

# DISTRIBUTED SYSTEMS

## NOTATIONS

- $\langle \text{Send} | q, m \rangle$ : send  $m$  to  $q$
- $\langle \text{Deliver} | p, m \rangle$ : delivers  $m$  from  $p$
- $\text{Crash}(p_j)$ : after this event the process  $p_j$  doesn't execute any local computation step ( $\text{no } \text{Exec}(j)$ ).
- $\text{Byz}(p_j)$ : // // //  $p_j$  behaves in an arbitrary way (when we have  $\text{Exec}(j)$ ) we don't follow anymore the automaton of  $p_j$
- **UNIFORM**: usually means that also the processes that are FAULTY participate in the property.
- **LOG**: APPEND ONLY List of Commands
  - $\text{ENTR}$   $x$   $\leftarrow$  term = period of time with same leader
  - $\text{add}$   $\leftarrow$  command
- **LOG CONSISTENCY**: If log entries on different servers have SAME INDEX and TERM, 1. they store the same command; 2. the logs are identical in all preceding entries.

# THEORY CONCEPTS

## TIME

### - HAPPENED-BEFORE RELATION

Two events  $e$  and  $e'$  are in a happened-before relationship ( $e \rightarrow e'$ ) if:

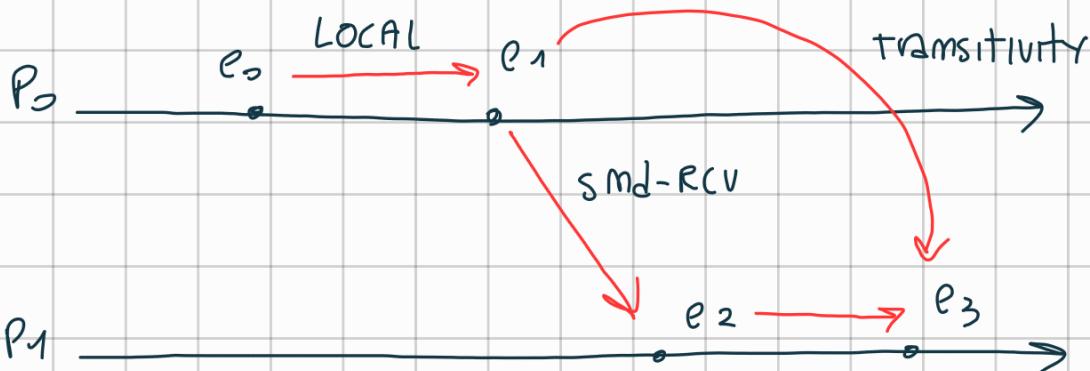
1. LOCAL ORDERING :  $\exists p_i \mid e \rightarrow_i e'$

2. SENDER-RECEIVING ORDERING :  $\forall m, \text{send}(m) \rightarrow \text{receive}(m)$

$e = \text{send}(m)$  is the event of sending  $m$

$e' = \text{receive}(m)$  is the event of receipt the same  $m$

3. Transitivity :  $\exists e'' : (e \rightarrow e'') \wedge (e'' \rightarrow e') \text{ then } e \rightarrow e'$



## FAILURE DETECTOR

- FAIL STOP = You can detect failures



- FAIL NOISY = You can make mistakes



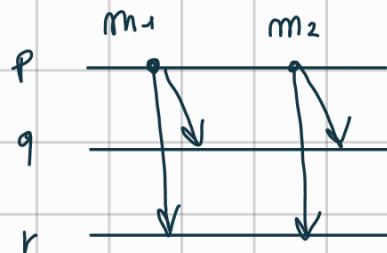
- FAIL SILENT = You can't detect failures

## BROADCAST

- QUORUM =  $m$  processes  $\rightarrow |P| = m$ , a quorum is a subset of  $P$  of size  $\geq m/2 + 1$
- BEST EFFORT BROADCAST : If the SENDER is NOT CORRECT someone can deliver the message, someone can not.
- CAUSAL ORDER BROADCAST : Is an extension of the happened-before relation: a message  $m_1$  have potentially caused another message  $m_2$  ( $m_1 \rightarrow m_2$ ) if :

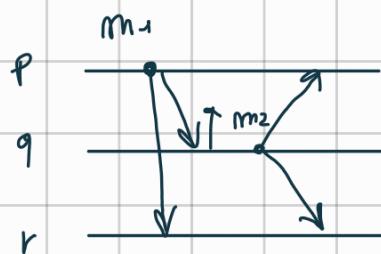
LOCAL

1. Some process  $p$  broadcast  $m_1$  before it broadcast  $m_2$ ;



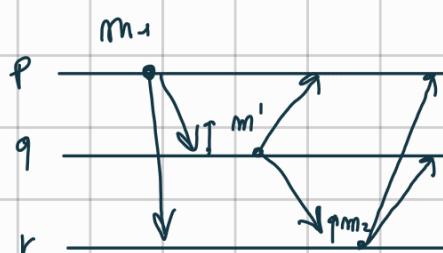
SENDER-RECEIVING

2. Some process  $p$  delivers  $m_1$  and subsequently broadcast  $m_2$ ;



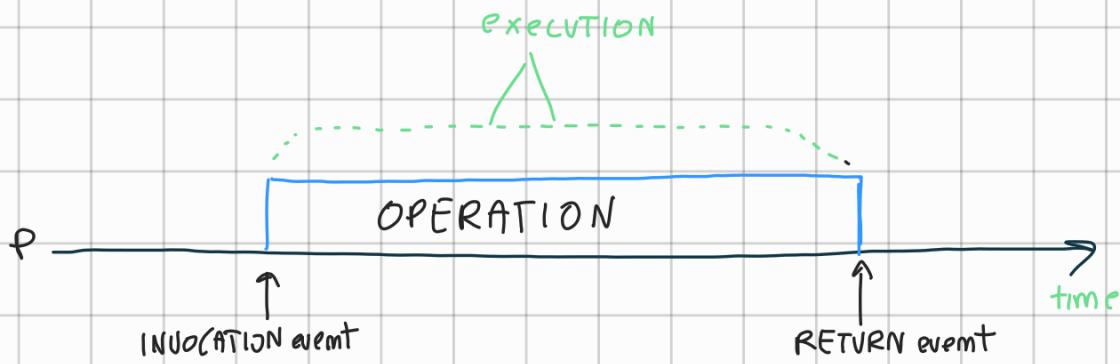
TRANSITIVITY

3. there exists some message  $m'$  s.t.  $m_1 \rightarrow m'$  &  $m' \rightarrow m_2$



## REGISTER

- **read () → v** : returns the "current" value  $v$  of the register
- **write( $v$ )** : Writes the value  $v$  in the register and returns true at the end of the operation.



