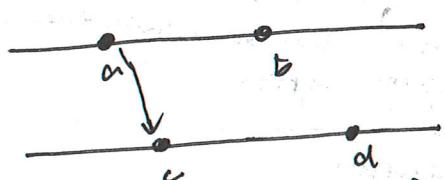
E_i = set of events in P_i

Happen-Before Relation / Causality Relation

local order: if $a, b \in E_i$, a occurs before b .

then

$a \rightarrow b$



$a \rightarrow b$ ①

$c \rightarrow d$ ①

$a \rightarrow c$ ②

message exchange: if a in E_i is sending msg.
 b in P_j is receiving msg.

then $a \rightarrow b$

transitivity: if $a \rightarrow b$, $b \rightarrow c$

then $a \rightarrow c$

Concurrency: if neither $a \rightarrow t$ nor $b \rightarrow a$, then a and b are concurrent (all b)

Causal past (history) of event a.

$$P(a) = \{ b \in E \mid b \rightarrow a \}$$

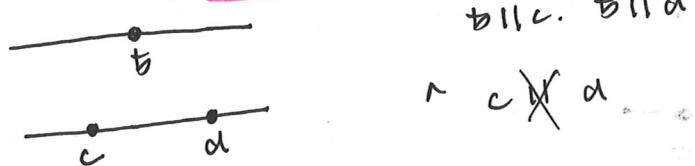
Events concurrent with a:

$$C(a) = \{ b \in E \mid a \parallel b \}$$

Causal Future of an event:

$$F(a) = \{b \in E \mid a \rightarrow b\}$$

concurrency isn't transitive



Logical Clocks

fical clocks
S: partially ordered set (e.g. natural numbers)

logical clock: $C: E \rightarrow S$

if $a \rightarrow b$. $c(a) < c(b)$

converse may not hold!

logical clock characterize HB:

$$c(a) < c(b) \text{ iff } a \rightarrow b$$

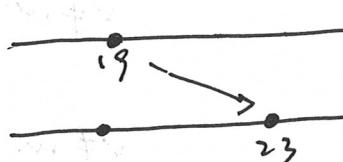
Conditions for Logical clock

1) single P. a, b in E_i . if $a \rightarrow b$, then $C_i(a) < C_i(b)$



2) transitivity

3) may ex.



PI now receiving later
than sending?

Scalar Clock Implementation

1) a in Ei. a isn't msg receive event.

~~then~~
cr(i integer counter in Pi)

① Pi increments ci

② $c(a) \leftarrow ci$

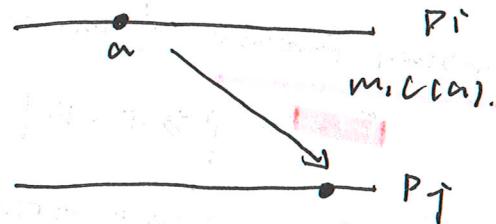
$$\frac{c(a)}{a} \quad \frac{c(b) = c(a) + 1}{b}$$

2) for send event

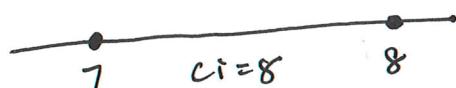
① Pi send $c(a)$ along with m

② Pj assign $c(j) = \max(c_j + 1, c(a) + 1)$

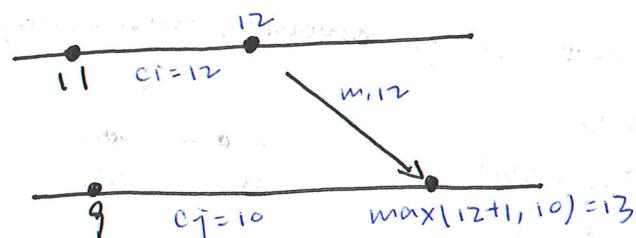
③ Pj assign $c(b) \leftarrow c(j)$



e.g. 1)



2)



Problem: ~~x~~ characterize HB relation.

e.g. 2 P x communicate



5 < 7



but a ~~<~~ b

Vector Clock

comparison of 2 vector v, w of length K

$v=w$ iff $v[i] = w[i]$ $i=1, 2, \dots, K$

$v \leq w$ iff $v[i] \leq w[i]$

$v < w$ iff $v[i] < w[i]$

$v \geq w$ iff $v[i] \geq w[i]$

$v > w$ iff $v[i] > w[i]$

$v \neq w$

$v \neq w$

$$\max(v, w)[i] = \max(v[i], w[i])$$

Vector Clock Implementation

Initial: $V_i = [0, 0 \dots 0]$

1) if a in E_i , a is internal event.

① P_i increments $V_i[i]$ + 1

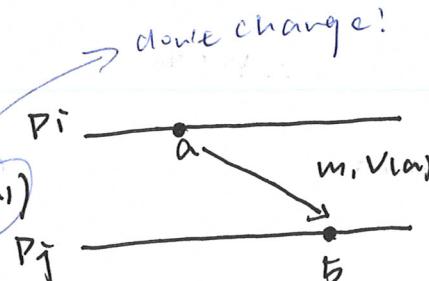
② $V(a) \leftarrow \cancel{V_i}$

2) if a is sending event in P_i
to receiving event in P_j

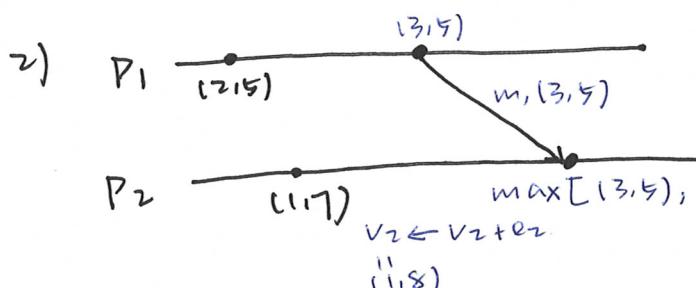
① P_i sends $V(a)$ along with m

② P_j assigns $V_j = \max(V_j + e_j, V(a))$

③ P_j sets $V(b) \leftarrow V_j$



e.g. 1)



Meaning of Vector Clock

1) $V_i[i] = \begin{cases} \text{local event } i = j \\ \text{the last event in } P_j \text{ that } P_i \text{ knows. } i \neq j \\ \text{either directly or indirectly} \end{cases}$

2) all b are concurrent

if $V(a)[i] > V(b)[i]$

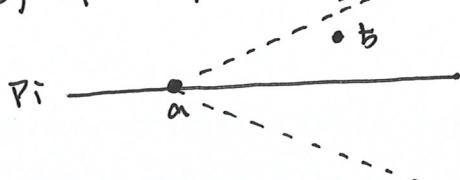
$V(a)[i] < V(b)[i]$

Vector Clock characterize HB

$$V(a) < V(b) \text{ iff } a \rightarrow b$$

1) if $a \rightarrow b$, then $V(a) < V(b)$

2) if $V(a) < V(b)$,



a, b in P_i , $V(a) < V(b)$

a in P_i , b in P_j

b has at least same knowledge about P_i as at event a

Theorem:

For k -dimensional vector clock to characterize HB relation,
with n processes, we need $K \geq n$.

Proof for $n=2$:

Suppose $K=1$. Take 2 p don't communicate.

x characterize ~~internal~~ events.

$$\therefore K \geq 2$$

Proof for $n \geq 2$:

