

Paradigms for Process Interaction in Distributed pgms.

Distributed computation — Concurrent pgms in which processes communicate by msg passing.

Process Interactions

1. networks of filters ↵ (merge sort)
2. clients & servers
3. heartbeat algns
4. probe/echo algns
5. broadcast algns
6. token-passing algns
7. decentralized servers
8. bags of tasks.

Message Passing

Decisions

- What kind of msgs
- When ?
- How

Different types of process

- | | input | Output |
|------------|--|---------|
| 1. filters | → "data transformer" (channel) | channel |
| 2. clients | → triggering process | |
| 3. Servers | → reactive process (non-terminating process) | |
| 4. peers | → two peers | |
- For Eg:-

Interprocess Communication

Programming Notation

channel - physical commⁿ

↳ Two primitives

(1) send

(2) receive

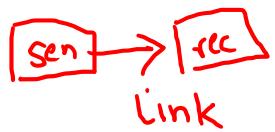
} global

oneway / two-way

Design choices

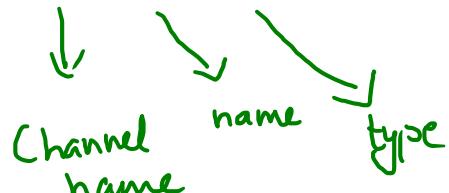
- ① Asynchronous message passing → non-blocking ↳ different chanls used for diff kind of mugs
- ② Synchronous message passing → blocking
- ③ Generative communication → processes share a single commⁿ channel ↳ tuple space
 - ↳ naming association
- ④ RPC (Remote Procedure Call)
- ⑤ Rendezvous

OS
1. pipes
2. SHM
3. message passing
4. Semaphore
5. Monitors.

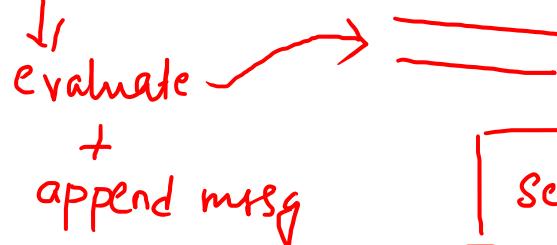


Chan input (char)

Chan ch ($f_1:t_1, \dots, f_n:t_n$)



Send 'ch (expr₁, ..., expr_n)'



channels are declared globally

```

chan input(char) ← stream of characters
chan output([1:MAXLINE] char) ← Line
Char_to_Line:::
var line[1:MAXLINE] : char
var i : int := 1
do true →
    receive input(line[i])
    do line[i] ≠ CR and i < MAXLINE →
        {line[1:i] contains last i input characters}
        i := i + 1; receive input(line[i])
    od
    send output(line); i := 1
od

```

busy-wait polling

receive → blocking primitive

unbounded

Send is nonblocking



Filters: A sorting network

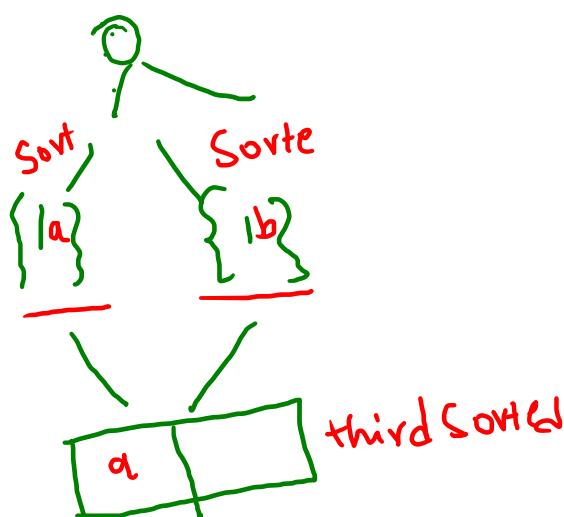
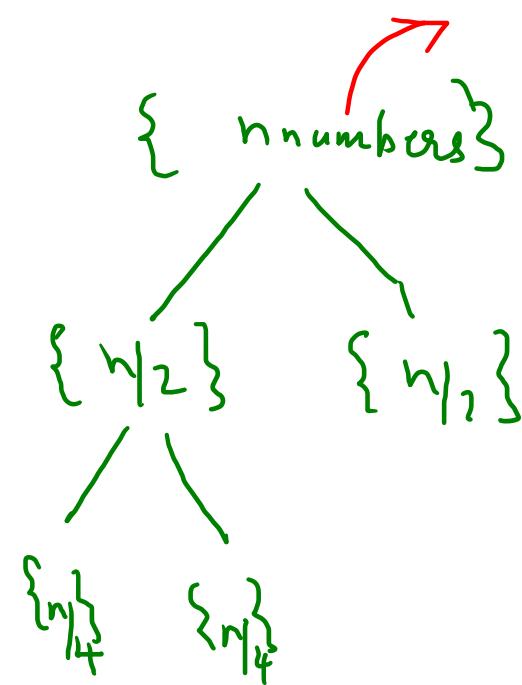
- data transformer
- $\text{fun}(\text{input}) \rightarrow \text{output}$
- " o/p of filter is a fn of its i/p"