$O$  IS CONCURRENT WITH  $O'$

- **REGISTERS SEMANTIC**: How the register behaves when multiple processes WRITE and READ.
- **( $X, Y$ )**: a register where  $X$  processes can WRITE and  $Y$  can READ
  - e.g.  $(1, 1) \rightarrow$  only a process can write & read

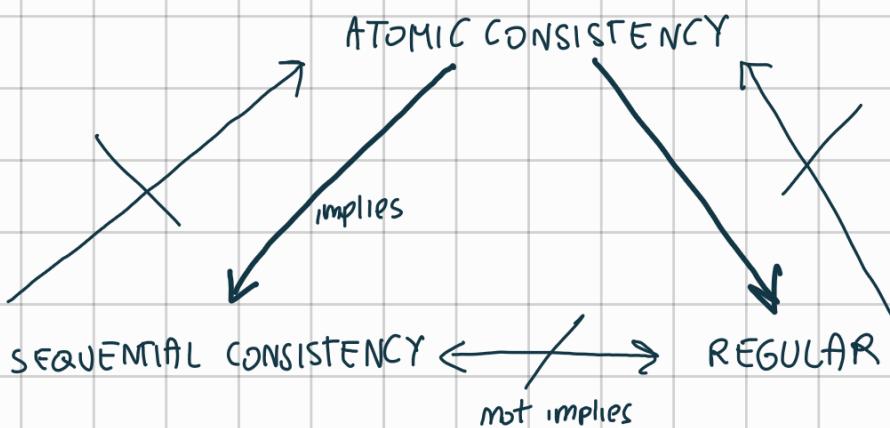
- **REGULAR CONSISTENCY**: A read operation returns the last value written or the value concurrently written.
- **SEQUENTIAL CONSISTENCY**: An execution is sequential consistency if  $\exists G$  GLOBAL ORDER  $G$  (serialized and sequential global order, a fictional one, not the one that happened in real life, we construct it) which justify the values that all the operation write or read

If  $\exists G \Rightarrow$  sequential consistency

Important: In each process you have to respect its LOCAL ORDER.

- **ATOMIC CONSISTENCY / LINEARIZABILITY**: An execution is ATOMIC if I can find a "fictional" point (LINEARIZATION POINT) in each operation in which the operation happen.

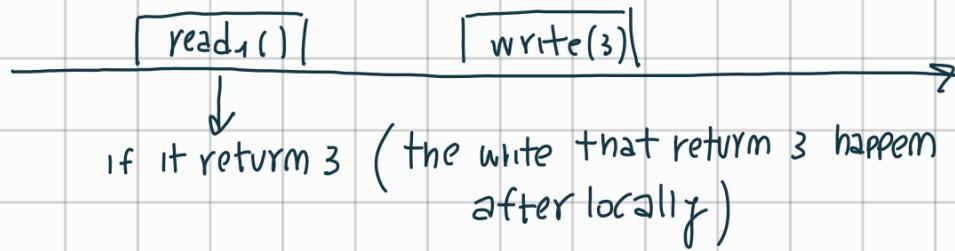
### • REGISTER CONSISTENCY RELATION



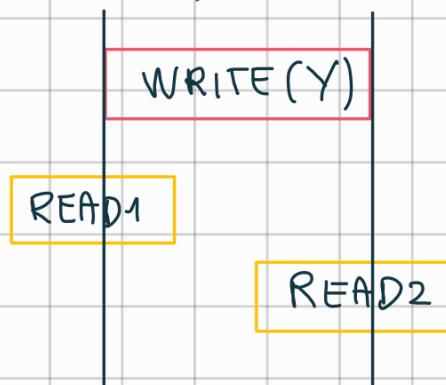
- Return a sequence of value that is **NOT SEQUENTIAL CONSISTENCE**

↓  
BREAK LOCAL ORDER (only in processes with 2 or more operation)

-e.g.



### • ATOMIC SOLUTION (EXERCISE)



with this pattern

Read1 influences Read2

IF READ1 = Y then

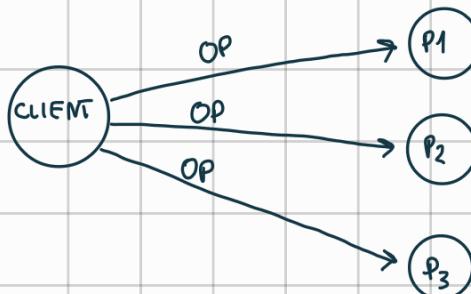
READ2 = Y or GREATER

### REPLICATION

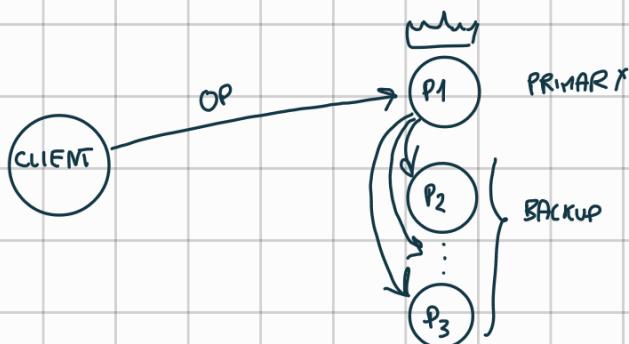
- **CONSISTENCY CRITERIA:** It defines the result returned by an operation → all acceptable runs.
- **LINEARIZABILITY:** An execution on our generic object  $O$  is LINEARIZABLE if for any operation  $OP$  on  $O$  we can find a fictional point (LINEARIZATION POINT) between the start and the end of  $OP$  in which the operation happens instantaneously and is visible to our entire system,

THERE ARE TWO MAIN TECHNIQUES TO IMPLEMENT REPLICATED LINEARIZABLE OBJECTS :

### • ACTIVE REPPLICATION



### • PRIMARY BACKUP (PASSIVE REPPLICATION)



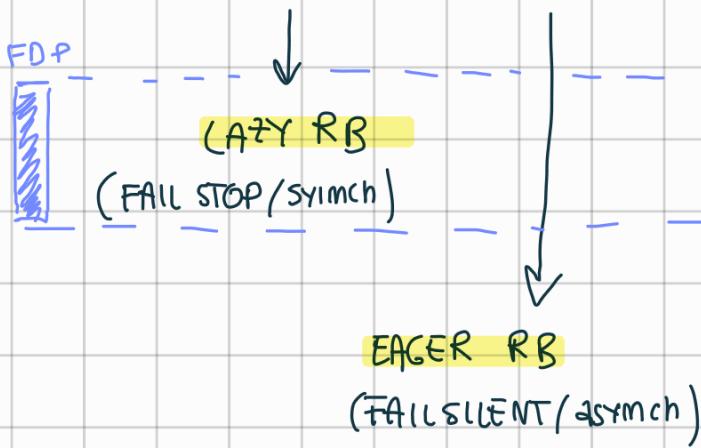
# ALGORITHMS

**BROADCAST**: Communication in a group of processes.

GENERAL SCHEME

- **BEST EFFORT BROADCAST (BEB)**

• **(REGULAR) RELIABLE BROADCAST (RB)**



- **UNIFORM RELIABLE BROADCAST (URB)**

All-ACK URB  
(FAIL STOP/synch)