6

$$v_1 = (x_1, y_1)$$

$$v_2 = (x_2, y_2)$$

$$v_3 = (x_3, y_3)$$

$$v_4 = (x_4, y_4)$$

$$v_5 = (x_5, y_5)$$

$$v_6 = (x_6, y_6)$$

$$v_7 = (x_7, y_7)$$

$$v_8 = (x_8, y_8)$$

$$v_9 = (x_9, y_9)$$

$$v_{10} = (x_{10}, y_{10})$$

$$v_{11} = (x_{11}, y_{11})$$

$$v_{12} = (x_{12}, y_{12})$$

$$v_{13} = (x_{13}, y_{13})$$

$$v_{14} = (x_{14}, y_{14})$$

$$v_{15} = (x_{15}, y_{15})$$

$$v_{16} = (x_{16}, y_{16})$$

$$v_{17} = (x_{17}, y_{17})$$

$$v_{18} = (x_{18}, y_{18})$$

$$v_{19} = (x_{19}, y_{19})$$

$$v_{20} = (x_{20}, y_{20})$$

$$v_{21} = (x_{21}, y_{21})$$

$$v_{22} = (x_{22}, y_{22})$$

$$v_{23} = (x_{23}, y_{23})$$

$$v_{24} = (x_{24}, y_{24})$$

$$v_{25} = (x_{25}, y_{25})$$

$$v_{26} = (x_{26}, y_{26})$$

$$v_{27} = (x_{27}, y_{27})$$

$$v_{28} = (x_{28}, y_{28})$$

$$v_{29} = (x_{29}, y_{29})$$

$$v_{30} = (x_{30}, y_{30})$$

$$v_{31} = (x_{31}, y_{31})$$

$$v_{32} = (x_{32}, y_{32})$$

$$v_{33} = (x_{33}, y_{33})$$

$$v_{34} = (x_{34}, y_{34})$$

$$v_{35} = (x_{35}, y_{35})$$

$$v_{36} = (x_{36}, y_{36})$$

$$v_{37} = (x_{37}, y_{37})$$

$$v_{38} = (x_{38}, y_{38})$$

$$v_{39} = (x_{39}, y_{39})$$

$$v_{40} = (x_{40}, y_{40})$$

$$v_{41} = (x_{41}, y_{41})$$

$$v_{42} = (x_{42}, y_{42})$$

$$v_{43} = (x_{43}, y_{43})$$

$$v_{44} = (x_{44}, y_{44})$$

$$v_{45} = (x_{45}, y_{45})$$

$$v_{46} = (x_{46}, y_{46})$$

$$v_{47} = (x_{47}, y_{47})$$

$$v_{48} = (x_{48}, y_{48})$$

$$v_{49} = (x_{49}, y_{49})$$

$$v_{50} = (x_{50}, y_{50})$$

$$v_{51} = (x_{51}, y_{51})$$

$$v_{52} = (x_{52}, y_{52})$$

$$v_{53} = (x_{53}, y_{53})$$

$$v_{54} = (x_{54}, y_{54})$$

$$v_{55} = (x_{55}, y_{55})$$

$$v_{56} = (x_{56}, y_{56})$$

$$v_{57} = (x_{57}, y_{57})$$

$$v_{58} = (x_{58}, y_{58})$$

$$v_{59} = (x_{59}, y_{59})$$

$$v_{60} = (x_{60}, y_{60})$$

$$v_{61} = (x_{61}, y_{61})$$

$$v_{62} = (x_{62}, y_{62})$$

$$v_{63} = (x_{63}, y_{63})$$

$$v_{64} = (x_{64}, y_{64})$$

$$v_{65} = (x_{65}, y_{65})$$

$$v_{66} = (x_{66}, y_{66})$$

$$v_{67} = (x_{67}, y_{67})$$

$$v_{68} = (x_{68}, y_{68})$$

$$v_{69} = (x_{69}, y_{69})$$

$$v_{70} = (x_{70}, y_{70})$$

$$v_{71} = (x_{71}, y_{71})$$

$$v_{72} = (x_{72}, y_{72})$$

$$v_{73} = (x_{73}, y_{73})$$

$$v_{74} = (x_{74}, y_{74})$$

$$v_{75} = (x_{75}, y_{75})$$

$$v_{76} = (x_{76}, y_{76})$$

$$v_{77} = (x_{77}, y_{77})$$

$$v_{78} = (x_{78}, y_{78})$$

$$v_{79} = (x_{79}, y_{79})$$

$$v_{80} = (x_{80}, y_{80})$$

$$v_{81} = (x_{81}, y_{81})$$

$$v_{82} = (x_{82}, y_{82})$$

$$v_{83} = (x_{83}, y_{83})$$

$$v_{84} = (x_{84}, y_{84})$$

$$v_{85} = (x_{85}, y_{85})$$

$$v_{86} = (x_{86}, y_{86})$$

$$v_{87} = (x_{87}, y_{87})$$

$$v_{88} = (x_{88}, y_{88})$$

$$v_{89} = (x_{89}, y_{89})$$

$$v_{90} = (x_{90}, y_{90})$$

$$v_{91} = (x_{91}, y_{91})$$

$$v_{92} = (x_{92}, y_{92})$$

$$v_{93} = (x_{93}, y_{93})$$

$$v_{94} = (x_{94}, y_{94})$$

$$v_{95} = (x_{95}, y_{95})$$

$$v_{96} = (x_{96}, y_{96})$$

$$v_{97} = (x_{97}, y_{97})$$

$$v_{98} = (x_{98}, y_{98})$$

$$v_{99} = (x_{99}, y_{99})$$

$$v_{100} = (x_{100}, y_{100})$$

$$v_{101} = (x_{101}, y_{101})$$

$$v_{102} = (x_{102}, y_{102})$$

$$v_{103} = (x_{103}, y_{103})$$

$$v_{104} = (x_{104}, y_{104})$$

$$v_{105} = (x_{105}, y_{105})$$

$$v_{106} = (x_{106}, y_{106})$$

$$v_{107} = (x_{107}, y_{107})$$

$$v_{108} = (x_{108}, y_{108})$$

$$v_{109} = (x_{109}, y_{109})$$

$$v_{110} = (x_{110}, y_{110})$$

$$v_{111} = (x_{111}, y_{111})$$

$$v_{112} = (x_{112}, y_{112})$$

$$v_{113} = (x_{113}, y_{113})$$

$$v_{114} = (x_{114}, y_{114})$$

$$v_{115} = (x_{115}, y_{115})$$

$$v_{116} = (x_{116}, y_{116})$$

$$v_{117} = (x_{117}, y_{117})$$

$$v_{118} = (x_{118}, y_{118})$$

$$v_{119} = (x_{119}, y_{119})$$

$$v_{120} = (x_{120}, y_{120})$$

$$v_{121} = (x_{121}, y_{121})$$

$$v_{122} = (x_{122}, y_{122})$$

$$v_{123} = (x_{123}, y_{123})$$

$$v_{124} = (x_{124}, y_{124})$$

$$v_{125} = (x_{125}, y_{125})$$

$$v_{126} = (x_{126}, y_{126})$$

$$v_{127} = (x_{127}, y_{127})$$

$$v_{128} = (x_{128}, y_{128})$$

$$v_{129} = (x_{129}, y_{129})$$

$$v_{130} = (x_{130}, y_{130})$$

$$v_{131} = (x_{131}, y_{131})$$

$$v_{132} = (x_{132}, y_{132})$$

$$v_{133} = (x_{133}, y_{133})$$

$$v_{134} = (x_{134}, y_{134})$$

$$v_{135} = (x_{135}, y_{135})$$

$$v_{136} = (x_{136}, y_{136})$$

$$v_{137} = (x_{137}, y_{137})$$

$$v_{138} = (x_{138}, y_{138})$$

$$v_{139} = (x_{139}, y_{139})$$

$$v_{140} = (x_{140}, y_{140})$$

$$v_{141} = (x_{141}, y_{141})$$

$$v_{142} = (x_{142}, y_{142})$$

$$v_{143} = (x_{143}, y_{143})$$

$$v_{144} = (x_{144}, y_{144})$$

$$v_{145} = (x_{145}, y_{145})$$

$$v_{146} = (x_{146}, y_{146})$$

$$v_{147} = (x_{147}, y_{147})$$

$$v_{148} = (x_{148}, y_{148})$$

$$v_{149} = (x_{149}, y_{149})$$

$$v_{150} = (x_{150}, y_{150})$$

$$v_{151} = (x_{151}, y_{151})$$

$$v_{152} = (x_{152}, y_{152})$$

$$v_{153} = (x_{153}, y_{153})$$

$$v_{154} = (x_{154}, y_{154})$$

$$v_{155} = (x_{155}, y_{155})$$

$$v_{156} = (x_{156}, y_{156})$$

$$v_{157} = (x_{157}, y_{157})$$

$$v_{158} = (x_{158}, y_{158})$$

$$v_{159} = (x_{159}, y_{159})$$

$$v_{160} = (x_{160}, y_{160})$$

$$v_{161} = (x_{161}, y_{161})$$

$$v_{162} = (x_{162}, y_{162})$$

$$v_{163} = (x_{163}, y_{163})$$

$$v_{164} = (x_{164}, y_{164})$$

$$v_{165} = (x_{165}, y_{165})$$

$$v_{166} = (x_{166}, y_{166})$$

$$v_{167} = (x_{167}, y_{167})$$

$$v_{168} = (x_{168}, y_{168})$$

$$v_{169} = (x_{169}, y_{169})$$

$$v_{170} = (x_{170}, y_{170})$$

$$v_{171} = (x_{171}, y_{171})$$

$$v_{172} = (x_{172}, y_{172})$$

$$v_{173} = (x_{173}, y_{173})$$

$$v_{174} = (x_{174}, y_{174})$$

$$v_{175} = (x_{175}, y_{175})$$

$$v_{176} = (x_{176}, y_{176})$$

$$v_{177} = (x_{177}, y_{177})$$

$$v_{178} = (x_{178}, y_{178})$$

$$v_{179} = (x_{179}, y_{179})$$

$$v_{180} = (x_{180}, y_{180})$$

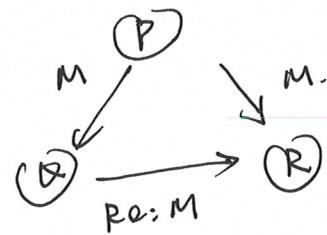
$$v_{181} = (x_{181}, y_{181})$$

$$v_{182} = (x_{182}, y_{182})$$

$$v_{183} = (x_{183}, y_{183})$$

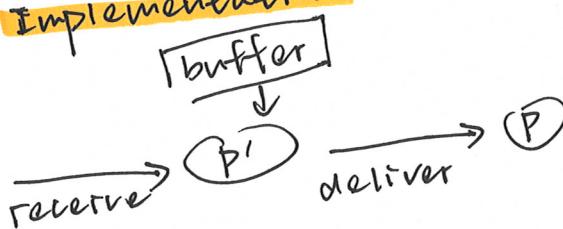
$$v_{184} = (x_{184}, y_{184})$$

Chapter 3 Message Ordering



use clock to enforce ordering of msg.

Implementation



incoming msg \rightarrow buffer

P1 check if it can be delivered

✓ deliver

✗ postpone delivery

$\text{Dest}(m)$: set of destinations

point-to-point: $|\text{Dest}(m)| = 1$

Multicast: $|\text{Dest}(m)| > 1$

Broadcast: $|\text{Dest}(m)| = \# \text{process}$

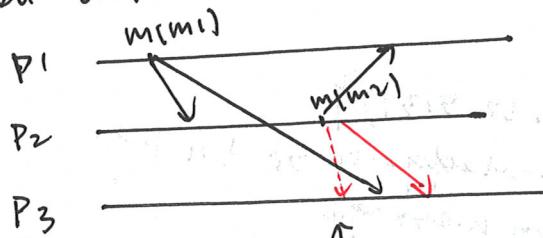
$m(m)$: event of multicasting a msg m.

$d_i(m)$: event of delivering a msg m to P_i

Causal Message Ordering

if $m(m_1) \rightarrow m(m_2)$

then $d_i(m_1) \rightarrow d_i(m_2)$ for all i in $\text{Dest}(m_1)$ and $\text{Dest}(m_2)$



↑
can't deliver,
put in buffer

Birman-Schiper-Stephenson Algorithm (Broadcast msg)

Idea:

1) receive broadcast msg from same process.
inspect number check if miss intermediate msg.

2) msg contains knowledge of process at the point of sending
r check whether missing msg from other processes

Implementations

use vector clock.

send m along with V_m (local vector clock).

condition for delivering: $P_i \xrightarrow{m} P_j$

$$D_j(m) = (V + e_j \geq V_m)$$

P_i was at least as up-to-date as P_j
when it send message

Correctness

- 1) same P .
- 2) diff P .

Causal order: point-to-point = Schiper-Eggle-Sandoz

Causal order: point-to-point = Schiper-Eggle-Sandoz

m_3 : heard P_1 told m_1

P_3 → P_2

m_2 : told m_1 to P_2

P_3 → P_1

m_1 : process identifier

P_1 → timestamp

buffer $S = f(P_j, V_j); \dots$

(local buffer: knowledge about all other process | process at least
(use for yourself) should be that point)

Send along with msg:
update knowledge for receiving process

e.g. in $P_1: S_1 = \{(2, 1, 3, 1, 4), (3, 1, 2, 5, 3)\}$

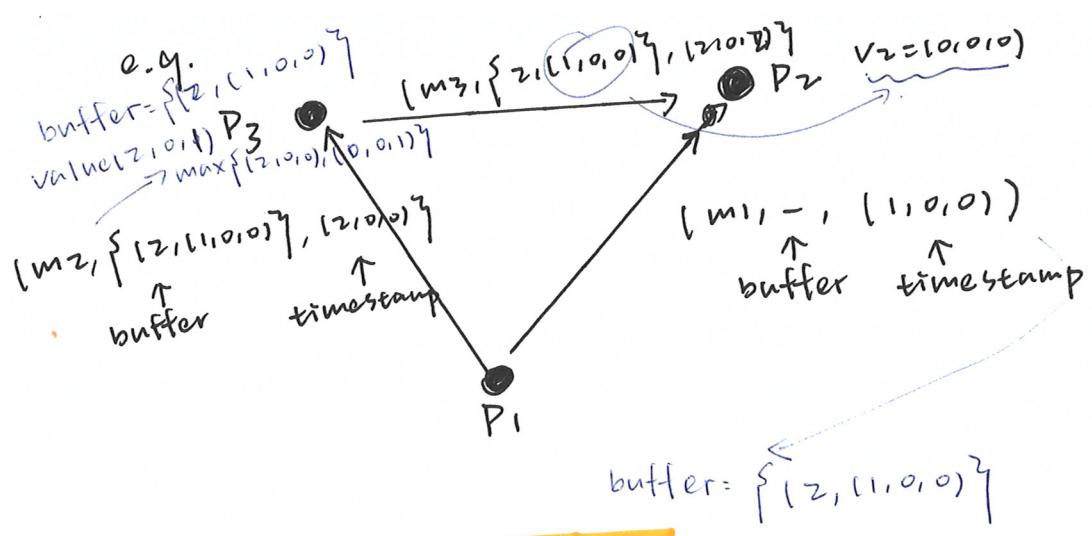
after sending msg, replace knowledge in S_1 for P_j
by the timestamp of sending message

Condition for delivering msg: $D_j(m)$ with accompanying buffer S_m in P_i

$(m, S_m) \rightarrow (P_i)$

① $(i, v) \times \text{exist in } S_m$ inspect etc not for m
i.e. sending process doesn't have any clue at all
the receiving process should be

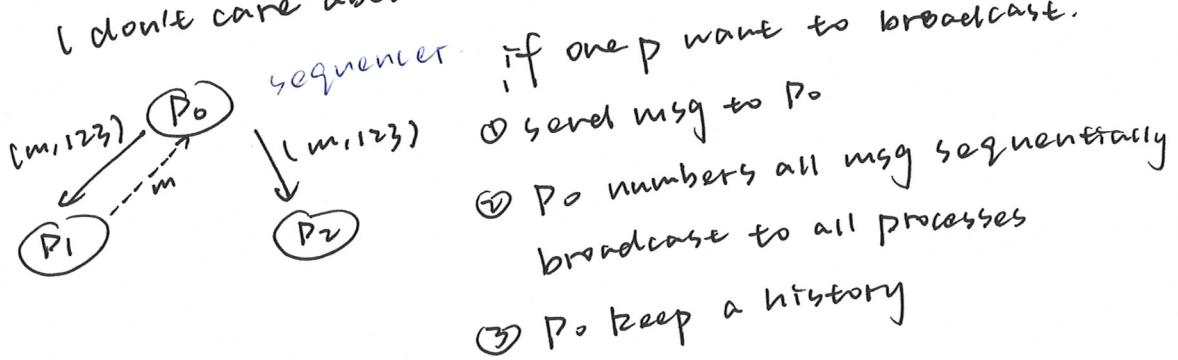
② (i, v) exist in S_m and $v \in V_i$
i.e. at least has up-to-date knowledge for all
components



check own clock for all component at least equal to pair in buffer

Total Message Ordering

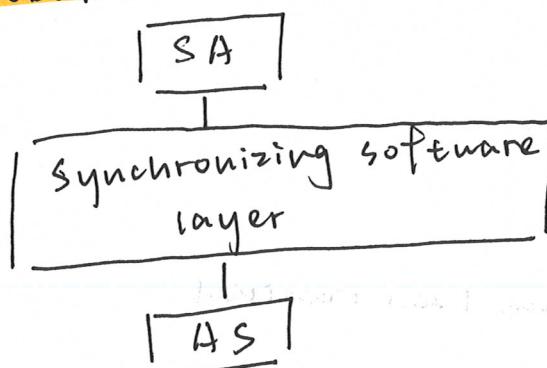
All processes receive all messages in the same order.
 (don't care about order, not necessarily causal order)



17/11/20

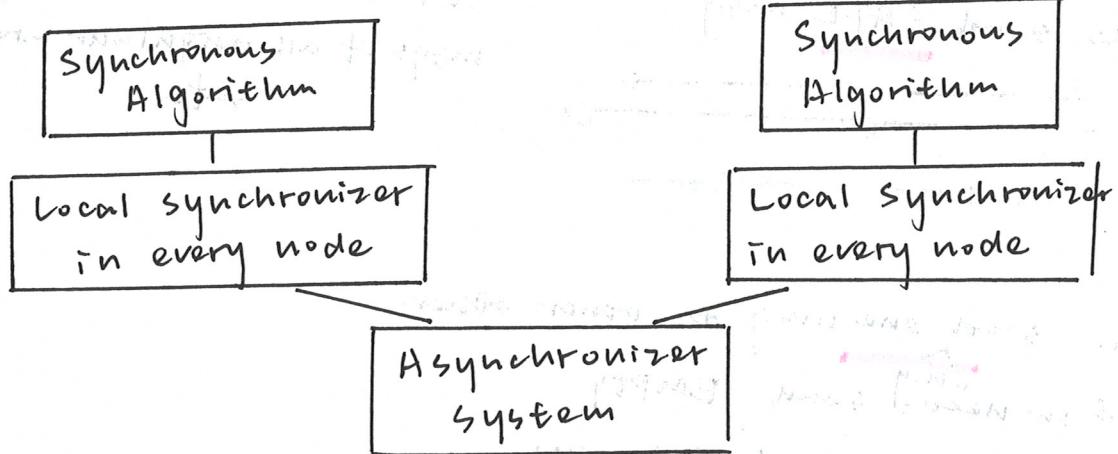
Synchronous systems more powerful than asynchronous system