SORT: $(\forall i: 1 \leq i < n: \text{sent}[i] \leq \text{sent}[i + 1]) \wedge$
values sent to *output* are a permutation of values received from *input*

An outline of *Sort* is as follows:

chan *input(int)*, *output(int)*

Sort:: declarations of local variables
receive all numbers from input
sort the numbers
send the sorted numbers to output



end sentinel

MERGE:

$in1$ and $in2$ are empty \wedge $sent[n + 1] = EOS$ \wedge
 $(\forall i: 1 \leq i < n: sent[i] \leq sent[i + 1]) \wedge$
values sent to out are a permutation of values received from $in1$ and $in2$

filter

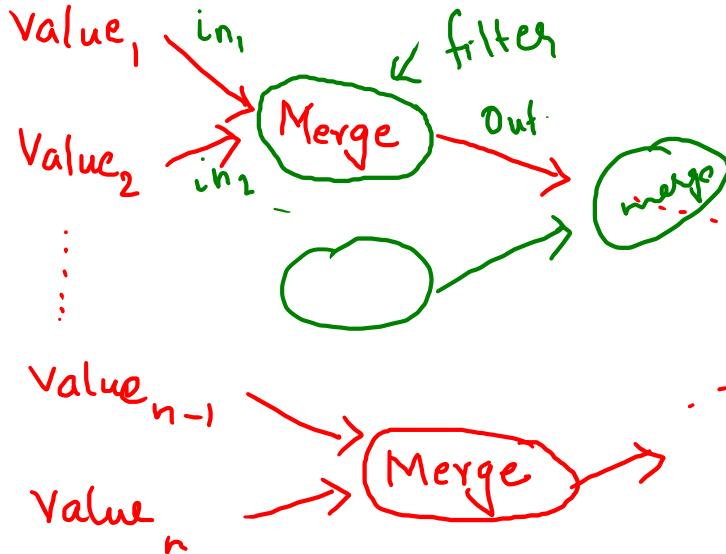
```

chan in1(int), in2(int), out(int)
Merge:: var v1, v2 : int
receive in1(v1); receive in2(v2)
do more input to process →
  send smaller of v1 and v2 to out
  receive another input value from
    in1 or in2
od
send out(EOS) { MERGE }