QUORUM

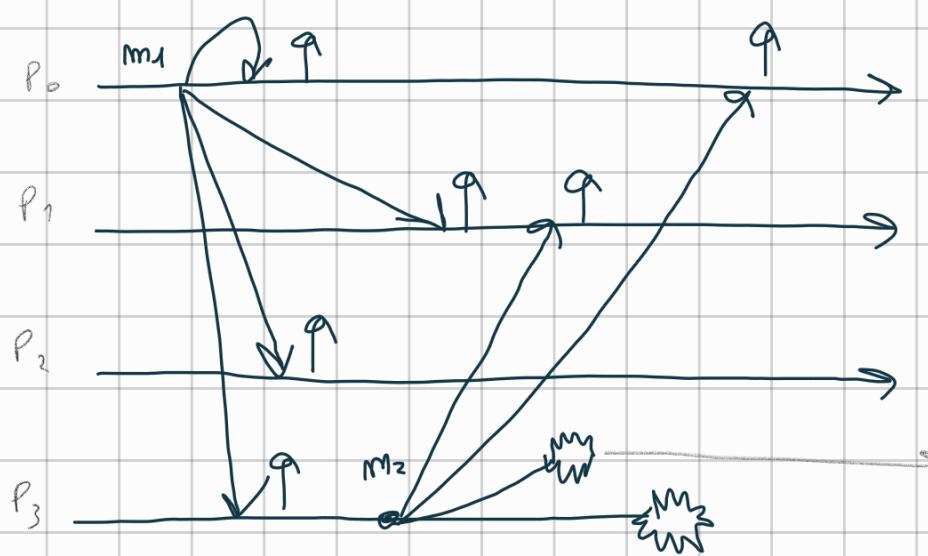
MAJORITY ACK URB  
(FAIL SILENT/asynch)

- **PROBABILISTIC BROADCAST (PB)**



EAGER PB

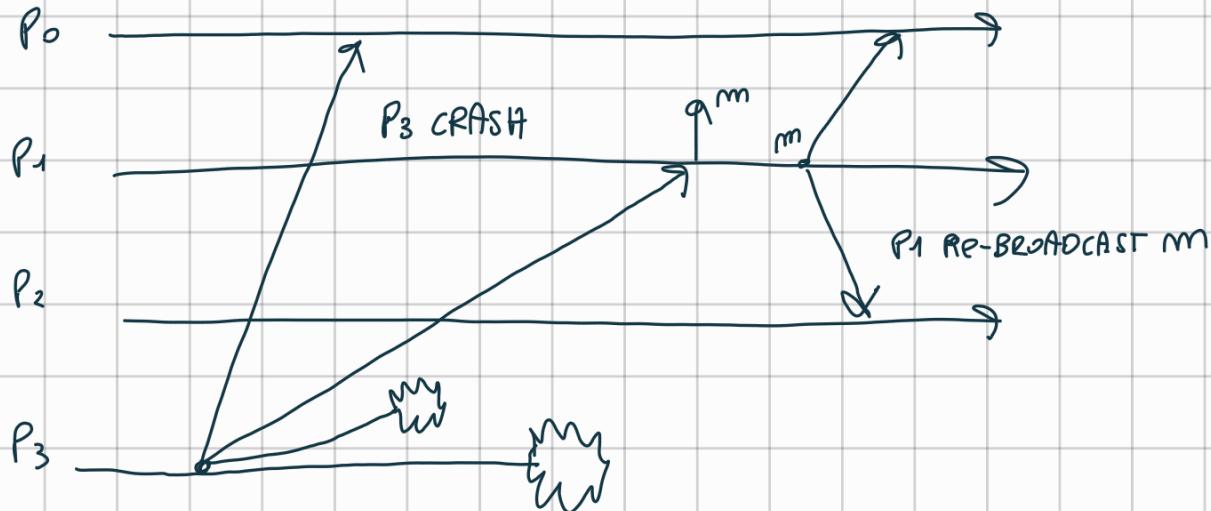
• **BEST EFFORT BROADCAST (BEB)**: In Asynchronous system, if a process want to broadcast, it send a message to all the processes. The delivery of messages is ensured as long as the sender does not fail.



Se un processo  
viene a播散ta,  
prima o poi tutti,  
i corretti deliveremo,  
se播散ta un  
faulty non importa

If  $P_3$  crash before  
to send the message  
 $P_2$  will never receive it.

- (**REGULAR**) **RELIABLE BROADCAST** (RB) : When I receive a message  $m$  from a process, if the process is NOT-CORRECT I RE-BROADCAST the message. When a process crash, all the messages of that process are BROADCASTED.



ALL-CORRECT PROCESS must deliver or nothing  
 ↓  
 (at least)

- **UNIFORM RELIABLE BROADCAST** (URB) ; Also the messages delivered by the FAULTY processes , are eventually delivered by every correct process .