Global simulation



outsider see a synchronized system

Local simulation



every process locally is a synchronous system
outsider see an asynchronous system

Synchronizer: make AS looks like SS

Problem of AS: no idea when msg arrive (finite ^ unbounded)

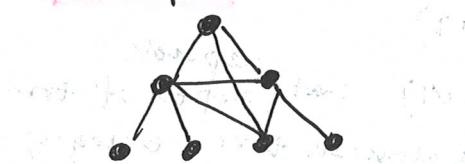
Basic Idea: processes run in rounds

P receive all msg it should receive in current round

then a clock pulse

Types of Synchronizers

α -Synchronizers



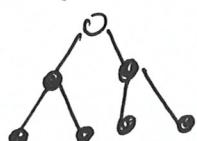
nodes have point-to-point connections
(not necessarily complete)

synchronizer take all msg (n^2)

communication - inefficient

time - efficient

B-Synchronizers



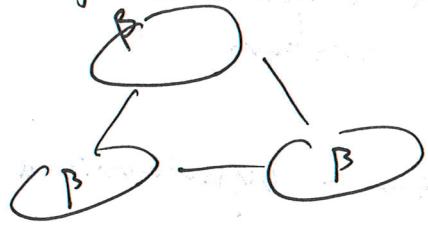
cut msg

create spinning tree

communication - efficient

time - inefficient (msg travel ↑↓)

V-Synchronizer



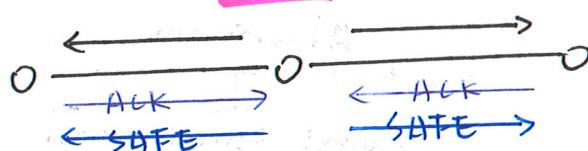
split network into cluster
intra-cluster: β | spinning trees
inter-cluster: α

Generating a pulse
x wait for all msgs

(-) acknowledge

node is safe: all msgs it sent has been received
detect by: reception of ACK.

safe node send SAFE msg



a node knows it receives all msgs if all neighbour are safe

(=) tricky

every node send one msg to other other.

{ if no need: send only EMPTY

{ if multiple: pack into one

a node should in every pulse receives single msg
from neighbor

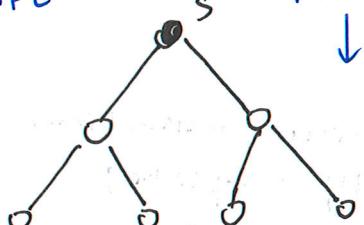
B-Synchronizers Impl

Initial: ① elect a leader

② create spinning tree

SAFE

PULSE



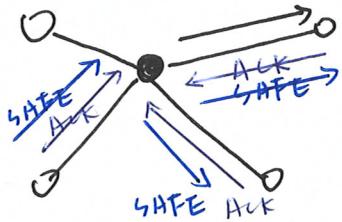
when a node finds itself & all descendants safe
it sends SAFE to its parent
when root receives SAFE, send PULSE down
when node receives PULSE, it can go to next round

complexity: Communication, $C(\beta) = D(V)$

Time:

$T(\beta) = O(V)$ depends on depth of tree
balanced tree: $O(\log n)$

α -Synchronizer Impl

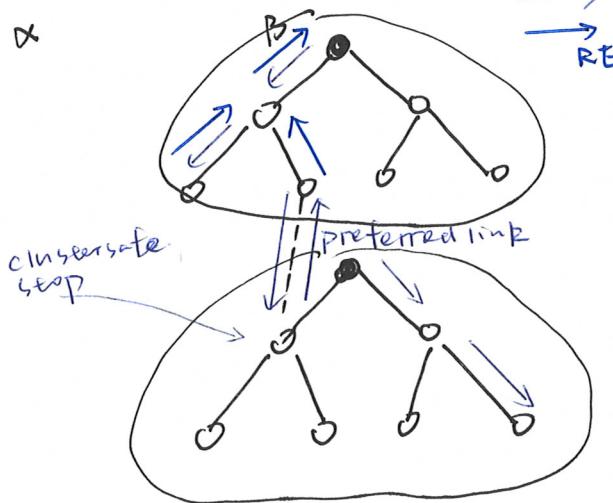


A node:

- find itself safe if receives all **ACK**
- send **SAFE** to all neighbors
- receive **SAFE** from neighbors
- generate local pulse

Complexity: communication: $C(\alpha) = O(|V|^2)$
time: $T(\alpha) = O(1)$ { 3 steps = constant time }

γ -Fmannal-Synchronizer Impl



Initial:

① partition nodes into clusters

② select leader in cluster

③ create spanning tree in cluster

④ create **preferred link** between clusters

β -Synchronizer: (1st wave)

- every node broadcast **PULSE** down
- convergecast **SAFE** upward

α -Synchronizer: (2nd wave)

when root knows whole tree is **SAFE**

broadcast **CLUSTER-SAFE**

down whole tree

 preferred link

② node converge cast **READY** upwards

when receive **READY** from all descendants

and **cluster-SAFE** along all **preferred links**

17/11/22 Global State

Problem: determining a global state in asynchronous system
(processes, ticks, msgs)

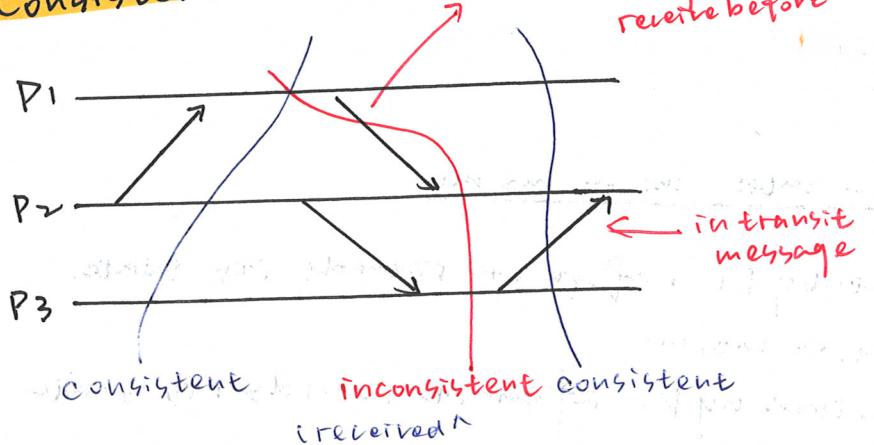
Applications: debugging applications (distributed)

detect local stable properties (deadlock)

e.g. transferring money in bank account

freeze the money in your account & in transit

Consistent cut



cut is a set of one event for every process $\{c_1, c_2, \dots, c_n\}$

consistent cut:

no e_i in P_i & e_j in P_j

that $e_i \rightarrow e_j$ & $c_i \rightarrow e_i$

& $e_j \not\rightarrow c_j$

use vector clock $V(c_i)$ of event c_i

$V(c)_{ij} = \max(V(c_1)_{ij}, V(c_2)_{ij}, \dots, V(c_n)_{ij})$ component-wise

↑
time of cut

c_i is consistent iff

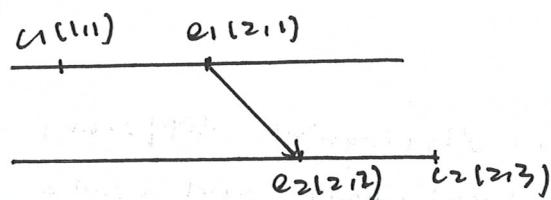
$V(c) = (V(c_1)_{ij}, V(c_2)_{ij}, \dots, V(c_n)_{ij})$

local component is maximum

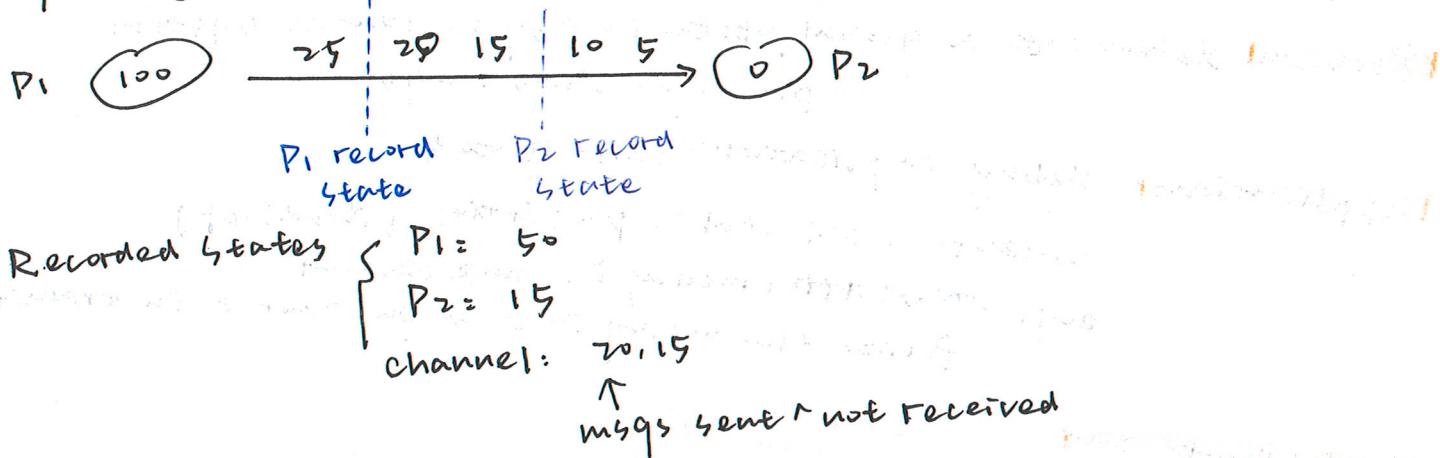
$$V(c) = \max((1, 1), (2, 3)) = (2, 3)$$

$$(V(c_1)_{ij}, V(c_2)_{ij}) = (1, 3)$$

e.g. violation

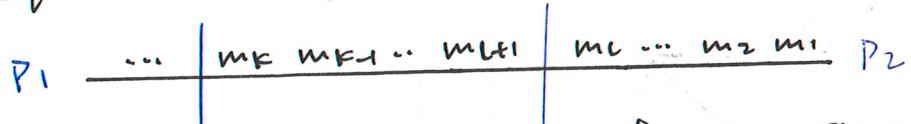


e.g. money transfer



State of channel (FIFO)

a sequence of msgs



state of c = msgs sent by P₁ before it records its state.

m_k, m_{k-1}, ..., m_{l+1}

= msgs received by P₂ after it records its state

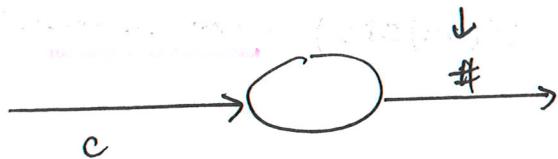
... m_k, m_{k-1}, ..., m_{l+1}

$k > l$: state of c = [m_k, m_{k-1}, ..., m_{l+1}]

$k = l$: empty

$k < l$: inconsistent

Chandy's and Lamport's marker



record local state

blue line

① record local state

② for outgoing channel: send #

incoming

create buffer: msg after recording

state & before receiving marker

wait for marker on incoming channel

Assumption: channel is FIFO

may spontaneously start

append (pid, seq-id) to # marker to distinguish different invocations

After algorithm, each process have local states and state

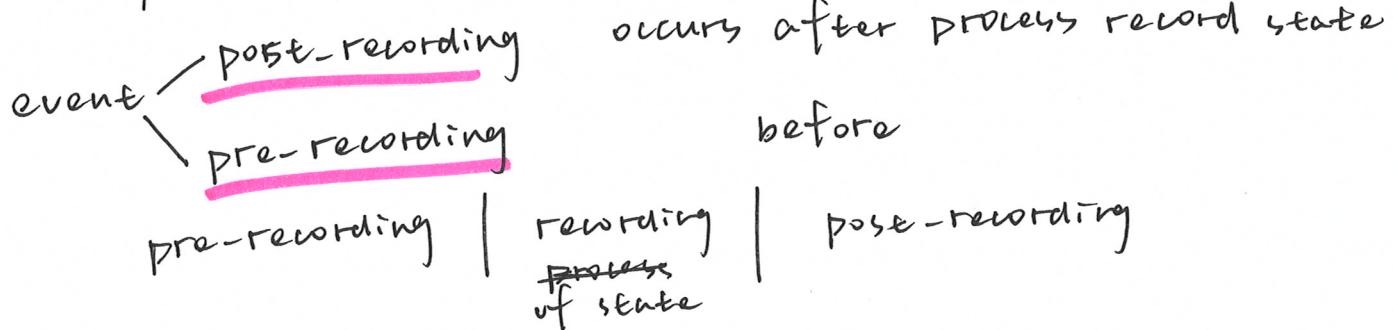
for incoming channel

Global state can be collected

Recorded State

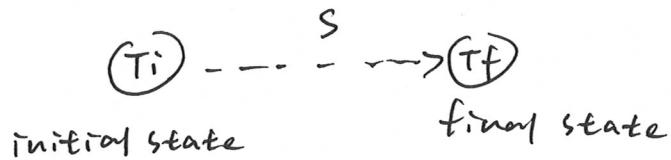
Q: what is the use of state if it hasn't occurred?