```

network of filters

Sorting network



chan $in1(int)$, $in2(int)$, $out(int)$

Merge:: var $v1, v2 : int$

receive $in1(v1)$; receive $in2(v2)$

do $v1 \neq \text{EOS}$ and $v2 \neq \text{EOS} \rightarrow$

if $v1 \leq v2 \rightarrow$

send $out(v1)$; receive $in1(v1)$

$\| v2 \leq v1 \rightarrow$

send $out(v2)$; receive $in2(v2)$

fi

$\| v1 \neq \text{EOS}$ and $v2 = \text{EOS} \rightarrow$

send $out(v1)$; receive $in1(v1)$

$\| v1 = \text{EOS}$ and $v2 \neq \text{EOS} \rightarrow$

send $out(v2)$; receive $in2(v2)$

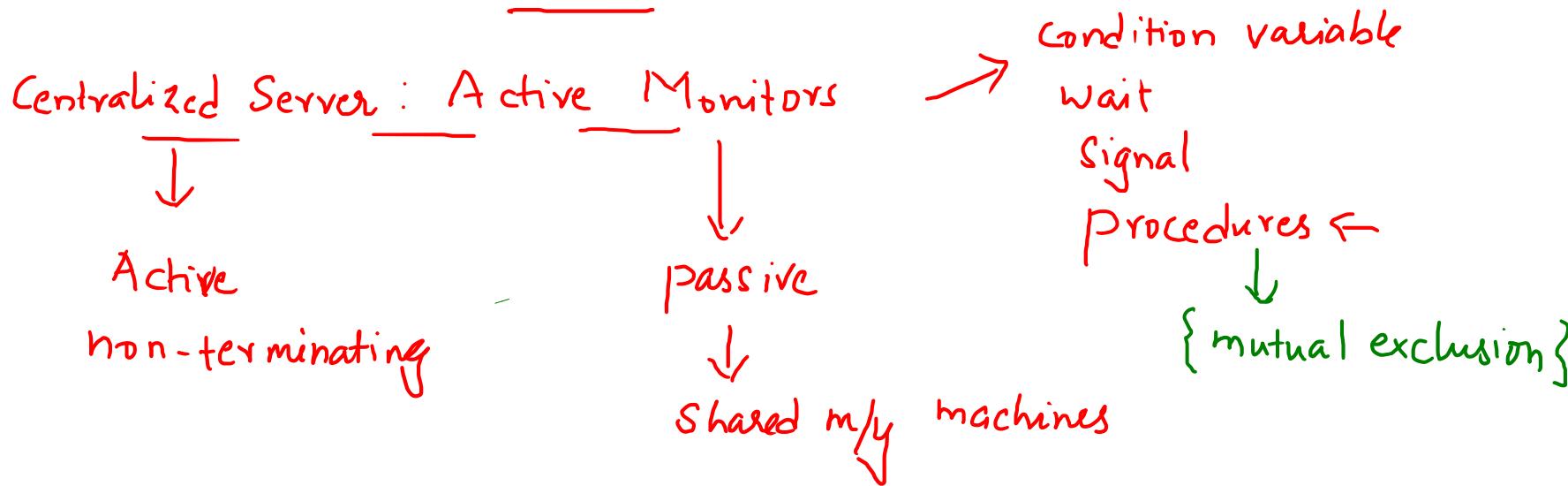
od

send $out(\text{EOS}) \{ MERGE \}$

exhaustive

Communication pattern \rightarrow tree

Clients & Servers.



```
monitor Mname # Invariant MI
  var permanent variables
    initialization code
  proc op1(formals1) body1 end
  .
  .
  proc opn(formalsn) bodyn end
end
```

```

monitor Resource_Allocator
  {ALLOC: avail ≥ 0 ∧ (avail > 0) ⇒ empty(free)}
  var avail : int := MAXUNITS, units : set of int,
    free : condition
  code to initialize units to appropriate values
  ① proc acquire( res id : int )
    if avail = 0 → wait(free)
    [] avail > 0 → avail := avail-1
    fi
    id := remove(units)
    end
  ② proc release( id : int )
    insert(id, units) ←
    if empty(free) → avail := avail+1
    [] not empty(free) → signal(free) ← awakens any delayed
    fi
    end
  end

```

Figure 2. Resource allocation monitor.



Static naming
Dynamic naming

```

type op_kind = enum(ACQUIRE, RELEASE)
chan request(index : int, op_kind, unitid : int)
chan reply[1:n](int)

Allocator:: var avail : int := MAXUNITS, units : set of int, pending : queue of int
var index : int, kind : op_kind, unitid : int
code to initialize units to appropriate values
{ALLOC: avail ≥ 0 ∧ (avail > 0) ⇒ empty(pending)}
do true → receive request(index, kind, unitid)
  if kind = ACQUIRE →
    if avail > 0 → # honor request now
      avail := avail-1; unitid := remove(units)
      send reply[index](unitid)
      [] avail = 0 → # save index of requesting process
        insert(pending, index)
    fi
    [] kind = RELEASE →
      if empty(pending) → # release unitid
        avail := avail+1; insert(units, unitid)
      [] not empty(pending) → # give unitid to first requester
        index := remove(pending); send reply[index](unitid)
    fi
  od
fi

Client[i: 1..n]:: var unitid : int
...
# acquire a unit of the resource
send request(i, ACQUIRE, 0) # unitid not needed
receive reply[i](unitid)
# use resource unitid and then later release it
send request(i, RELEASE, unitid)

```

Figure 3. Resource allocator and clients.

Monitor-Based Programs

permanent variables
procedure identifiers
procedure call
monitor entry
procedure return
wait statement
signal statement
procedure bodies

Message-Based Programs

local server variables
request channel and operation kinds
send request; receive reply
receive request
send reply
save pending request
retrieve and process pending request
arms of case statement on operation kind

Figure 4. Duality between monitors and centralized servers.

Disk scheduling

Shortest Seek time (SST)

moving-head disk

PF → cylinder no.:, track number
an offset

whose cylinder is closest to
the current position

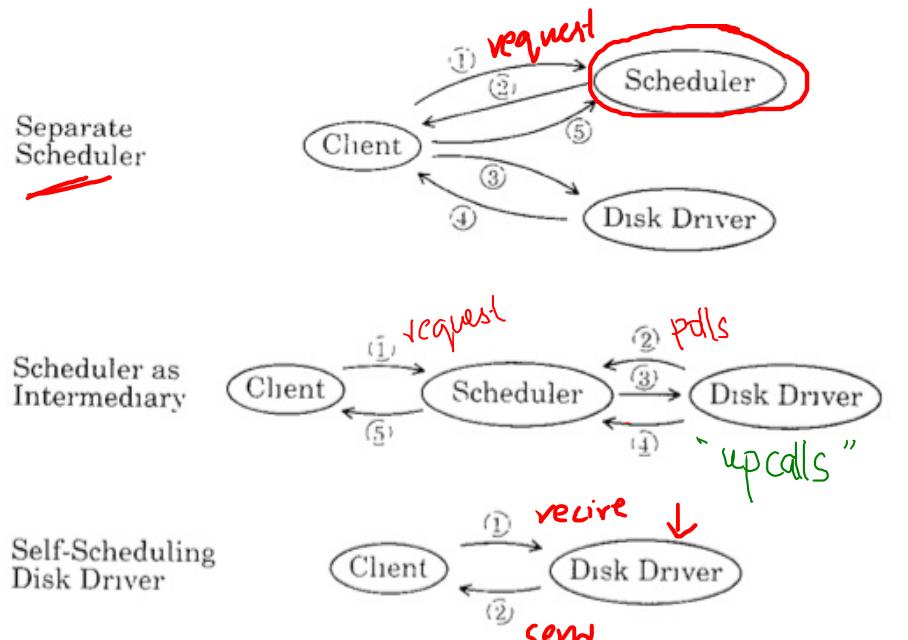


Figure 5. Disk scheduling structures with message passing.