set of messages delivered  
by a **CORRECT** PROCESS



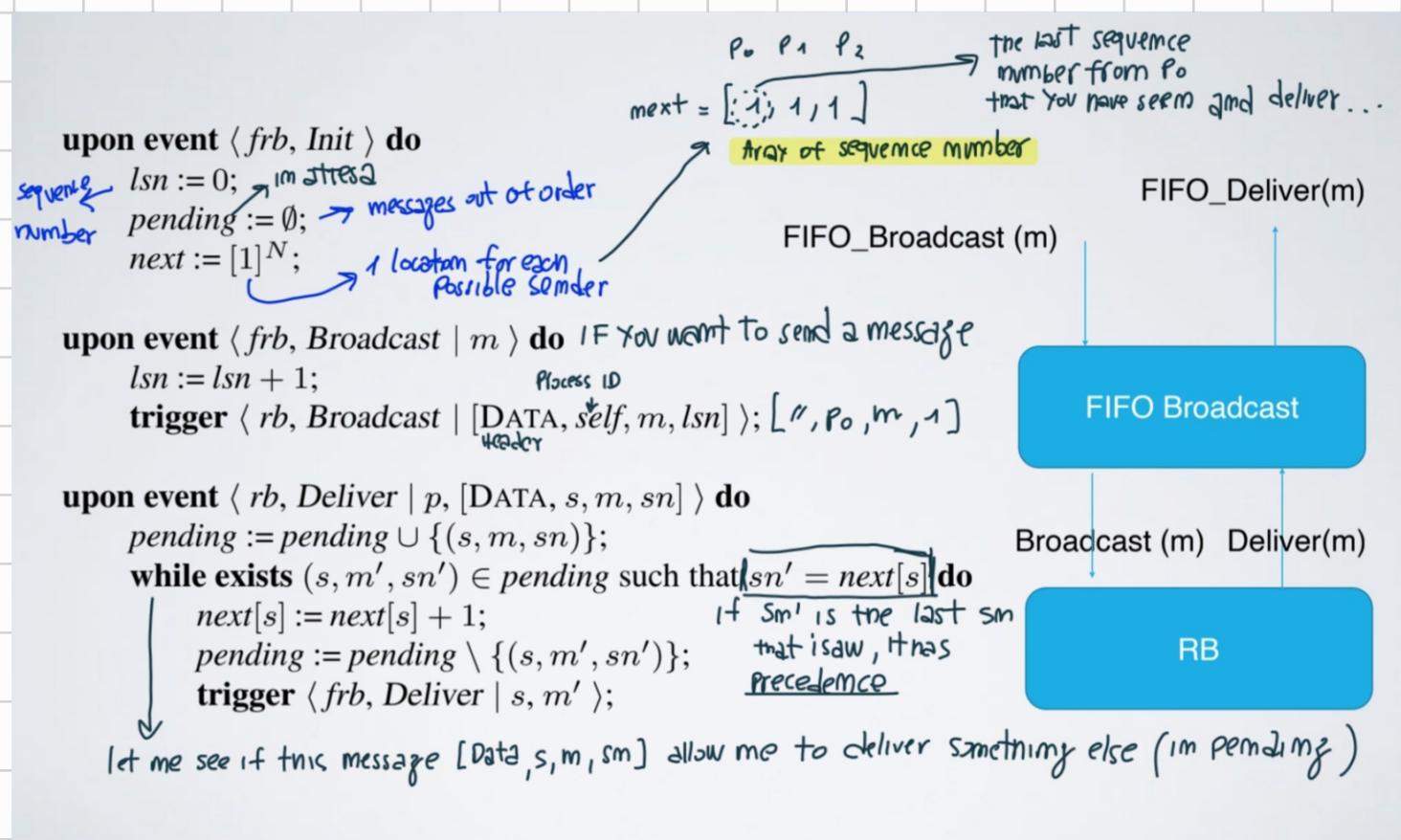
set of messages of the same  
delivered by a **FAULTY** PROCESS

# ORDERED COMMUNICATIONS

To define guarantees about the order of deliveries inside group of processes

- FIFO RELIABLE BROADCAST (FRB)**: The order is just within the single sender, not among different senders.  
FIFO allows to ORDER the messages from the SAME SOURCE.

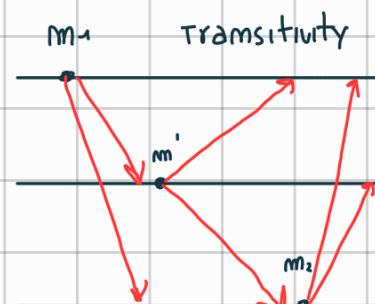
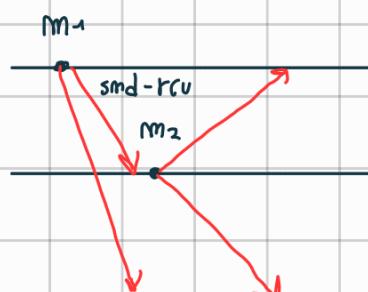
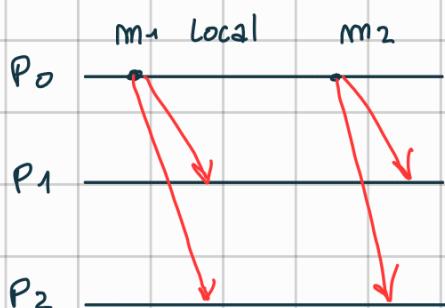
IMPLEMENTATION:



FIFO IS not enough! :

CAUSAL ORDER = FIFO + LOCAL

We NEED CAUSAL ORDER, which guarantees that messages are delivered s.t. they respect all CAUSE-EFFECT relations. IT IS an extension of the HAPPENED-BEFORE RELATION :



→ NO WAIT when you receive the msg (NO PENDING SET)

- **NO-WAIT CAUSAL RELIABLE BROADCAST (crb)**: Every time I broadcast a message I send also the history of past messages (CAUSAL HISTORY) and then add the message to that history. When I deliver a message  $m$ , if in the PAST of the SENDER there is a message  $m'$  that I missed, I deliver it and I add it in my PAST. Only them I can deliver the original message  $m$ .

Properties:

CRB1-CRB9: Same as RB1-RB6 in (regular) reliable broadcast.

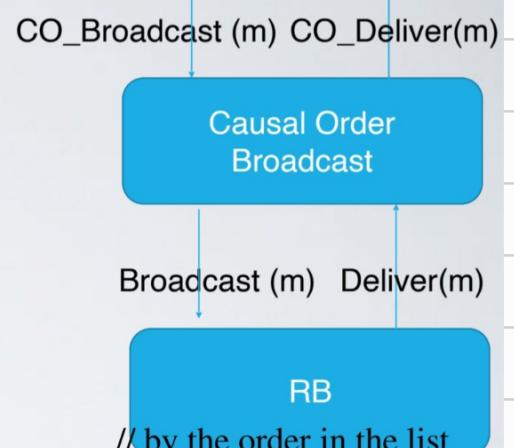
CRB5: CASUAL DELIVERY: For any message  $m_1$  that potentially caused a message  $m_2$ , i.e.  $m_1 \rightarrow m_2$  no process delivers  $m_2$  unless it has already delivered  $m_1$ .

IMPLEMENTATION:

```
upon event < crb, Init > do
    delivered := ∅;
    past := [] // All the messages that you have seen

upon event < crb, Broadcast | m > do
    trigger < rb, Broadcast | [DATA, past, m] >;
    append(past, (self, m));
    ↑ SENDER PAST

upon event < rb, Deliver | p, [DATA, mpast, m] > do
    if  $m \notin$  delivered then if I haven't already delivered it
        forall  $(s, n) \in mpast$  do Scan the list in ORDER
            if  $n \notin$  delivered then → I didn't deliver
                trigger < crb, Deliver | s, n >;
                delivered := delivered ∪ {n};
                if  $(s, n) \notin$  past then → my PAST list
                    append(past, (s, n));
    trigger < crb, Deliver | p, m >; → only them can deliver the original message m
    delivered := delivered ∪ {m};
    if  $(p, m) \notin$  past then
        append(past, (p, m));
```



PROBLEM: Each message travels together with all its CAUSAL HISTORY  
↳ REQUIRES INFINITE MEMORY.

- **IMPROVED NO-WAIT CAUSAL**: Maintaining all the messages in PAST is expensive, so I can discard a message that has been delivered by all correct processes.