Might occur if sequence of concurrent event is different.



If want to record global state

S = sequence of events



$s = \alpha \alpha \alpha \alpha \alpha \alpha | \alpha \alpha \alpha \alpha$
↑ ↑
Pre-recording Post-recording

For multiple ~~even~~ processes, sequence might be out of order.

Reorder events

order events post pre
↑ ↑
S = e e e e op eq e e e

- ① esp. esp. both send/receive events
related v.

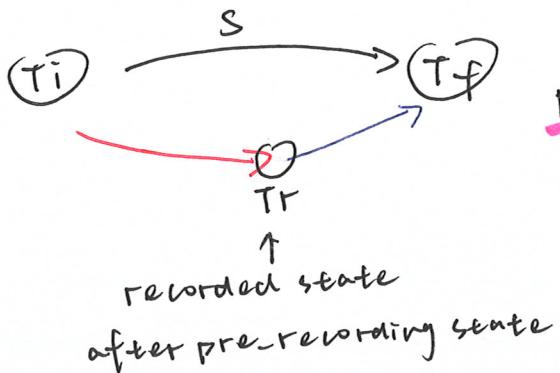
⑦ QP send. QP receive. → cause interchange.

but it's impossible

~~QF~~ = QF = post-rewording (x)

(P) $\xrightarrow{m \#}$ (Q)

reordered states of processes are still same (\vdash : same events)
used to check
s. original \Rightarrow sealable properties



7/11/22

Global State

Termination Detection

Problem: detect the application is terminated

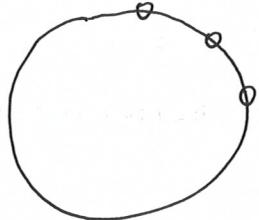
A distributed application consists process in different processor.
Some P finish spontaneously as done all local work.
active \rightarrow passive

Passive nodes may become active again if receive msg carry
additional work.

problem: some messages are still in transit
may re-activate process

Termination Detection in a unidirectional ring. (FIFO links)

Naive Solution:



Token: travel along ring

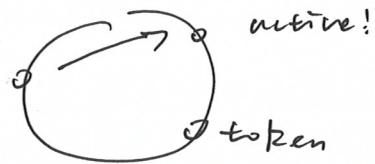
capture the knowledge active/passive
/re-activated?

when token hits process
passive. \rightarrow continue

active \rightarrow wait until becomes passive

when come back to originating process, report termination

Problem: process can be re-activated



Solution:

introduce color (black, white) for token and process

Process turn black if send message to a higher numbered process
(send from left to right)

\rightarrow spot the passiveness

Token: start with white

turn black if arrives at black process

(continue black to initiating process)

when return \swarrow black = no conclusion
white

correctness

$loc_t = \text{location of token}$

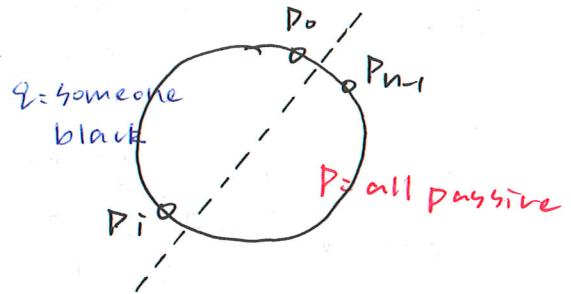
$loc_t = i$ if { token in P_i

token in transit. P_i is the last process token present

Predicates:

P : for all i , $1 \leq i < n$, $\text{state}_i = \text{passive}$

Q : There exist an i , $0 \leq i \leq n-1$,
 $\text{color}_{-Pi} = \text{black}$.



Γ : $\text{color}_{-t} = \text{black}$.

At least 1 predicate is true.

Proof: $P \sqcap Q \sqcap \Gamma$ is invariant.

① If token in location 0 (white)

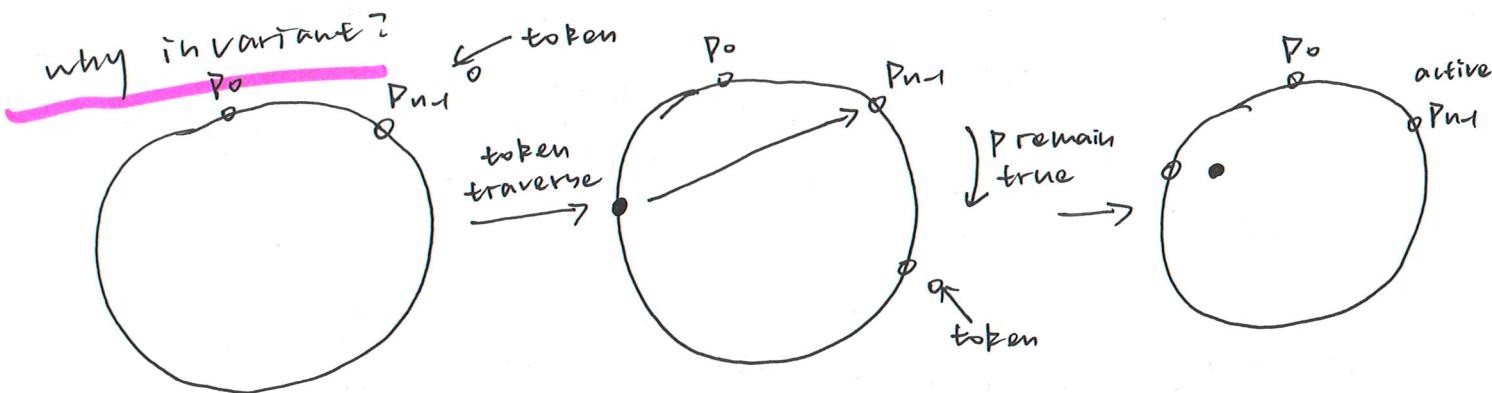
$$\text{Xnew} = P_0.$$

$$\vdash Q \times$$

$$\Gamma \times$$

P must hold (means all process on the ring is passive)
+ $P_0 = \text{passive} \Rightarrow$ conclude termination

why invariant?



$P = \{0 \leq i < n \mid \text{empty set}\}$

\uparrow
 $n-1$
empty set is
always true

$$P = X$$

$Q = V$
(token travel and pass P_0).
 P_0 also know about sport.

$$P = X$$

$$Q = X$$

$$\Gamma = V$$

General termination detection

No special structure for network.

All processes and messages have non-negative weights

One special process

$$\xrightarrow{\text{initialization}} w(P_i) = \frac{1}{n}$$

Initial weight

when process terminates

(P_i)

(P)

special process = 1

all other = n .

when send a message

before w_i

$$m(w_i)$$

w

before w_i

$$\frac{w_i}{2}$$

$$w_j$$

$$w_j + \frac{w_i}{2}$$

after 0

$$w + w_i$$

Termination Condition: $w(P) = 1$

Problem:

Each process request one resource.

They are waiting for each other because of resource request.

Resource model

Every resource is represented by process = process RA represents resource A.

Processes are waiting for each other.

Process P requests Resource A



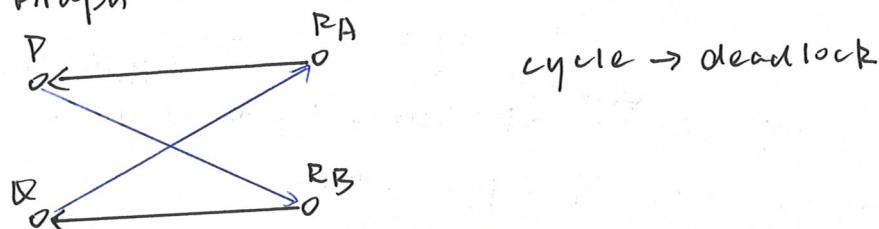
After request has been granted:



After resource has been released:



Wait For Graph

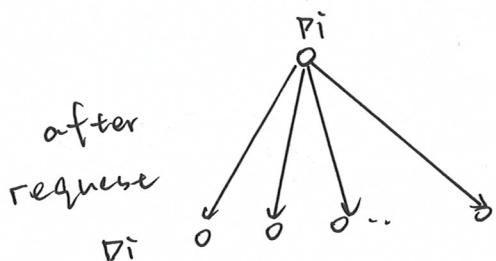
**Communication Model**

Idea: Process blocks and waits for one message from bunch of msgs.

P_i

D_i : dependent set

A Process turns active when receives a single message from any message in its dependent set.



After reception of a single message D_i

$D_i = \{0, 0, \dots, 0\}$

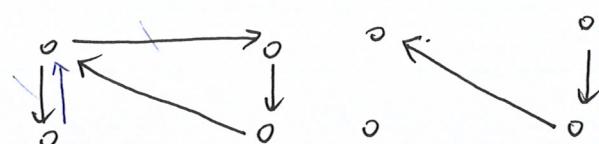
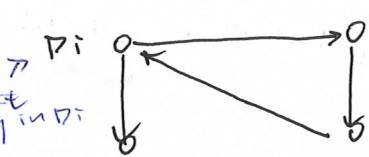
Deadlock and WFG graph

Resource model: a cycle



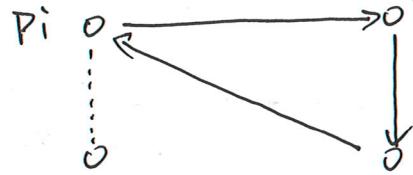
it receive msg in D_i
all arrow gone

communication model:
a cycle doesn't necessarily indicate a deadlock



Knot: A set of nodes in directed graph

with a path from every process to every other process (not need complete graph.)
no edge from a node in S to a node not in S



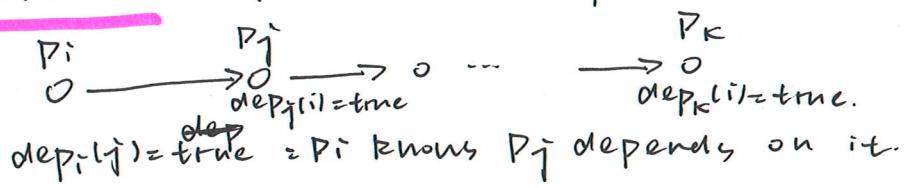
There is a deadlock if there is a knot in WFG.

Chandy - Misra - Haas: Resource Model

Idea: when a process suspects deadlock (wait for long time), it sends a special message to all processes it waits for.
(PROBE)
other processes propagate message

Deadlock condition: if receives one of own messages

P_i is dependent on P_j if there is a path in WFG from P_i to P_j



$\text{dep}_{ij} = \text{true} \Leftrightarrow P_i \text{ knows } P_j \text{ depends on it.}$

P_i is deadlock if dep_il = true

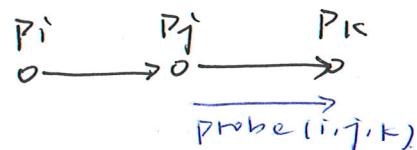
Implement.

Probe msg: P_j sends probe(i, j, k)

if P_j is blocked

P_j waits for P_k

P_j knows P_i is dependent on it.



Chandy - Misra - Haas: Communication Model

Idea: when a process suspects deadlock, it sends a special message to all processes it is waiting for

those processes propagate these msgs.

when a process already receive the msg, send reply immediately.

when a process exhaust all links, sends back a reply.

Deadlock if receives a reply from all processes it send query to

Query / Reply Msg (i, m, j, k)

i: initiator of the query

m: query sequence number

P_j sends msg to P_k

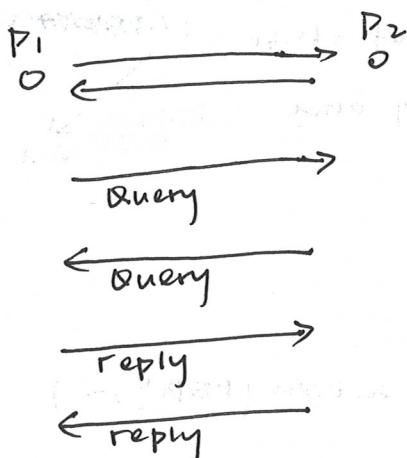
P_i → P_j → P_k

Engager: the node from which the process receives first Query msg.

whenever a node turns active (resources are granted), the whole algorithm stops

A process initiates a query is deadlocked if and only if it receives a REPLY to every Query it sends.

Example



Bracha and Toney: Requests

N-out-of-M request

A request is targeted at M processes

A process can be satisfied with a lower number N.

OR: N=1

AND: N=M

Example

in distributed database, transaction runs on different location

1) AND: when commit a transaction, needs agreement from all component.