lower []
 ↙
 Current Head position
 ↙ next request

higher []

```

→ chan request(index : int, cylinder : int, other argument types)
  # other arguments indicate read or write, disk block, memory buffer, etc.
→ chan reply[1:n](result types)
Disk_Driver:: var lower, higher : ordered queue of (int, int, other argument types)
  # contain index, cylinder, and other arguments of pending request
var headpos : int := 1, nsaved := 0
var index : int, cyl : int, args : other argument types
{ SST: (lower is an ordered queue from largest to smallest cyl ∧
  all values of cyl in lower are ≤ headpos) ∧
  higher is an ordered queue from smallest to largest cyl ∧
  all values of cyl in higher are ≥ headpos ∧
  (nsaved = 0) ⇒ (both lower and higher are empty) }

do true →
  do not empty(request) or nsaved = 0 →
    # wait for first request or receive another one, then save it
    // receive request(index, cyl, args)
    if cyl ≤ headpos := insert(lower, (index, cyl, args))
    [] cyl ≥ headpos := insert(higher, (index, cyl, args))
    fi
    nsaved := nsaved+1
  od
  # select best saved request; there is at least one in lower or higher
  if size(lower) = 0 → (index, cyl, args) := remove(higher)
  [] size(higher) = 0 → (index, cyl, args) := remove(lower)
  [] size(higher) > 0 and size(lower) > 0 →
    { remove element from front of lower or higher depending
      on which saved value of cyl is closer to headpos
    fi
    headpos := cyl; nsaved := nsaved-1
    access the disk
    2= send reply[index](results)
  od

```

Figure 6. Self-scheduling disk driver.

Heart beat Algm

Model
 Processor → process parallel

Communication
links

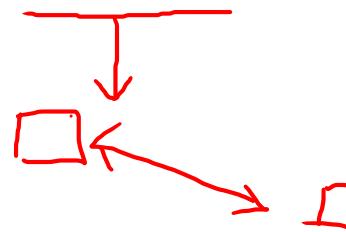
→ Shared
Channels

Pblm:- Each processor to determine the topology of entire n/w

```
var top[1:n, 1:n] : bool := ([n*n] false) # topology to be stored as set of all links
Node[p: 1..n]:= var links[1:n] : bool
# initialized so links[q] is true if q is a neighbor of Node[p]
top[p, 1:n] := links # fill in p'th row of top
{ top[p, 1:n] = links[1:n] }
```

Figure 8. Network topology using shared variables.