PROPERTIES: Same as above

IMPLEMENTATION:

SAME AS BEFORE

```

upon event < crb, Broadcast | m > do
    trigger < rb, Broadcast | [DATA, past, m] >;
    append(past, (self, m));

upon event < rb, Deliver | p, [DATA, mpast, m] > do
    if m ∉ delivered then
        forall (s, n) ∈ mpast do
            if n ∉ delivered then
                trigger < crb, Deliver | s, n >;
                delivered := delivered ∪ {n};
                if (s, n) ∉ past then
                    append(past, (s, n));
            trigger < crb, Deliver | p, m >;
            delivered := delivered ∪ {m};
            if (p, m) ∉ past then
                append(past, (p, m));

```

```

upon event < crb, Init > do
    delivered := ∅;
    past := [];
    correct := Π;
    forall m do ack[m] := ∅;

upon event < P, Crash | p > do
    correct := correct \ {p};
    when I see a message, immediately broadcast ACK to everyone
    upon exists m ∈ delivered such that self ∉ ack[m] do
        ack[m] := ack[m] ∪ {self};
        trigger < rb, Broadcast | [ACK, m] >; SEND ACK
    if the process p crash

upon event < rb, Deliver | p, [ACK, m] > do
    ack[m] := ack[m] ∪ {p};
    collect the ACKs from other people
    when I see the ACK from all correct processes
    upon correct ⊆ ack[m] do
        for a certain msg m
            forall (s', m') ∈ past such that m' = m do
                remove(past, (s', m));
            add the source in the ack set

```

After UNBOUNDED physical time, I remove a msg from PAST.

MESSAGE DELAY: 2 UNIT of time.  
(logical time)

- **WAITING CAUSAL RELIABLE BROADCAST (CRB)**: Instead of send the PAST, I can create a SUMMARY of the PAST which is smaller than PAST. I can use a VECTOR CLOCK of the message => This creates a TRADE OFF => I send smaller messages but I have to WAIT the VECTOR CLOCK tells me that there exists something in the system that I have not seen (like sm of FIFO)

PROPERTIES: Same as above

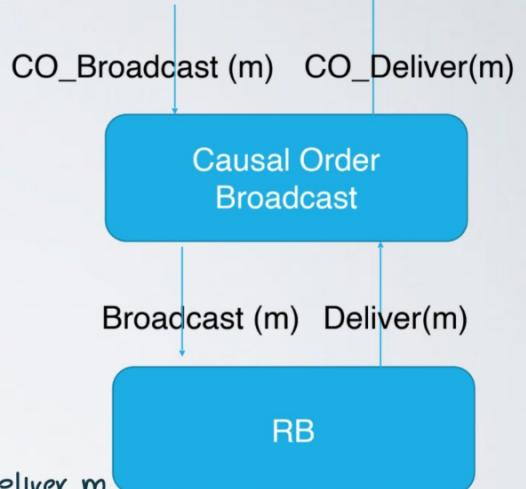
# IMPLEMENTATION:

**upon event**  $\langle \text{crb}, \text{Init} \rangle$  **do**  
 Vector  $V := [0]^N$ ;  $V = [-, -, \underbrace{\dots}_{\substack{\# \text{ of msgs} \\ \text{seen by process } i}}, N]$   
 $\text{clock} \quad \text{lsn} := 0; \rightarrow \text{local sequence number}$   
 $\text{pending} := \emptyset;$

**upon event**  $\langle \text{crb}, \text{Broadcast} \mid m \rangle$  **do**  
 $W := V; \nearrow \text{look for my location}$   
 $W[\text{rank}(\text{self})] := \text{lsn};$   
 $\text{lsn} := \text{lsn} + 1;$   
**trigger**  $\langle \text{rb}, \text{Broadcast} \mid [\text{DATA}, W, m] \rangle;$

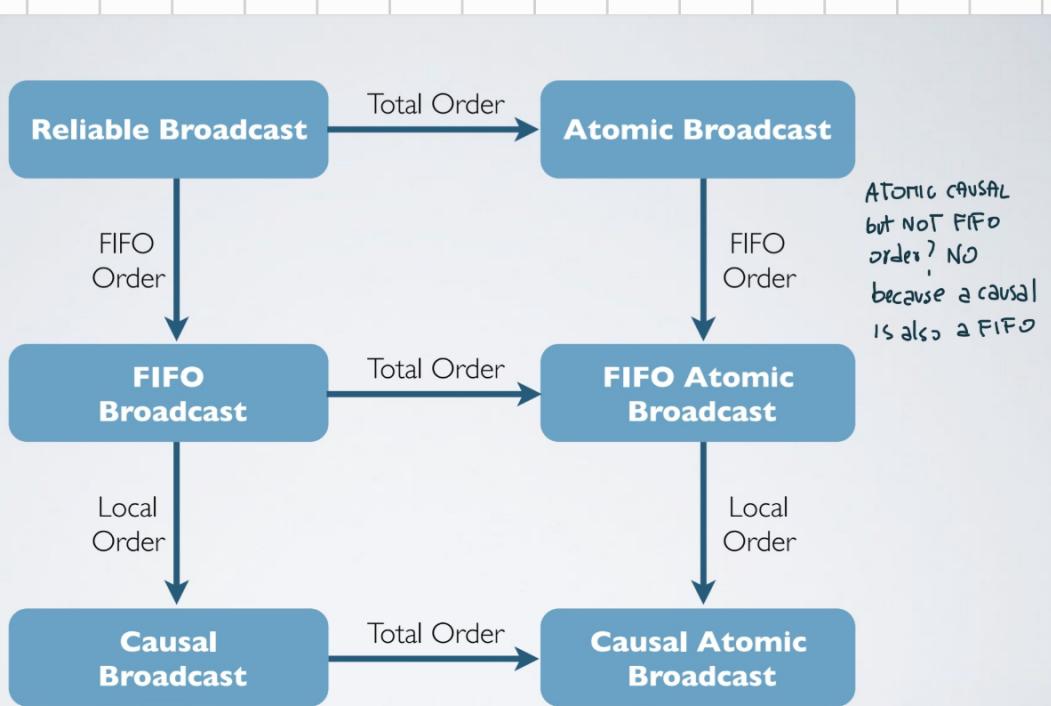
**upon event**  $\langle \text{rb}, \text{Deliver} \mid p, [\text{DATA}, W, m] \rangle$  **do**  
 $\text{pending} := \text{pending} \cup \{(p, W, m)\}; \rightarrow \text{You can not deliver } m \text{ immediately (waiting)}$   
**while exists**  $(p', W', m') \in \text{pending}$  such that  $W' \leq V$  **do**  
 $\text{pending} := \text{pending} \setminus \{(p', W', m')\}; \rightarrow \text{every component must be } \leq$   
 $V[\text{rank}(p')] := V[\text{rank}(p')] + 1; \rightarrow \text{I've seen something from } p', \text{ so I increase its location.}$   
**trigger**  $\langle \text{crb}, \text{Deliver} \mid p', m' \rangle;$

FAIL SILENT



- **TOTAL ORDER (RELIABLE) BROADCAST**: orders all messages, even those from different senders and those that are not causally related. Orthogonal with respect to FIFO and CAUSAL ORDER.