2) N-out-of-M: quorum-based replication ($N > \frac{M}{2}$)

Idea:
Any process that suspects a deadlock can initiate the algo.

Algo consists of 2 phases

{ notify state of algo (wave-out)
simulate granting of resources

when initiator remains blocked after exec, it is deadlocked

Processes
active could request and then blocked
blocked wait for request to be satisfied

Active process can do N-out-of-M request.
2-out-of-3 require

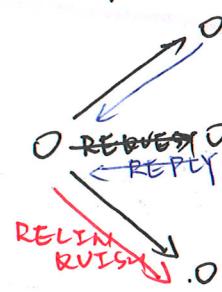
1) send a REQUEST to all M

2) become blocked

3) wait for positive REPLY from N of those

4) send RELINQUISH to remaining ones

5) become active



FIFO Channel

WFA

A change in WFA correspond on an action (REPLY, ...)

Schedule: sequence of operations

Scheduled applied to $G = G - G'$

A node is deadlocked in G if there is no schedule that takes G to G' in which v is active

Static System: msg are instantaneous

In set = OUTset (\because no msg in transit)

Transformation in WFA

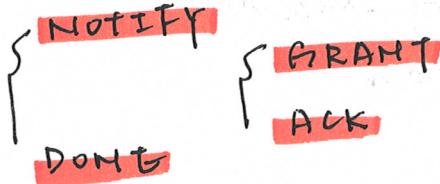


① Add K outgoing msgs

$\rightarrow r$ -out-of- K Request

② Receive REPLY \rightarrow delete an edge

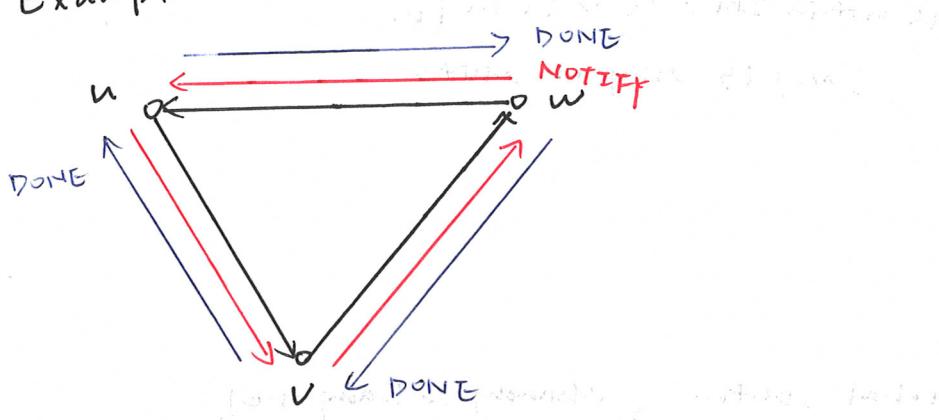
Bracha and Toneg: With instantaneous messages



Initiator is not deadlocked if $\text{free} = \text{true}$

(in other algo. one is deadlocked \rightarrow others also deadlocked in communication model)

Example



Bracha and Toneg with messages in transit

use color to represent links

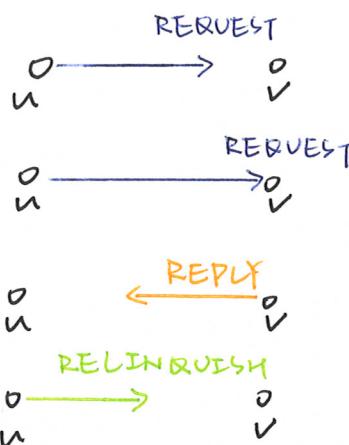
color $(u, v) =$

gray:

black:

white:

translucent



grey edges
will eventually
become black

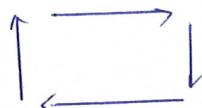
Idea: ignore grey

white & translucent will disappear if msg is received

map colored WFG to colorless WFG

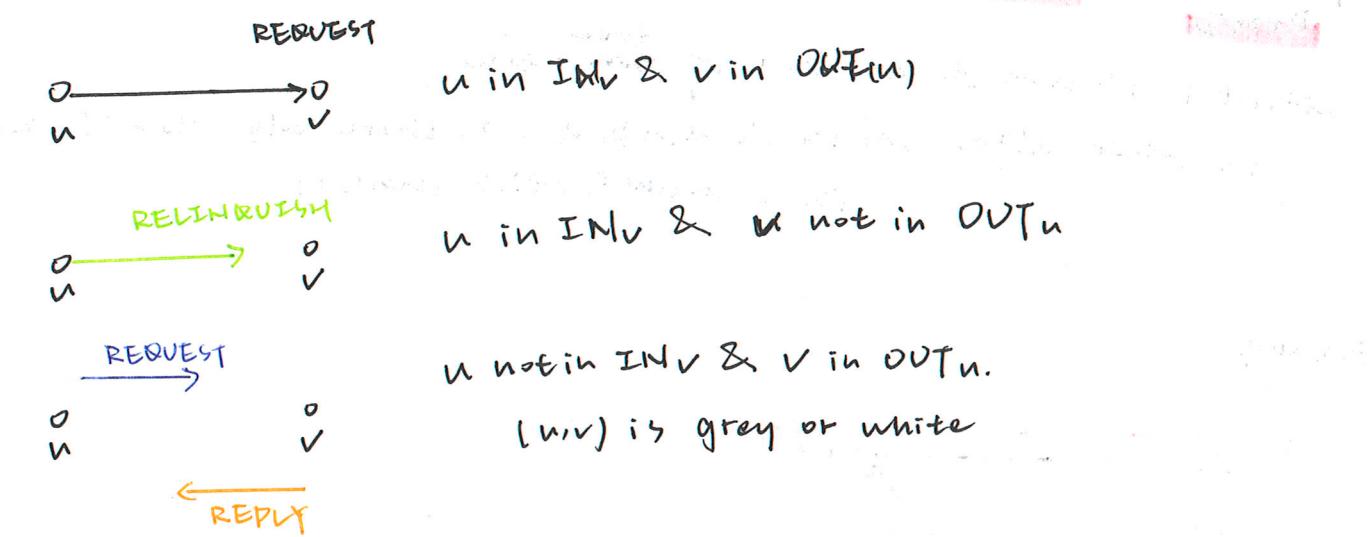
A node is active if $n \leq \# \text{grey write}$

can detect a deadlock about to happen



REQUEST

Nodes need to know color of edges
by exchanging COLOR messages containing IN & OUT set
compare IN & OUT set → figure out color



DYNAMIC SYSTEMS

No static WFG

Trick: detect a global state. (Chandy & Lamport).
apply algo for static system

17/11/29 Mutual Exclusion

Problem: grant exclusive privilege

Requirements: no deadlock
no starvation
fairness

Most algorithms are centralized

Lamport's mutex algorithm

Channels are FIFO. use scalar clock (+ process id)

Each process maintains request queue (ts is ordered)

P_i request access to CS
broadcast (REQUEST, t_s, i) to all P_j (P_i incl.)

P_j receive request

send (REPLY, t_s')

put (t_s, j) in request queue

enter CS condition

receive REPLY from all

request at head

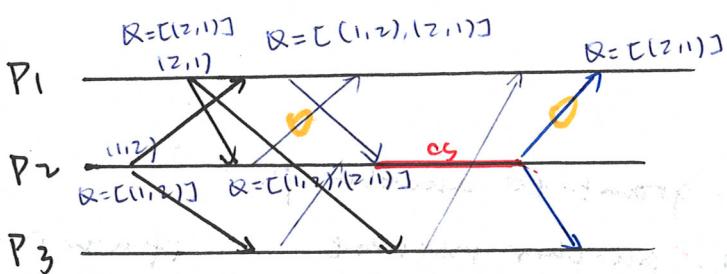
P_i Release CS

send (RELEASE, r) to all (P_i incl.)

Receive Release from P_j

remove (t_s, j) from request queue

RQ is equal



REQUEST

REPLY

optimization

$3(n-1)$

suppress unnecessary REPLY

links FIFO needed

Ricart-Agrawala's mutex algorithm

Optimization of Lamport: combine RELEASE and REPLY

defer REPLY when request has lower priority

links FIFO needed?

Morikawa's mutex algorithm

Idea: process sends request to subset of msg

Request set

reply: give exclusive access to CS

don't give another REPLY to other process unless receive RELEASE

Feature of request sets:

request sets has same sets

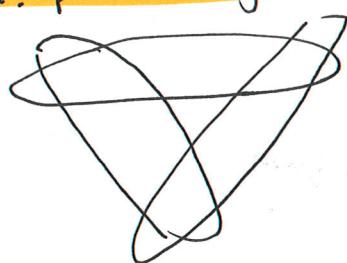
Share the load

every P is contained in the same number of request sets

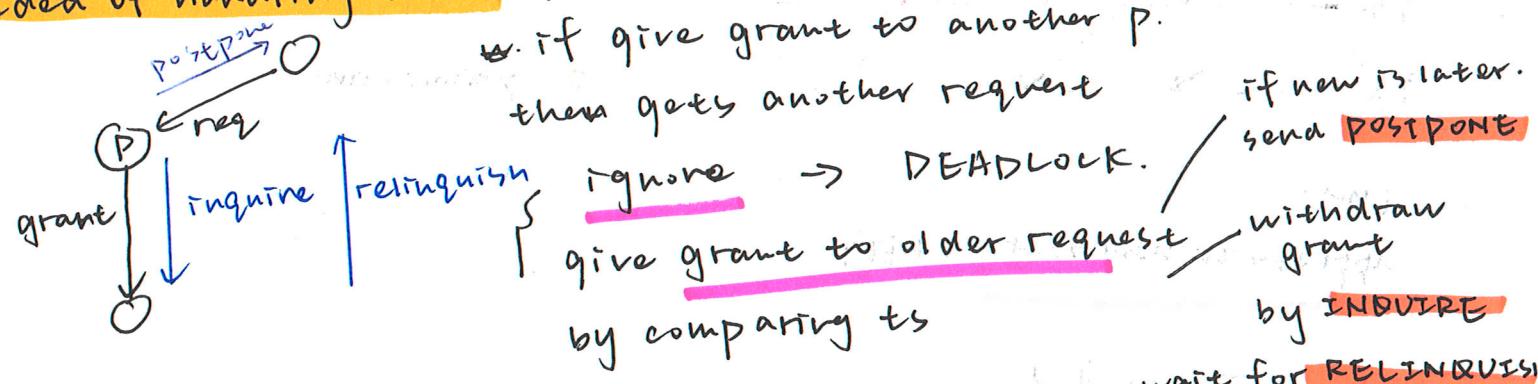
any two request sets have a non-empty intersection
(grant permission to only one p).

every p in own request set

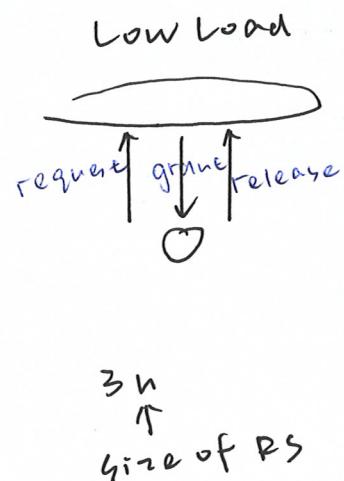
Problem: possibility of deadlock



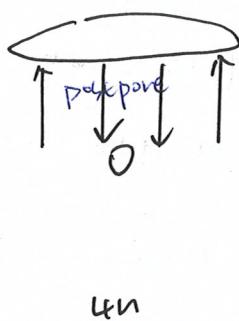
Idea of handling deadlock



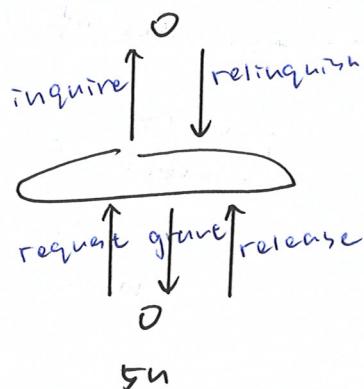
Complexity:



Heavy load



old request

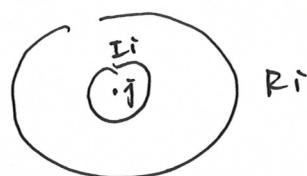


Generalized mutex algorithm

common structure for assertion-based algorithm

each process maintains 3 sets:

{ Request set:
inform set:
status set:



RELEASE: not sent to everyone in R_i.
to a smaller set I_i^r

P_i inform P_j



P_j records status of P_i

Lemma: minimum size of request set is $\lceil \frac{D}{K} \rceil$

K: size of request set

D: the time of every process appear in request set

$$\text{Process number} = \frac{\text{number of req set} \times \text{size}}{\text{number of appearance}}$$