{graph
DS}



Solution

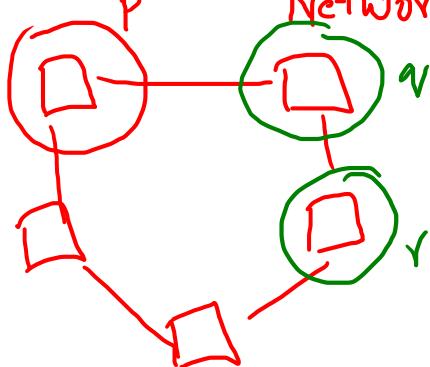
(1) Shared memory access

(2) distributed computation by replicate

heart beating : first expand, send infoⁿ out

then contract, gather new infoⁿ in

Network Topology :- Distributed Solution



Single process T
 ↑
 Other nodes send
 its value of links

Drawback

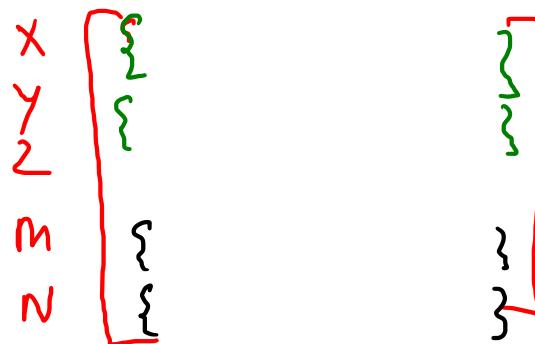
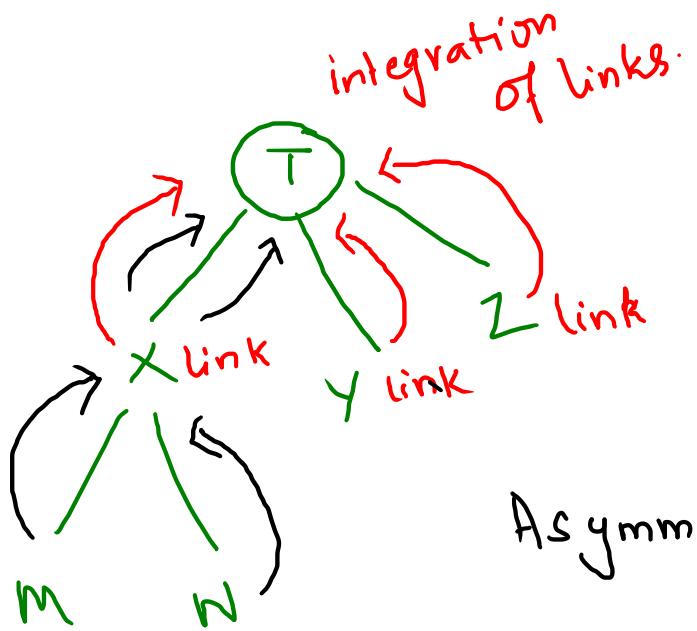
T is not one of the
 neighbor. {Symmetric algm}

```

chan topology[1:n](top[1:n, 1:n] : bool) # one channel per node
Node[p: 1..n]:: var links[1:n] : bool [0, 0, 1, 1, 0]
                                              # initialized so  $links[q]$  is true if  $q$  is a neighbor of  $Node[p]$ 
var top[1:n, 1:n] : bool := ([n*n] false) # local view of topology
var r : int := 0
var newtop[1:n, 1:n] : bool
top[p, 1..n] := links # fill in row for my neighbors
{top[p, 1:n] = links[1:n]  $\wedge$  r = 0}
{ROUND: (  $\forall q: 1 \leq q \leq n: (dist(p, q) \leq r) \Rightarrow top[q, *]$  filled in )}
do r < D →
  # send local knowledge of topology to all neighbors
  fa q := 1 to n st links[q] → send topology[q](top) af
  # receive their local topologies and or it with top
  fa q := 1 to n st links[q] →
    receive topology[p](newtop)
    {top := top or newtop}
  af
  r := r+1
od
{ROUND  $\wedge$  r = D} {TOPOLOGY}
  
```

Figure 9. Heartbeat algorithm for network topology; first refinement

"Diameter of the
 network" $\rightarrow D$
 "distance b/w the farthest
 pair of nodes"



Symmetric :- Each node needs to be able to compute for itself the entire topology



Node p \rightarrow knows about the links to its neighbors

After one round of msgs

\hookrightarrow P will have learned about topology within two links of it

after two rounds of msg

\hookrightarrow each node will know about the topology within 3 links of it.

\vdots \vdots
r rounds

$$\left\{ \forall q_r : 1 \leq q_r \leq n : \text{dist}(p, q_r) \leq r \right. \\ \left. \text{top} [q_r, *] \text{ filled in} \right\}$$

```

chan topology[1:n](sender : int, done : bool, top[1:n, 1:n] : bool)
Node[p: 1..n]:: var links[1:n] : bool
    # initialized so that links[q] true if q is a neighbor of Node[p]
    var active[1:n] : bool := links      # neighbors who are still active
    var top[1:n, 1:n] : bool := ([n*n] false)      # local view of topology
    var r : int := 0, done : bool := false
    var sender : int, qdone : bool, newtop[1:n, 1:n] : bool
    top[p, 1..n] := links      # fill in row for my neighbors
{ top[p, 1:n] = links[1:n]  $\wedge$  r = 0  $\wedge$   $\neg$ done }
{ ROUND  $\wedge$  (done  $\Rightarrow$  all rows in top are filled in) }
do not done  $\rightarrow$ 
    # send local knowledge of topology to all neighbors
    fa q := 1 to n st links[q]  $\rightarrow$  send topology[q](p, false, top) af
    # receive their local topologies and or it with top
    fa q := 1 to n st links[q]  $\rightarrow$ 
        receive topology[p](sender, qdone, newtop)
        top := top or newtop
        if qdone  $\rightarrow$  active[sender] := false fi
    af
    if all rows of top have some true entry  $\rightarrow$  done := true fi
    r := r+1
od
{ ROUND  $\wedge$  all rows in top are filled in } { TOPOLOGY }
# send topology to all neighbors who are still active
fa q := 1 to n st active[q]  $\rightarrow$  send topology[q](p, true, top) af
# receive one message from each to clear up message queue
fa q := 1 to n st active[q]  $\rightarrow$  receive topology[p](sender, qdone, newtop) af

```

dynamically //

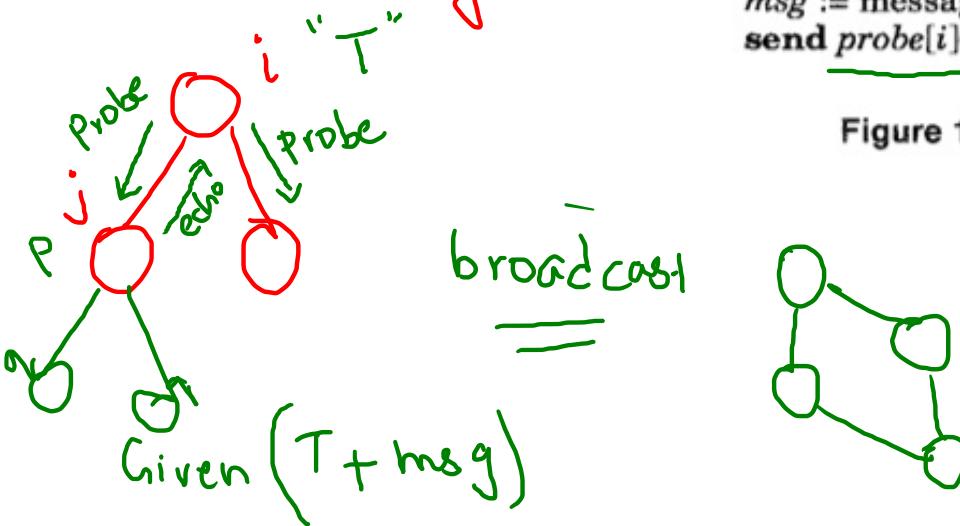
Figure 10. Heartbeat algorithm for network topology; final version.