BROADCAST RELATIONSHIP :



# DISTRIBUTED REGISTERS:

READ-ONE WRITE-ALL

- **(1, N) REGULAR REGISTER:** Each process stores a locally copy of the register. When I want to WRITE a value  $v$  into a register, I BEB  $v$  to every one (WRITE-ALL). When I DELIVER  $v$  I update my local copy of the register (READ-ONE) and I SEND an ACK. A WRITE completes when the writer receives an ACK from all the correct processes.

PROPERTIES :

ONRR1 : Termination: If a correct process invokes an operation, then the operation eventually completes,

ONRR2; Validity: A read that is not concurrent with a write returns the last value written; a read that is concurrent with a write returns the last value written or the value concurrently written.

IMPLEMENTATION: (FAIL-STOP)

```
upon event < onrr, Init > do
    val := ⊥;
    correct := Π;
    writeset := ∅; → ACK set
```

```
upon event < P, Crash | p > do
    correct := correct \ {p};
```

```
upon event < onrr, Read > do
    trigger < onrr, ReadReturn | val >;
```

```
upon event < onrr, Write | v > do
    trigger < beb, Broadcast | [WRITE, v] >;
```

```
upon event < beb, Deliver | q, [WRITE, v] > do
    val := v; → UPDATE register
    trigger < pl, Send | q, ACK >;
```

```
upon event < pl, Deliver | p, ACK > do
    writeset := writeset ∪ {p};
```

```
upon correct ⊆ writeset do
    writeset := ∅;
    trigger < onrr, WriteReturn >;
```

NEED FOR PFD for guarantee VAUITY

- **MAJORING VOTING REGULAR REGISTER**: Each process stores a locally copy of the register. When the writer BEB send also a TIMESTAMP, when the reader deliver a write with an old TIMESTAMP he discard it. The WRITER complete the write when he receives ACKS from a majority of processes (QUORUM). To READ the value in the register the reader asks the value from all the processes (BEB READ), Then he takes the value with the greatest timestamps among the processes in quorum.

PROPERTIES: same as above

IMPLEMENTATION: (FAIL SILENT)

```

upon event < onrr, Init > do
  (ts, val) := (0, ⊥);
  wts := 0; → to ts the unrequest
  acks := 0;
  rid := 0; → ID given to each read operation
  readlist := [⊥]N; → contain all the values
    that you want read from other processes
when the writer want to write
upon event < onrr, Write | v > do
  wts := wts + 1;
  acks := 0;
  trigger < beb, Broadcast | [WRITE, wts, v] >;
upon event < beb, Deliver | p, [WRITE, ts', v'] > do
  if ts' > ts then → to discard old values
    (ts, val) := (ts', v'); → ts of the write
    trigger < pl, Send | p, [ACK, ts'] >;
upon event < pl, Deliver | q, [ACK, ts'] > such that ts' = wts do
  acks := acks + 1;
  if acks > N/2 then a QUORUM
    acks := 0;
    trigger < onrr, WriteReturn >;
  
```

```

upon event < onrr, Read > do
  rid := rid + 1;
  readlist := [⊥]N;
  trigger < beb, Broadcast | [READ, rid] >;
  
```

I WANT to READ

```

upon event < beb, Deliver | p, [READ, r] > do
  trigger < pl, Send | p, [VALUE, r, ts, val] >;
  
```

to not accept the old read in the system

```

upon event < pl, Deliver | q, [VALUE, r, ts', v'] > such that r = rid do
  readlist[q] := (ts', v');
  if #(readlist) > N/2 then
    v := highestval(readlist); → you take the value with the highest timestamp
    readlist := [⊥]N;
    trigger < onrr, ReadReturn | v >;
    ↓
    # of location of readlist ≠ ⊥
  
```

PERFORMANCE :

- COMPLEXITY: WRITE/READ  $2N$
- MSG DELAYS: WRITE/READ 2 STEPS

## • **(1, N) ATOMIC REGISTER (OOAR)** : (USING REGULAR REGISTER)

The algo consists of two phases :

PHASE 1 :  $(1, N) \text{ RR} \rightarrow (1, 1) \text{ AR}$

PHASE 2 : SEVERAL  $(1, 1) \text{ AR} \rightarrow (1, N) \text{ AR}$

### **PHASE 1 : FROM $(1, N)$ -RR TO $(1, 1)$ -AR**

$P_1$  is the writer and  $P_2$  is the reader of the  $(1, 1)$  atomic register we want to implement. We use a  $(1, N)$  RR where  $P_1$  is the writer and  $P_2$  the reader.

When  $P_1$  writes, writes the pair  $(\text{value}, \text{timestamp})$  into the underlying RR. The reader  $P_2$  tracks the timestamp of previously read values to avoid to read something old.

### PROPERTIES :

ONAR<sub>1</sub> - ONAR<sub>2</sub> : Same as ONRR<sub>1</sub> - ONRR<sub>2</sub> of  $(1, N)$  RR

ONAR<sub>3</sub> : ORDERING : If a read returns a value  $v$  and a subsequent read returns value  $w$ , then the write of  $w$  does not precede the write of  $v$ .