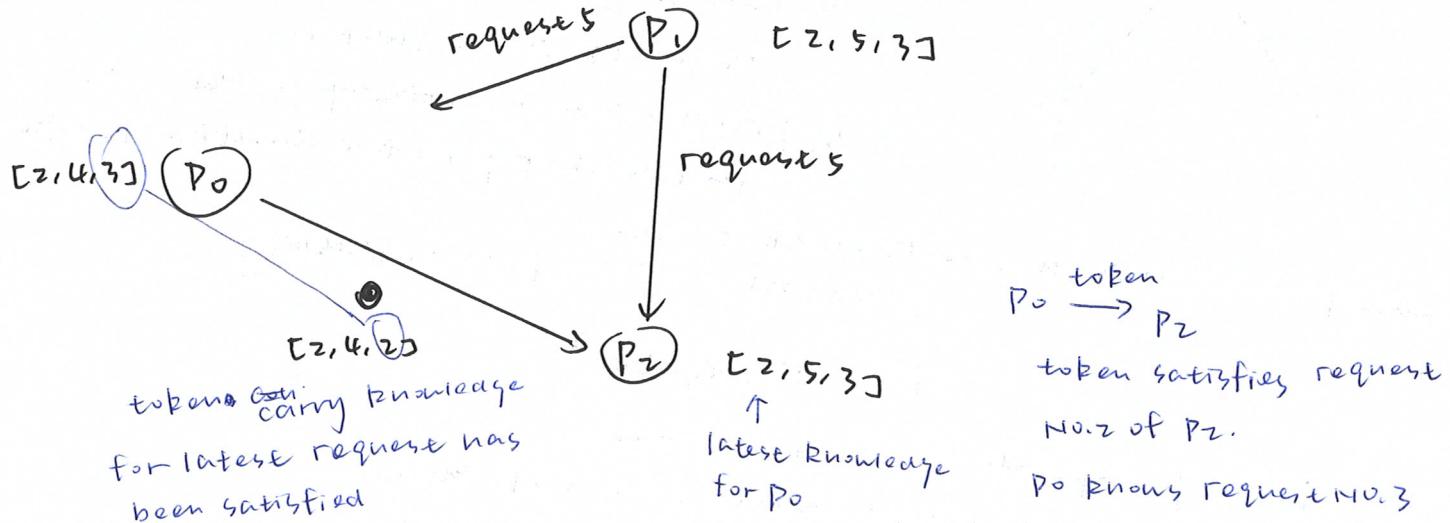
$$n = n \times \frac{K}{D} \quad K=D$$

Suzuki-Kasami's mutex algorithm

token-based: a p can enter CS when possesses token

Idea: number request sequentially
broadcast request

token transfer to knowledge & to req has been satisfied



Problem: process with high number may get starvation
as token may be circulated in low number of processes

Variation: In token maintaining a queue

Singhal's mutex algorithm

Improvement: not send request to all processes

send to p which may possess token

Implementation: each p maintain state array $s[1 \dots n]$

possible states: { R: requesting token
E: executing CS
H: keeping token & outside CS
(don't know where to go)
O: other }

Initialization:

suppose token is located in P_1

in P_1 : $s[1] = H$.

$s[j] = O$. $j = 2 \dots n$.

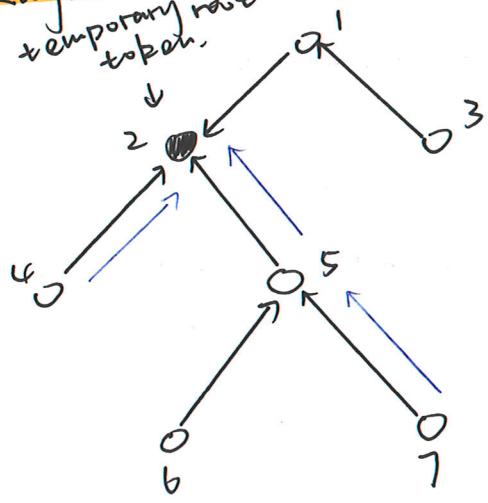
in P_i : $s[i] = R$. $i = 1 \dots n-1$

$s[j] = O$. $j = i+1 \dots n$

token carries T_N and T_S
 \uparrow request num \uparrow process states

$T_N[i] = 0$. $T_S[i] = 0$

Raymond's token-based algorithm



b) each process maintain a request queue
 { own ID
 | id of neighbor

p is requesting

- Idea: "suppose a spanning tree"
- 1) process holding the token is root
 - 2) process sends request in the direction towards root.
 - 3) let a single process represent all requests in a subtree
 - 4) process has a notion of holder, which is the id of p in the direction of token

Msg complexity: depends on depth of tree
 token does a tree traversal: travel each link twice.
 every token step is preceded by a request step

17/2/14 Election

Problem: every process has an integer.

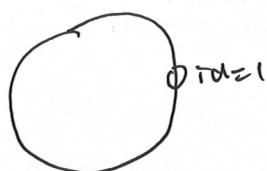
Election a process with largest/smallest integer.

Trivial solution: every process sends id to every other process

Anonymous network: nodes initially don't have ids
Election is impossible in anonymous network

Comparison-based algorithm

Non-comparison-based algo in synchronous ring



all process know current round number

round 1: every process inspect own id = 1
if so, send msg along ring

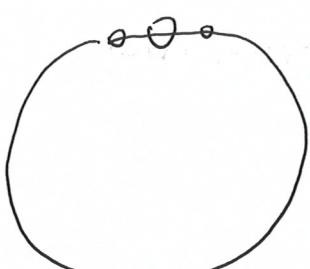
if no. in first n round, no msgs.

round n+1: $ID = 2?$

msg / time complexity: $O(n)$.

Compare round number
& id.

Bidirectional ring (asym) $O(n \log n)$



Idea: a process check if its id is largest among two neighbors. — if so, survive
— if no. stop

in round K, check 2^{K-1} process on both sides.

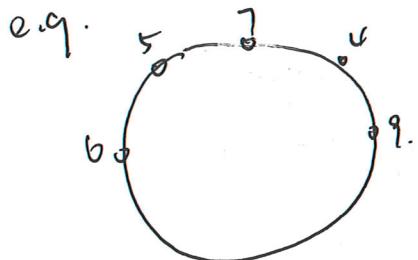
-
- A diagram of a process node with three outgoing arrows. One arrow points left with a label "probe", another points right with "id = 1", and a third points right with "hop count". A fourth arrow points right with "discard". Below the node, four numbered steps are listed:
- 1) own id larger
 - 2) own id smaller
 - 3) hop count > 0.
propagate msg. with hop--
 - 4) reply ok.

A process starts phase $K+1$, if receives $2^{phase-K}$ OK msgs.

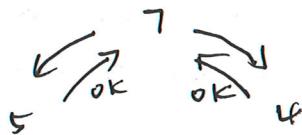
A process is elected. — don't know size of ring.

① receive own Probe msg.

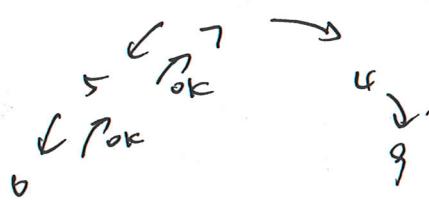
② 1 process receives same id from either side
stop at K round, K is $2^{phase} - 1$.



Round
phase 1:



phase 2:



Bidirectional = solution 2

in first round:
exchange ids with 2 neighbors

a process remain active if own id is larger
else passive.

every next round: repeat round 1 in virtual ring.

A process is elected if receives own id.

msg complexity: $2n \log(n)$.

Unidirectional ring: Peterson's

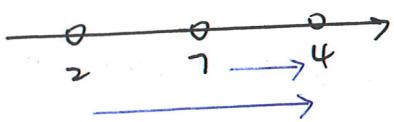
A simulation in unidirectional ring of solution 2 in bidirectional ring.
 $ttd=id$.

nttd: upstream neighbor

nttd: $\max(ttd, nttd)$

if $nttd > ttd$, $nttd = nttd$, remains active

$ttd=nttd$



else turn passive

every subsequent round, repeat round 1 in virtual ring

A process is elected if receives own id.

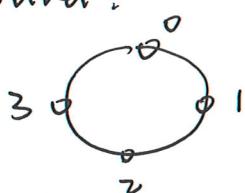
Active nodes will cut half every round,
at least

msg complexity: $2n \log(n)$ (at most)

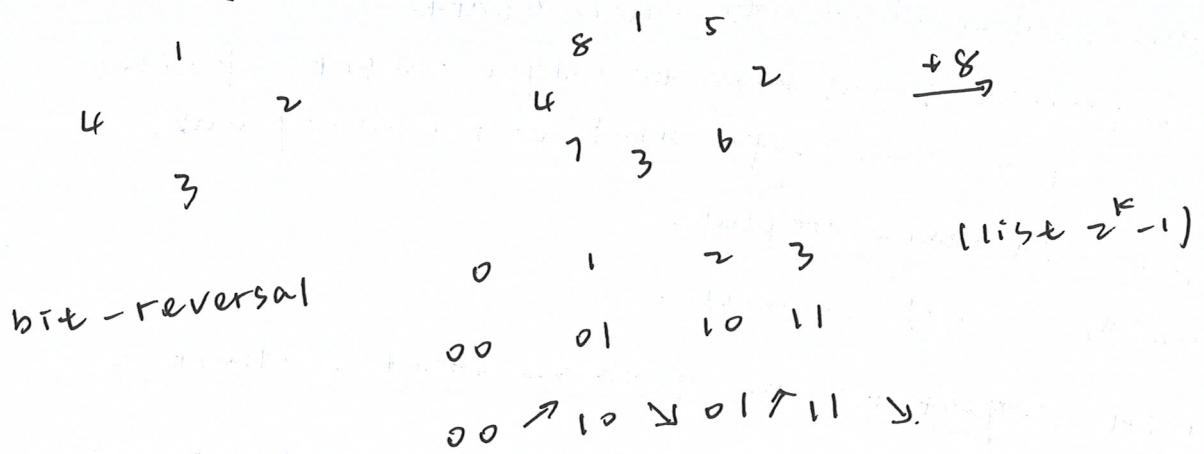
$2n$: in every round 2 msgs sent along one link
 $\log(n)$: number of rounds at most

Q1: when algo terminates after one round?

Increasing id. along ring.
(decreasing)



Q2: Ring of size $n=2^k$ use k rounds?



17/12/06

Assertion: $\overbrace{P_i}^{\text{if id of } P_i \text{ still survives in some active process } P_j, \text{ all processes between } P_i \text{ and } P_j \text{ are passive.}}$

Time complexity: $2n-1$

Complete Network Synchronous

straightforward: every process sends its id to everybody else.
 n^2 msgs.

Afek's and Gafni's

Idea: successively send ID to ever larger sets of processes
and wait for Ack. from all of them

If a process receives an ID larger than its own, send Ack.

A process is killed if it doesn't receive all Ack's it expects.

A process adopts the largest ID it sees. \downarrow syn

Alive candidate process send msg to larger subsets in successive rounds.

candidate process

Node { ordinary process = awake by reception of msg
(and doesn't want to be elected.)
spawns only spawns ordinary process

candidate: send msgs to remaining P.

ordinary : receives msgs from candidate process and compare.

msg: (level, id)

level = the number of rounds since start

candidate process sends msg to larger subsets of sizes powers of 2, and keeps track of remaining links

msg compared in lexicography =

(level, id) $(2, 5) > (1, 300)$

A node with largest id among those start earliest wins.

A node is elected if receives ACK from all other processes.

levels are not strictly needed.

Q: disadvantage of not using levels?
time may ↑.

Time complexity: $\log(n)$

Msg : $n \log n$

Afek's and Gafni's asynchronous algorithm

It doesn't make sense for a node to wait for all msgs.

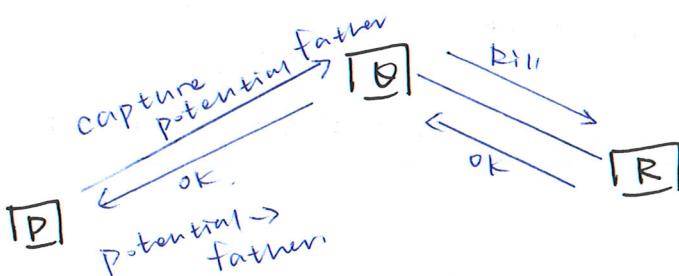
It reacts once receives a single msg.

If one node captures a node previously captured by another node, that node may also to be captured.

A node keeps two special links

father = link to owner

potential father: to potential new owner



Time complexity: $O(n)$

Msg : $n \log n$

Exercise Afek's and Gafni's

d. all processes start algorithm in round 1.

1st round. $8 \rightarrow 0$

$9 \rightarrow 1$

:

$15 \rightarrow 7$.

2nd round $12 \rightarrow 8 \quad \{8, 0\}$

$13 \rightarrow 9 \quad \{9, 1\}$

:

$15 \rightarrow 11 \quad \{11, 3\}$.

3rd round. $14 \rightarrow 12 \quad \{0, 8, 4, 12\}$

$15 \rightarrow 13 \quad \{1, 5, 9, 13\}$

4th round $15 \rightarrow 14 \quad \{0, 2, 4, 6, 8, 10, 12, 14\}$

Chp 4. Minimum Weight Spanning Trees

Problem: a weighted undirected graph $G = (V, E)$
make a selection of links that connect the network
completely
and total weight of links is smallest among all possibilities

Difficulty: compute in distributed way

Lemma: There is a unique MST.