Probe/Echo Alg.

Tree, Graph

{DFS →}

Probe → a msg sent

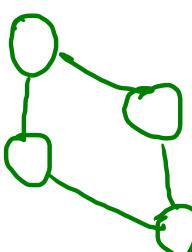


```

chan probe[1:n](span_tree[1:n, 1:n] : bool, message_type)
Node[p: 1..n]:: var span_tree[1:n, 1:n] : bool, msg : message_type
receive probe[p](span_tree, msg)
fa q := 1 to n st q is a child of p in span_tree →
send probe[q](span_tree, msg)
af
Initiator:: var i : int := index of node that is to initiate broadcast
var top[1:n, 1:n] : bool # initialized with topology of the network
var span[1:n, 1:n] : bool, msg : message_type
compute spanning tree of top rooted at i and store it in span
msg := message to be broadcast
send probe[i](span, msg)

```

Figure 11. Broadcast using a spanning tree.



$$G = (V, E)$$

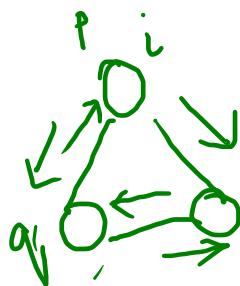
$$G_{\text{MST}} = (V, E' \subset E)$$

First node i

Sends the msg

to all its neighbor

"Form tree"



chan probe[1:n](message_type)

Node[p: 1..n]:: var links[1:n] : bool := neighbors of node p
var num : int := number of neighbors of p, msg : message_type
receive probe[p](msg)
send msg to all neighbors
fa q := 1 to n st links[q] → send probe[q](msg) af
receive num-1 redundant copies of msg
fa q := 1 to num-1 → receive probe[p](msg) af

Initiator:: var i : int := index of node that is to initiate broadcast
var msg : message_type := message to be broadcast
send probe[i](msg)

Figure 12. Broadcast using neighbor sets

```

const source = i    # index of node that initiates the algorithm
chan probe[1:n](sender : int)
chan echo[1:n](links[1:n, 1:n] : bool)    # contents are part of the topology
chan finalecho(links[1:n, 1:n] : bool)    # final echo to Initiator

Node[p: 1..n]:: var links[1:n] : bool := neighbors of node p
                    var localtop[1:n, 1:n] : bool := ([n*n] false)
                    localtop[p, 1:n] := links    # put neighbor set in localtop
                    var newtop[1:n, 1:n] : bool
                    var parent : int    # will be node from whom probe is received
                    receive probe[p](parent)
                    # probe on to all other neighbors, who are p's children
                    fa q := 1 to n st links[q] and q ≠ parent → send probe[q](p) af
                    # receive echoes for all children and union them into localtop
                    fa q := 1 to n st links[q] and q ≠ parent →
                        receive echo[p](newtop); localtop := localtop or newtop
                    af
                    if p = source → send finalecho(localtop)    # this node is the root
                    l p ≠ source → send echo[parent](localtop)
                    fi

Initiator:: var top[1:n, 1:n] : bool    # network topology as set of links
                send probe[source](source)
                receive finalecho(top)

```



Figure 13. Probe/echo algorithm for topology of a tree.

Broadcast Algs

Processors = 'n'

{ send
receive

ch[1:n] = array of channels

ith process to

p[i]