### IMPLEMENTATION :

#### **Algorithm 4.3:** From $(1, N)$ Regular to $(1, 1)$ Atomic Registers

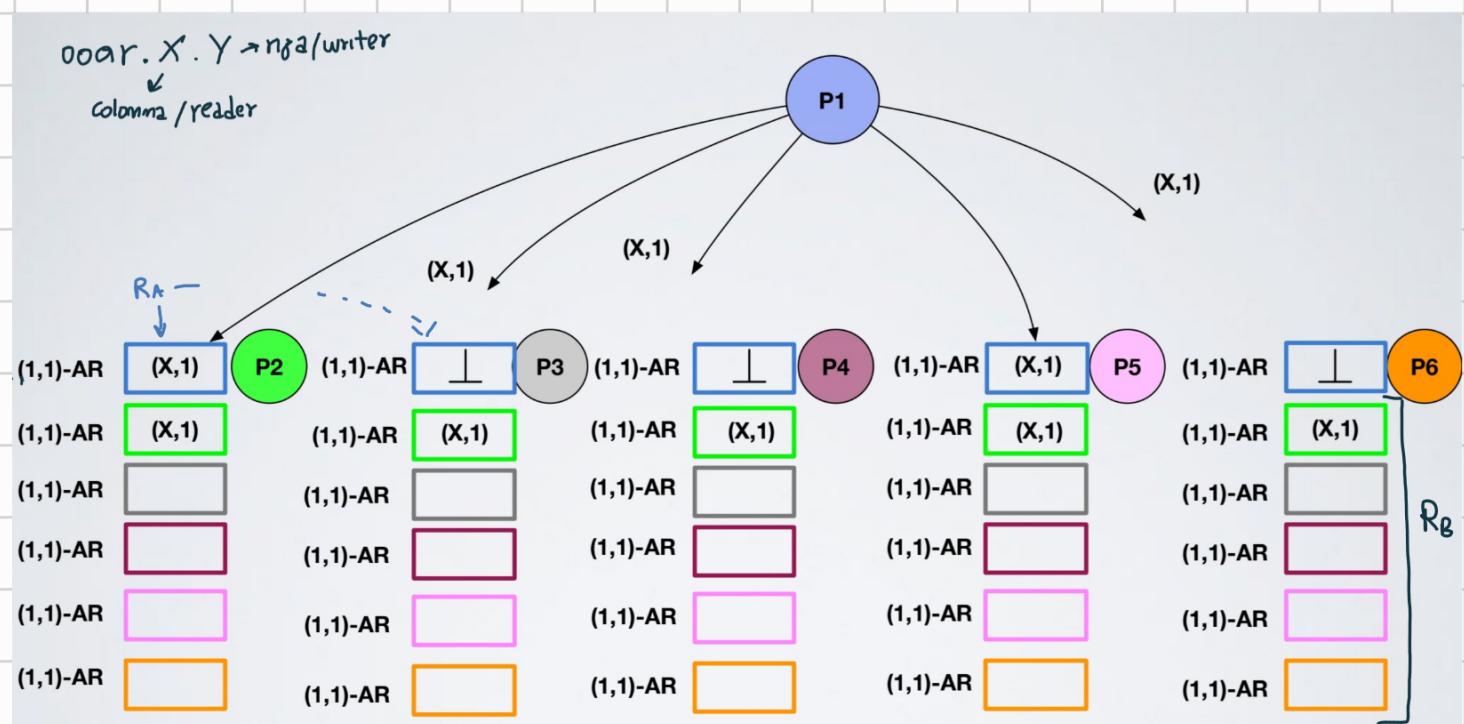
##### Implements:

```
(1, 1)-AtomicRegister, instance ooar.    upon event < onrr, WriteReturn > do
                                         trigger < ooar, WriteReturn >;
```

##### Uses:

```
(1, N)-RegularRegister, instance onrr.    upon event < ooar, Read > do
                                         trigger < onrr, Read >;
upon event < ooar, Init > do
  (ts, val) := (0, ⊥);
  wts := 0;
upon event < ooar, Write | v > do
  wts := wts + 1;
  trigger < onrr, Write | (wts, v) >;
                                         upon event < onrr, ReadReturn | (ts', v') > do
                                         if ts' > ts then
                                           (ts, val) := (ts', v'); → linearization point
                                         trigger < ooar, ReadReturn | val >;
```

## • PHASE 2 - FROM $(1,1)$ -AR TO $(1,N)$ -RR :

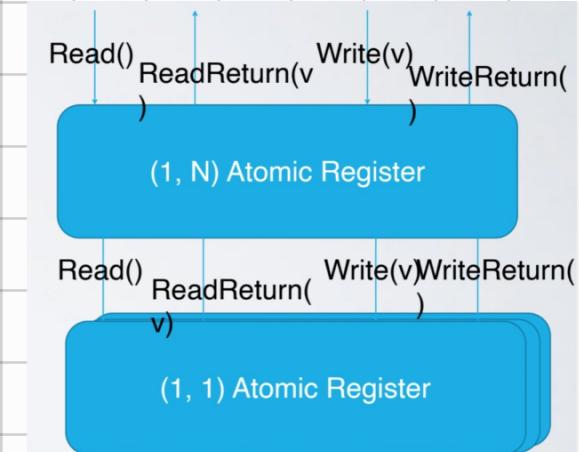


The same case could happen for other processes, so we need a matrix of registers.

$R_A : N(1,1)$  registers that connect the writer to the readers .

$R_B : N^2(1,1)$  registers that connect readers .

The WRITER writes val in registers  $R_A$ , the Reader reads from 1  $R_A$  register and from  $N R_B$  registers (on its column, VERTICAL). He chooses the val with the largest timestamp and writes val in all the  $N R_B$  registers on its row (horizontal). Then he waits for the ACKs and return the value val. When the WRITER receives all the ACKs he completes the WRITE.



PROPERTIES : Same as above

## IMPLEMENTATION:

**upon event**  $\langle \text{onar}, \text{Init} \rangle$  **do**

$ts := 0;$

$acks := 0;$

$writing := \text{FALSE};$

$readval := \perp;$

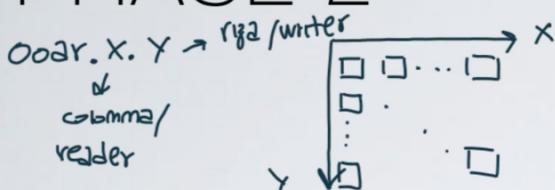
$readlist := [\perp]^N;$

**forall**  $q \in \Pi, r \in \Pi$  **do**

to the process  $q$  (reader) | write  
in the register  $r$  (writer)

Initialize a new instance  $ooar.q.r$  of  $(1, 1)$ -AtomicRegister  
with writer  $r$  and reader  $q$ ;

## PHASE 2



**upon event**  $\langle \text{onar}, \text{Write} | v \rangle$  **do**

$ts := ts + 1;$

$writing := \text{TRUE};$  for all the process  
**forall**  $q \in \Pi$  **do** write in the register self  
of process  $q$  when the original writer  
write in all register

**trigger**  $\langle ooar.q.self, \text{Write} | (ts, v) \rangle;$

**upon event**  $\langle \text{onar}, \text{Read} \rangle$  **do**

**forall**  $r \in \Pi$  **do**

**trigger**  $\langle ooar.self.r, \text{Read} \rangle;$

read in

the

column

X fixed, change Y

**upon event**  $\langle ooar.self.r, \text{ReadReturn} | (ts', v') \rangle$  **do**

$readlist[r] := (ts', v');$

**if**  $\#(readlist) = N$  **then** → read all your column

$(maxts, readval) := \text{highest}(readlist);$

$readlist := [\perp]^N;$

**forall**  $q \in \Pi$  **do**

write on the row

**trigger**  $\langle ooar.q.self, \text{Write} | (maxts, readval) \rangle;$

**upon event**  $\langle ooar.q.self, \text{WriteReturn} \rangle$  **do**

$acks := acks + 1; \# \text{ of completed rum}$

**if**  $acks = N$  **then**

$acks := 0;$

. **if**  $writing = \text{TRUE}$  **then**

**trigger**  $\langle \text{onar}, \text{WriteReturn} \rangle;$

$writing := \text{FALSE};$

**else**

**trigger**  $\langle \text{onar}, \text{ReadReturn} | readval \rangle;$

To discriminate between original writer  
and reader who write and to use  
same code for both case

can tolerate  
any number of failures

- **(1, N) ATOMIC REGISTER : READ-IMPOSE WRITE-ALL** (USING MESSAGES PASSING):  
(FAIL STOP)

The algorithm is a modified version of the Read-One Write-All. The main idea is that the READ operation WRITE: A READ op. IMPOSES to all correct processes to update their local copy of the register with the value read,

unless they store 2 more recent value, when I WRITE i increase WTS and I send  $m = (WTS, VAL)$  to anyone. I MUST WAIT all the ACKS to terminate the WRITE. When I READ locally a  $(WTS, VAL)$ , I also WRITE  $(WTS, VAL)$  to anyone (IMPOSE). A PROCESS accept a WRITE if its TS is less than the one received. SEND ACK in ANY CASE.