Assume all weights are different

Proof: suppose \geq MSTs: T_1 and T_2
there is some edge e that is part of
 T_1 but not T_2

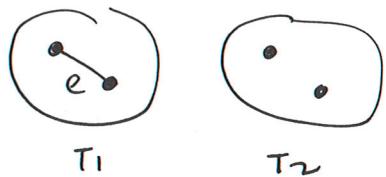
take one with lowest weight: e

add edge to T_2 then gets a cycle

T_1 has no cycles. $\Rightarrow T_2$ contains an edge e' that is not in T_1 .

$$w(e') > w(e)$$

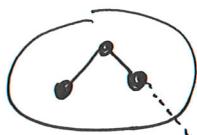
throw e' and keep e . T_2 has a lower sum of weight
than original T_2 .



Notion:

fragment: subtree of final MST.

outgoing edge: one node in fragment and
one node outside.



MOE | Minimum-weight outgoing edge

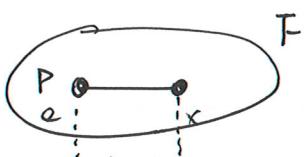
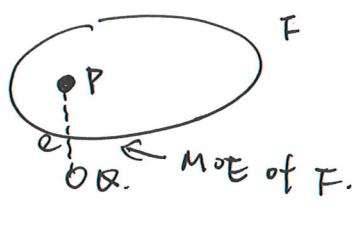
Lemma: If fragment F . MOE. e .

F with e and node incident with e is also a fragment.

Proof: suppose not.

so in final MST there must be path
connects P and X .

so there is a cycle if add e to
final MST.



- at least one edge x is also outgoing edge of F .

$\therefore w(x) > w(e)$.

replace x by e yield an lower weight MST.

\therefore contradiction.

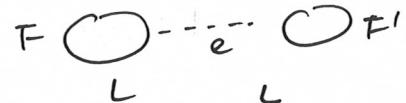
Notion:

level: fragment has level.

Fragment with single node has level 0.

Merge = two fragments with same level.

have common MOT, level



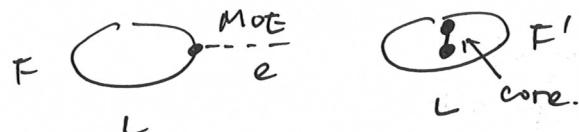
yield a new fragment whose ^{merge of} name is the core

↑
last edge used for merge

level increased by 1.

Absorb = a fragment of a lower level into fragment of a higher level

level stay same.
core same
name same



(Q: Minimum size of fragment level k ?

binary tree

2^k .

\geq level 0 make level 1

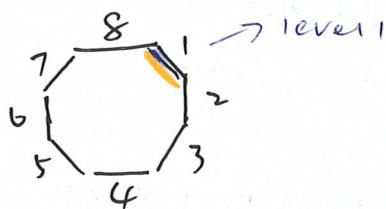
$\geq \geq$ for level 2, needs \geq level 1, 4 nodes.

(Q: Is it possible for MST have level 1?

First step is to merge.

link weight as a clock)

Rest all absorb (increase / decrease)



17/12/11

Data Structure

edge-state {
 ? in-MST
 in-MST
 not-in-MST}

node-state {
 sleeping initial state
 find
 found}

Message type {
 initiate from core find MoE
 report report candidate MoE to core
 test test an edge for being candidate
 MoE
 accept
 reject
 change-root change direction for core
 connect}

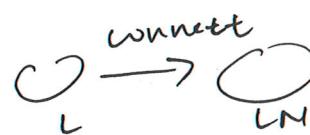
Algorithm

wakeup ()

check MoE.
MoE in MST
connect (MoE)

receive connect

if sleeping
wakeup ()



if $L < LN$.

IN-MST.

Send initiate msg.

find-count.

else

if - ? in-MST
 append queue

else

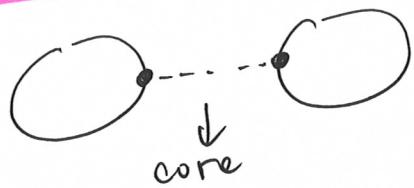
send initiate message

Complexity

Reason absorb: merge takes a lot effort from fragment
to find MST
absorb only takes small fragment searching, doesn't
take large one.

Message complexity: M nodes and E edges
 $2E + 5N \log N$

Election in MST



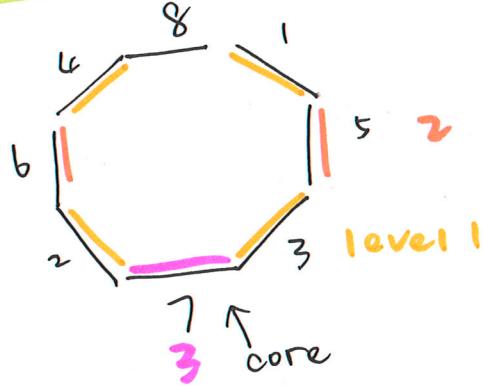
pick node in core
with highest/lowest weight

- Q: Is the core of final MST always same
if build on same network several times?
multiple
/ Same edges for merges & absorbs?

Exercise MST.

- a. final result of MST has level 1?
start with merge, and then all absorbs
- b. Maximum level of MST with $n=2^k$ nodes.
 K .

an example when $n=8$?



Chap 5. Fault Tolerance and Consensus

Example in real world.

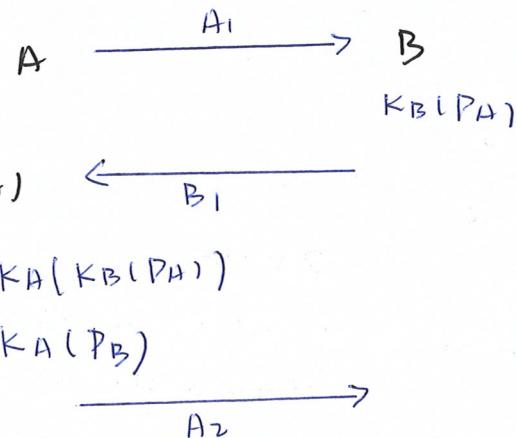
Alice and Bob want to make an appointment.

Alice sends msg A_1 .

Bob receives A_1 . Knows PA: $K_B(P_A)$

Bob sends back B_1 .

Alice receives B_1 .



Problem: messages have arbitrary delay
and may get lost.

Assume the problem has a solution of minimum number of exchanges
last msg may be lost

Common knowledge:

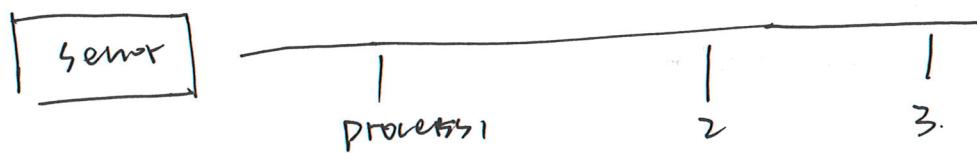
(everyone knows) (everyone knows)

Fault Tolerance & Consensus
processors decide on the value of single bit (Yes or No).
without errors, exchange value \rightarrow majority \rightarrow
failure in system e.g. x tell proposal
tell different value to diff process

Application

- 1) transaction in distributed DB
all sites agree on completion of transaction.
- 2) replicated DB
if want to modify data
in favor of changing?

Original Motivation 1980s
compute direction based on temperature
use redundancy to deal with faults



Sensor give all same value?
→ to agree on input.

Fault classification

fail-stop = just stop

omission = skip a msg

performance: slow

byzantine: execute random behavior.

Model properties

synchronous v.s. asynchronous

easy.

→ differentiate delay & failure

authentication = processor sign msgs.

Network connectivity: complete network

Byzantine



Army want to attack the city simultaneously.
Need to reach an agreement attack / not attack?

17/12/13

Byzantine agreement problem

① agreement.

② validity: correct processes agree on correct initial value with binary value

③ termination: within finite time

source/commander.

$\rightarrow 0 \text{ or } 1$.

Variations of problem

All processes starts with value v_i

1) agree on a vector with a value for all
 (v_1, v_2, \dots, v_n)

2) single value

majority (v_1, v_2, \dots, v_n)

Simple solution for stopping failure

Every process starts with value v .

collects all decision values received in set W .
 $W = \{v\}$ initially

In each round, broadcast W to all other P.
received W_j .
 $W = \text{Union}(W_j)$.

if P fail., its own value circulate in system.
(achieve validity)

Finally < if W contain single v . $\rightarrow v$.
default / majority?

Q: why need round? f+1 round. f processors.

At least 1 round without failure.

Optimization.

1) $|W|=1$ or $|W|>1$?

only broadcast { first initial
first different value \rightarrow decide default

Byzantine Solution.

number of processes: n .

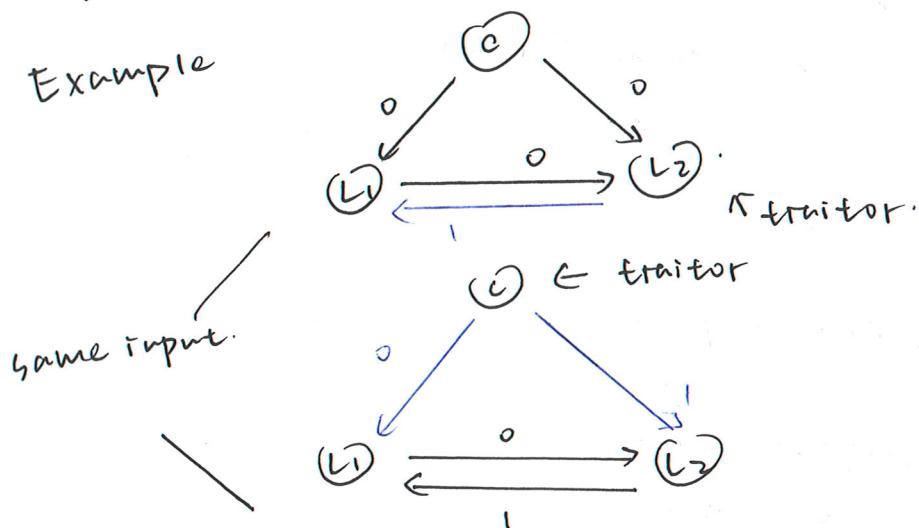
Possible failing processes: f .

$$f < \frac{n}{3}$$

Minimal number of rounds: $f+1$
in deterministic solution.

Some randomized solution have lower expected number of rounds.

Example



Synchronous Byzantine Solution
Recursive Oral Message: No signature

Bottom case = OM (0)
no failure.
commander broadcast initial value
every other P accepts

OM (f) . $f > 0$
commander broadcast initial
number process. commander = 0
lieutenant $1, \dots, n-1$

v_i = received from commander

Recursive.
every P execute $\text{for } o(f-1) \text{ - act as commander}$
decides on majority $(v_1, \dots, v_i, \dots, v_{n-1})$
from frame \uparrow when v_n is commander
decision take majority (\dots)

Number of executions:

exponential

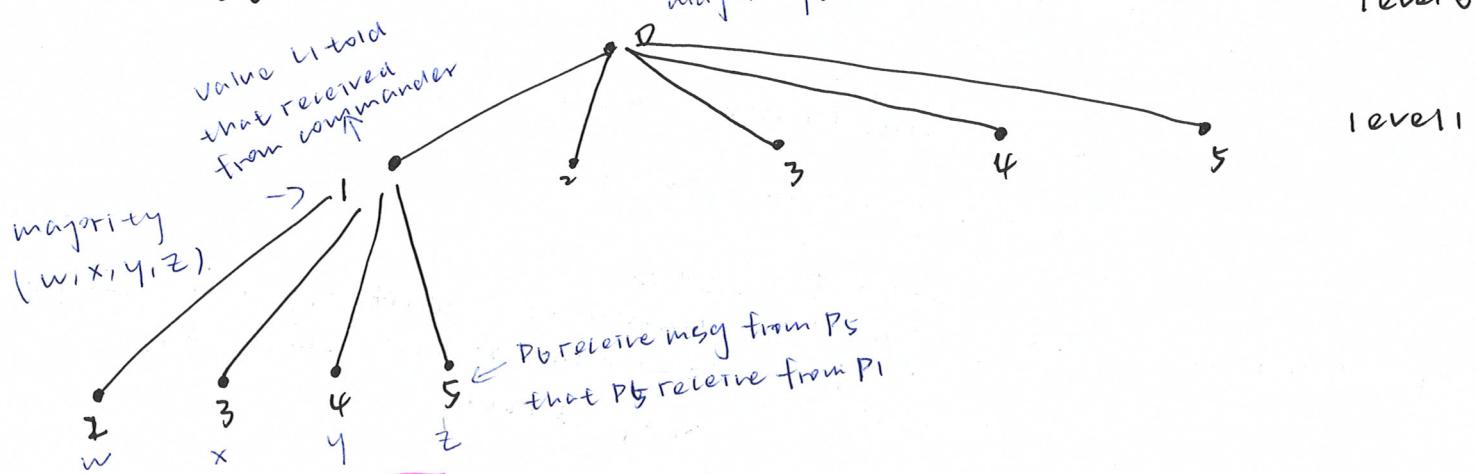
$$OM(f) = 1$$

$$OM(f-1) = (n-1)$$

$$OM(K) : (n-1)(n-2) \dots (n-f+K) \text{ times}$$

Example. Bb local tree

$$\text{bb. } n=7, f=2, i=b. \text{ majority}(V_b, \dots)$$



Decision Taken:

Leaf level: decide on value received $OM(1)$

higher level: take the majority of decisions made in child node

level 0: final decision include value received from commander.