{ broadcast ch(m) } → 'n' send messages

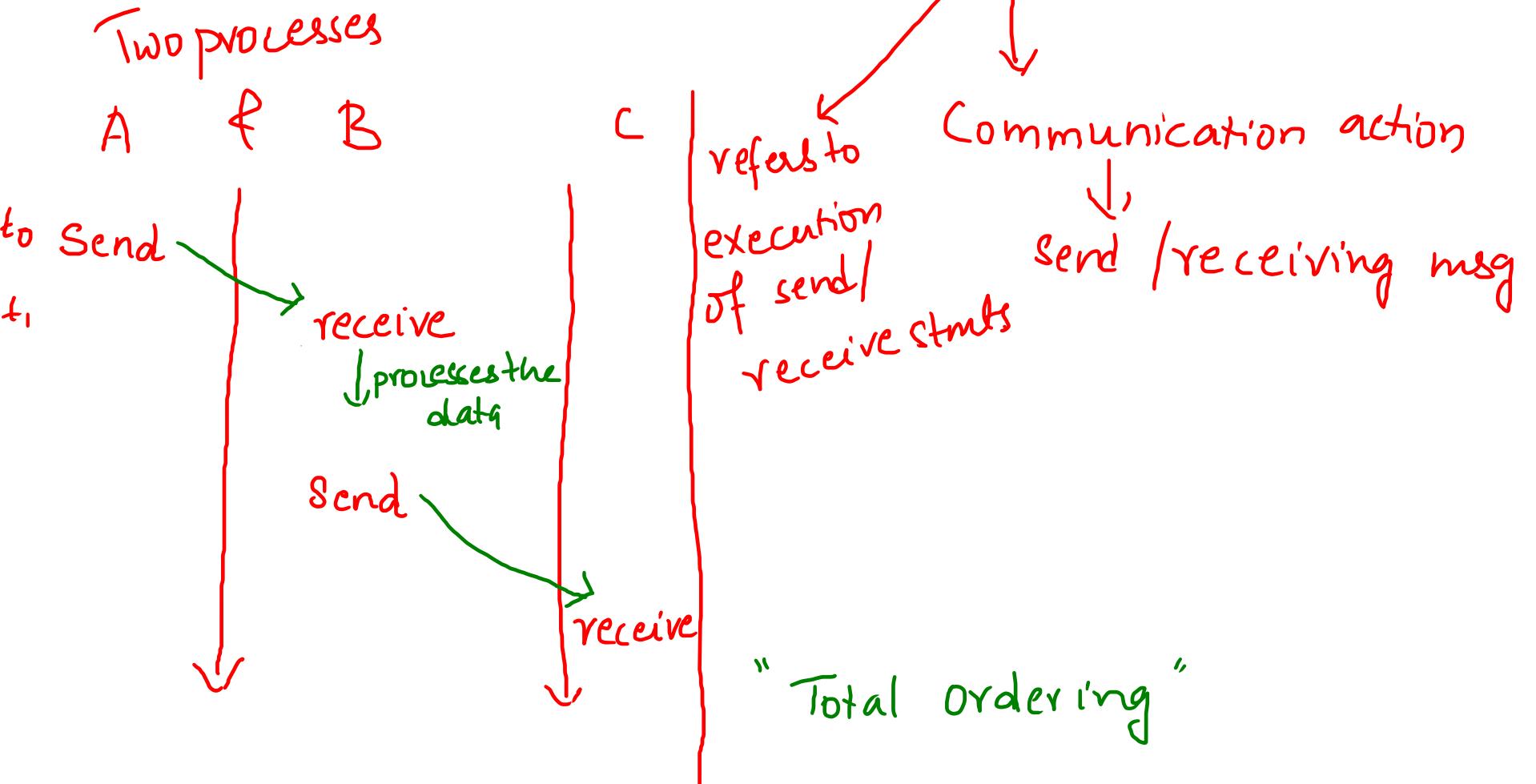
receive

Two processes
A & B

\ /

different orders

Logical clocks & Event Ordering



Logical clock

is an integer counter that is
incremented when events occur.

Process = logical clock + msg contains a timestamp.

lc in process A.

(1) Send — {current value of lc} ^{timestamp} and increment lc by 1

(2) receive — $lc = \max(lc, ts+1)$ then increment lc by 1

Distributed Semaphores

$P(s)$ $V(s)$
↓ ↓
 $\left\{ \begin{array}{l} \text{if } s \text{ is true} \\ \text{decrements}(s) \end{array} \right.$ Increment s

Semaphore invariant

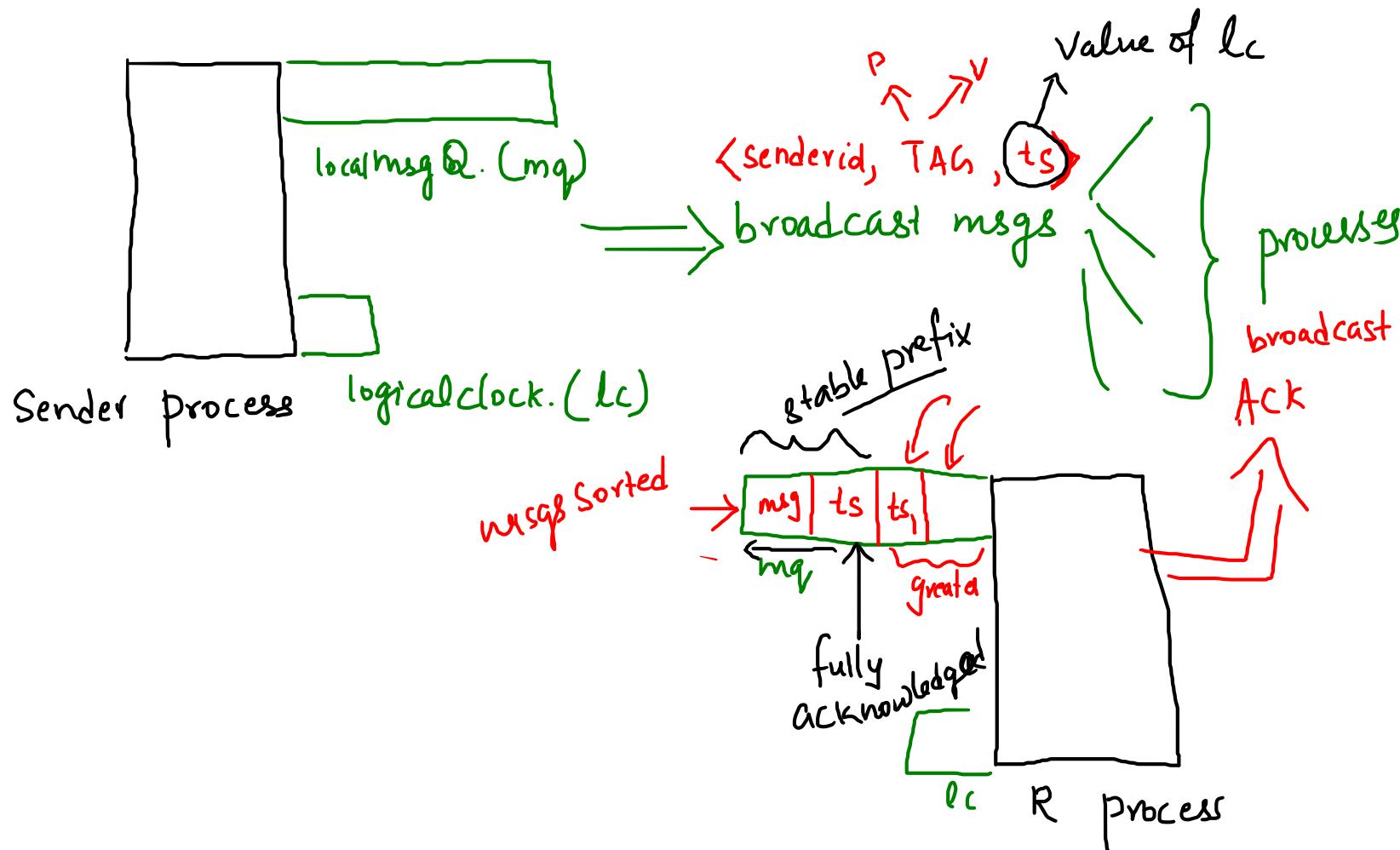
The no: of completed P Op^{ns} is at most the no: of completed V Op^{ns} + the Semaphore's initial value.

(1) how to represent the value of a semaphore

(2) and maintaining the semaphore invariant.

- * Way to count P and V op^{ns}
- * way to delay P op^{ns}

Broadcast msgs



local counters

n_P , n_V

of fully acknowledged
 $P \& V$ ops

execute $P \& V$ op^{ws}
stored in the stable
prefix of mq

$\{$ $n_V > n_P$
increment n_P } Increment n_V

— Delete the V or P msg from the mq

```

type kind = enum(V, P, ACK)
chan sem[1:n](sender: int, kind, timestamp: int)
chan go[1:n](timestamp : int)

User[i: 1..n]:: var lc : int := 0      # logical clock
                    var ts : int        # timestamp in go messages
                    # execute a V operation
                    broadcast sem(i, V, lc); lc := lc+1
                    ...
                    # execute a P operation
                    broadcast sem(i, P, lc); lc := lc+1
                    receive go[i](ts); lc := max(lc, ts+1); lc := lc+1
                    ...

Helper[i: 1..n]:: var mq : queue of (int, kind, int)    # ordered by timestamps
                    var lc : int := 0                      # logical clock
                    var nV : int := 0, nP : int := 0        # semaphore counters
                    var sender : int, k : kind, ts : int   # values in messages
                    do true → {loop invariant DSEM}
                        receive sem[i](sender, k, ts); lc := max(lc, ts+1); lc := lc+1
                        if k = P or k = V →
                            insert (sender, k, ts) at appropriate place in mq
                            broadcast sem(i, ACK, lc); lc := lc+1
                        [] k = ACK →
                            record that another ACK has been seen
                            fa fully acknowledged V messages →
                                remove the message from mq; nV := nV+1
                            af
                            fa fully acknowledged P messages st nV > nP →
                                remove the message from mq; nP := nP+1
                                if sender = i → send go[i](lc); lc := lc+1 fi
                            af
                        fi
                    od

```

Token Passing Algs

Token → Special kind of msg

- Convey permission
- gather info

1. Critical Section Problem

2. Termination Detection in distributed computation

Token Passing Algs

Token → Special kind of msg

- Convey permission
- gather info

1. Critical Section Problem

2. Termination Detection in distributed computation

Distributed Mutual Exclusion

Critical Section problem :- It ensures that at most one process at a time executes code that accesses a shared resource

1. Active monitors → nonreplicated files
2. Distributed Semaphores → decentralized solution
" huge number of msgs to be exchanged "
3. Token ring → decentralized solution.

$P[r : n]$

Critical
Section
Codes

+ noncritical
Section
Codes

access ↓

{ Shared resource }

only one
process →

↑
mutual
exclusion

↓

They can
execute
concurrently.

entry }
exit } protocols.

Tokens → to control the entry to critical section.



Token is circulated
b/w the Helper.