Byzantine agreements with authentication.

Each msg carry a signature, which cannot be forged.

Do by commander.

Allow any number of traitor, f.

Structure of msg: $(V: s_0: s_1: \dots: s_k)$

Every process maintain a set of order V . for deciding (majority...) \uparrow
all potential value in msg received

Commander: send (V, s_0) to all other.

Lientenant: Append signature. Union V with v .

Send msg to all processes do not received.

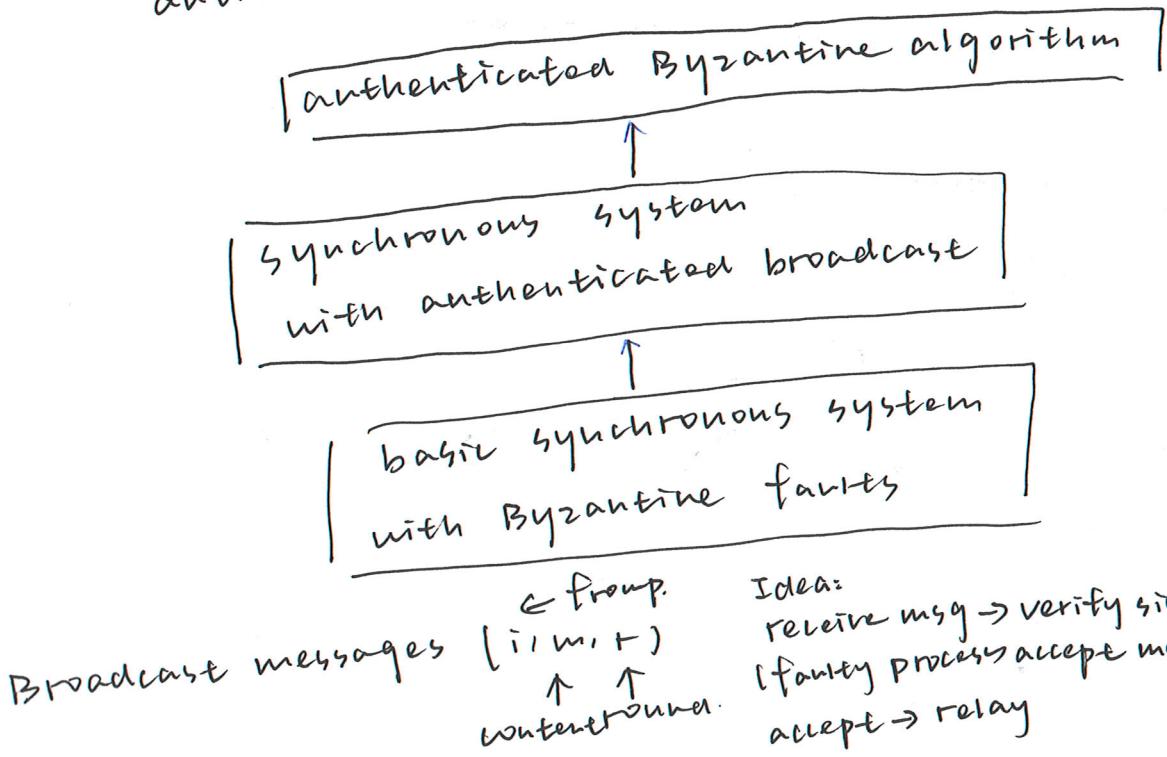
Example:



Efficient solution to Byzantine Agreement

Polynomial

Make communication look authenticated by
authenticated broadcast / consistent broadcast



Broadcast messages ($i \in \mathbb{N}$)

↑ ↑ ← from group
content round

Idea:
receive msg \rightarrow verify signature \rightarrow accept
(faulty process accept much later)
accept \rightarrow relay

Properties of authenticated broadcast:

1) correctness: correct process accept msg from
... ... in same round.

2) unforgeability: if a process is correct and doesn't
broadcast ($i \in \mathbb{N}$), no correct process will accept it.

3) relay: if a correct P receives msg in round $r' \geq r$,
every other accepts it in later round.
at latest the next

Algorithm

commander $\stackrel{\text{Po.}}{\text{start with}} v = 1$

All other P start with $v = 0$

Only $v=1$ are sent

if $v=1$ & never do broadcast before
broadcast (r, i, r)

if in previous round received 1 from Po
and received 1 from other processes $r-1$ times

set $v=1$
divide(v)

Randomized Byzantine Agreement (exponential)

Solution for asynchronous & synchronous system

$n > 5f$

Every process starts with v .

Alg proceeds in rounds consisting 3 phases.

notification phase = propagate msg

proposal phase

decision phase = check enough proposal for final decision

→ if no. flip coin to decide v start in next round

If wait for msgs from everybody else, stop after receiving $n-f$ msgs.

At least $n-f$ correct.

Implement:

broadcast NOTIFICATION to anyone else.

notification
phase

WAIT for $(n-f)$ Notify msg.

→ If gather enough notifications for one value.

broadcast that value $\rightarrow \frac{n-f}{2}$

$\frac{n+f}{2} - f = \frac{n-f}{2}$

proposal

→ If not.

broadcast ? (don't have proposal).

half
reliable
process

→ If decided in previous round.

STOP

→ wait $(n-f)$ proposal msgs

If receive support $> f$ (at least 1 correct supports)

$v \leftarrow w$. (adopt proposal)

decision
phase

if $> 3f$ msgs,

take decision on w

$v \leftarrow$ random.

Lemma 1: No simultaneous contradicting proposals by correct processes in same round.

Proof.: receive $\geq \frac{n-f}{2}$ proposal msgs.

more than $\frac{n-f}{2} - f = \frac{n-f}{2}$ from correct processes
which is majority

Lemma 2: When all correct processes have same value, immediate decision.

If all correct start with v in round 1, all decide v in round 1.

Proof.: receive $n-f$ notification. At least $n-2f$ from ~~correct~~ correct.

$$\begin{aligned} \because n > 5f \quad \therefore n-2f &= \frac{n}{2} + \frac{n}{2} - 2f \\ &> \frac{n}{2} + \frac{5f}{2} - 2f = \frac{n-f}{2} \end{aligned}$$

↑
correct processes propose v

$$\therefore n > 5f \quad \therefore n-2f > 3f.$$

Lemma 3: If a correct process decides v in round r ,
all correct processes propose v in round $r+1$.

Proof.: From Lemma 2, all correct process proposal v in round $r+1$.
If decides v in r , must receive $> 3f$ proposals on v .
of which is correct. $m > 2f$.

Correct processes receive $m-f > f$ proposals on v ,
which is the condition for adopting value.

2017/7/18 Stabilization

A stabilizing algorithm:

deal with transient faults in an implicit way

(operate as usual, if meet some errors, resolve it)

can be started in any configuration

Real world problem:

orchestra play without conductor

player listen what neighbor playing

restart if find play at wrong page, restart..

Types of transient fault:

possible causes: 1) no connection

2) memory crashes

3) transmission errors

Definition

stable predicate P .

legal states is defined by predicate.

First stabilizing Algorithm: Dijkstra.

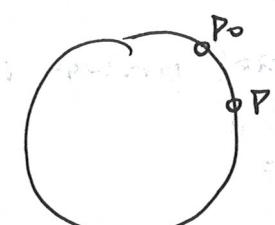
Initial problem: mutual exclusion in unidirectional ring

shared memory, P_i can inspect value in P_0 & its own

from predecessor

composite atomicity: reading and writing in

one atomic/single step



After, P_i has an integer x_i modulo k .

Idea:

Ordinary process P_i : compare value with predecessor

P_{i-1} P_i if different → copy it.

0 → ?

8 → ?

Special process P_0 :

compare value with predecessor
if equal → increment its own value

$$x_0 = (x_0 + 1) \bmod k$$

P_{i-1}	P_i
0	0
8	8
8	9

Q: What is the maximum number of steps before only one process is enabled?

$$(N-1) + (N-2) + \dots + 2+1 = \frac{(N-1)N}{2}$$

Q: What is legal states of ring?

↑
only one process can enter critical section

K, K, K, \dots, K, K (P_0 can enter)

$K, K, K, \dots, K, K-1, K-1, K-1, \dots$ (P_0 can't enter)

Correctness:

Idea: missing label concept:

All processes except P_0 copy value.

P_0 introduces new value

if $K > N$, P_0 go through $0, 1, 2, \dots, K-1$

∴ sometime introduce new value in system

if $K = N$

if $K = N-1$

if K is smaller, don't allow a lot different values in system
don't work. (ex. 3)

Remark:

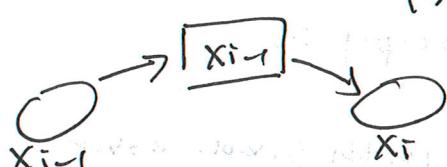
If number of processes is not a prime, a special process is needed.

If shared memory with read/write atomicity:

in one step ^d can only read or write

local integer + intermediate process

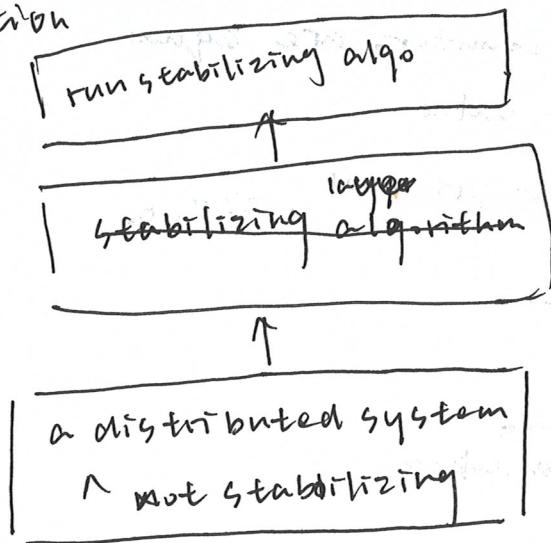
(send by one & receive by other)



$K > 2N - 2$

requires stabilizing communication

Stabilization



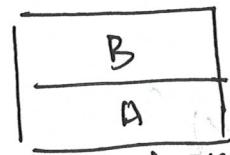
make it stabilizing

Fair protocol composition

Building a stabilizing algo on two different layers.

If have a combined protocol executing both,

do one protocol alternatively.



State transition functions:

for P_1 : $f_1: A \rightarrow A$

for P_2 : $f_2: A \times B \rightarrow B$

↑ only modify B
can access lower level

A is stabilizing layer.

B .

P is also stabilizing layer

by combining 2.

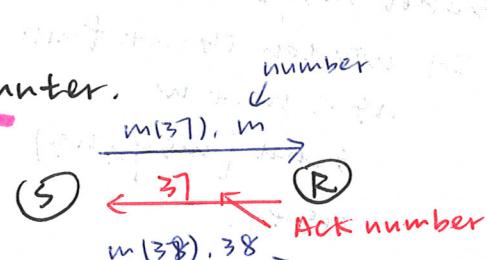
A protocol P is a fair composition/combination if executing states transitions in alternate fashion.

A stabilizing datalink algo: stop-and-wait

Idea: send a msg, and wait for acknowledgement to send next msg.

Sender and receiver maintain counter.

Send msg with msg number.



Whenever a receiver receive a msg,

if it is a new msg, adopt value in msg,
then send Ack back.

TimeOut = If sender doesn't receive ACK.

resend msg

Legal state: all counters are identical. (sender, receiver, msg)

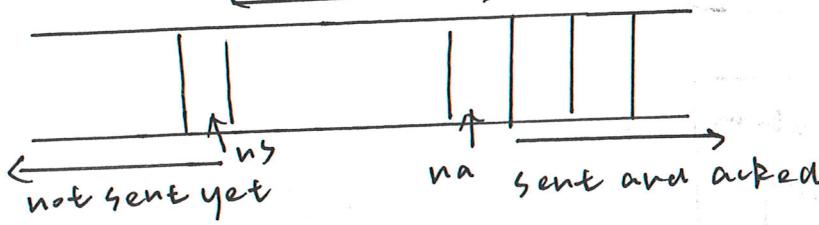
Remark:

System in legal state if all counters are equal.

Sender introduce new higher value

A stabilizing algo: sliding - window protocol TCP like

sent but not acked msg.



Sender { w: window size

 ns: the number of next msg to send

 first non-acked msg

 na: next msg to be received

Receiver = NR:

Stable predicate P

(na = nr) and (nr ≤ ns) and (ns ≤ na + w)

can't have ACK
if don't received

can't receive
if don't send

window has size

for each (message, i) in channel SR, i < ns

and

for each (ack, i) in channel RS, i < nr
can't have ACK if no received

Implement (Not stabilizing version)

Send msg:

If window is not full:

$ns < na + w$

 send (msg, ns)

$ns++$.

Receiving a msg:

check if it's the next msg expected
adapt nr.

send ACK. (whether expected or not)
(reserve all previous ones)

Receive ACK:

If is new ACK possibly
adapt na, not FIFO

Timeout in sender

Resend all msgs in window
unacked