### PROPERTIES:

Same as JNR1-ONRR2 of  $(1, N)$  RR

ORDERING: to complete a read operation, the reader process has to be sure that every other process has in its local copy of the register a value with timestamp bigger or equal of the timestamp of the value read. In this way, any successive read could not return an old value

### IMPLEMENTATION:

```

upon event < onar, Init > do
  (ts, val) := (0, ⊥);
  correct := ∏;
  writeset := ∅;
  readval := ⊥;
  reading := FALSE;

upon event < P, Crash | p > do
  correct := correct \ {p};

upon event < onar, Read > do
  reading := TRUE;
  readval := val;
  trigger < beb, Broadcast | [WRITE, ts, val] >; ↑ WRITE ALL

upon event < onar, Write | v > do
  trigger < beb, Broadcast | [WRITE, ts + 1, v] >; ↑
```

```

upon event < beb, Deliver | p, [WRITE, ts', v'] > do
  if ts' > ts then
    (ts, val) := (ts', v');
    trigger < pl, Send | p, [ACK] >; ↗

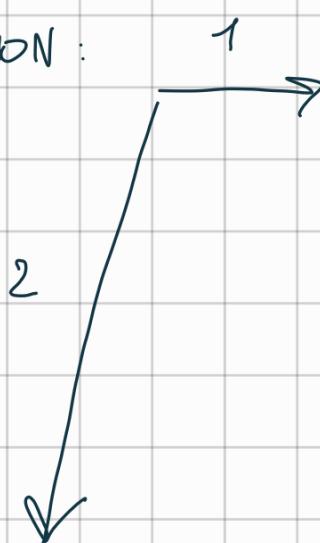
upon event < pl, Deliver | p, [ACK] > then
  writeset := writeset ∪ {p};

upon correct ⊆ writeset do
  writeset := ∅; ↗ If I'm the "original" reader (1:N)
  if reading = TRUE then
    reading := FALSE;
    trigger < onar, ReadReturn | readval >;
  else
    trigger < onar, WriteReturn >; ↙ is reading the writer register (1:N)
```

• **READ-IMPOSE WRITE-MAJORITY (FAIL SILENT)**: The algo is a variation of the Majority voting ( $1, N$ ) RR. The structure is the same as the one for regular and fail-silent but a READ, before to IMPOSE its value to any-one has a QUERY PHASE where he ask everyone the value to read. Then he waits for a MAJORITY of ACKS and he takes the val with the highest timestamp. Only then the READ can impose and terminate the READ after receiving again a majority of ACKS.

PROPERTIES: Same as above

IMPLEMENTATION :



```

upon event < onar, Init > do
  (ts, val) := (0, ⊥);
  wts := 0;
  acks := 0;
  rid := 0;
  readlist := [⊥]N;
  readval := ⊥;
  reading := FALSE;

upon event < onar, Read > do      QUERY PHASE
  rid := rid + 1; local timestamp
  acks := 0;
  readlist := [⊥]N;
  reading := TRUE;
  trigger < beb, Broadcast | [READ, rid] >;
```

↑ the response of my consecutive query

```

upon event < beb, Deliver | p, [READ, r] > do
  trigger < pl, Send | p, [VALUE, r, ts, val] >;
```

```

upon event < pl, Deliver | q, [VALUE, r, ts', v'] > such that r = rid do
  readlist[q] := (ts', v');
  if #(readlist) > N/2 then QUORUM
    (maxts, readval) := highest(readlist); the fresh ONE
    readlist := [⊥]N;           → IMPOSE
    trigger < beb, Broadcast | [WRITE, rid, maxts, readval] >;
```

*the original writer*

```

upon event < onar, Write | v > do
  rid := rid + 1;
  wts := wts + 1;
  acks := 0;
  trigger < beb, Broadcast | [WRITE, rid, wts, v] >;
```

```

upon event < beb, Deliver | p, [WRITE, r, ts', v'] > do
  if ts' > ts then
    (ts, val) := (ts', v');           ↓ "global" timestamp
  trigger < pl, Send | p, [ACK, r] >;
```

*receive ACK*

```

upon event < pl, Deliver | q, [ACK, r] > such that r = rid do
  acks := acks + 1;
  if acks > N/2 then QUORUM
    acks := 0;
    if reading = TRUE then if I'm the reader
      reading := FALSE;
    trigger < onar, ReadReturn | readval >;
  else I'm the register who write (the original ONE)
    trigger < onar, WriteReturn >;
```

PERFORMANCE :

- WRITE : at most  $2N$  MSGS
- READ : // //  $\approx N$  //
- MSG DELAY : WRITE 2 STEPS  
READ 4 STEPS

# ALL REGISTERS PERFORMANCE :

## MESSAGE PASSING IMPLEMENTATIONS!

		(1,n)-Regular Register	(1,n)-Atomic Register		
		READ	WRITE	READ	WRITE
Fail-Stop (P)	Messages	0 (local)	2N	2N	2N
	Steps	0 (local)	2	2	2
	Resiliency	$f < N$		$f < N$	
Fail-Silent	Messages	2N	2N	4N	2N
	Steps	2	2	4	2
	Resiliency	$f < N/2$		$f < N/2$	

## - (N,N) ATOMIC REGISTER (NNAR) :

PROPERTIES : NNAR1: TERMINATION : Same as ONAR1 of (1,N)-AR.

ATOMICITY : Every read operation returns the value that has written most recently in a hypothetical execution, where every failed operation appears to be complete or does not appear to have been invoked at all, and every complete operation appears to have been executed at some instant between its invocation and its completion.

## IMPLEMENTATION:

```

upon event < nnar, Init > do
    val := ⊥;
    writing := FALSE;
    readlist := [⊥]N;
    forall q ∈ Π do
        Initialize a new instance onar.q of (1, N)-AtomicRegister
        with writer q;
upon event < nnar, Write | v > do
    val := v;
    writing := TRUE;
    forall q ∈ Π do
        trigger < onar.q, Read >;
upon event < nnar, Read > do
    forall q ∈ Π do
        trigger < onar.q, Read >;
upon event < onar.q, ReadReturn | (ts', v') > do
    readlist[q] := (ts', rank(q), v');
    if #(readlist) = N then
        (ts, v) := highest(readlist);
        readlist := [⊥]N;
    if writing = TRUE then
        writing := FALSE;
        trigger < onar.self, Write | (ts + 1, val) >
    else
        trigger < nnar, ReadReturn | v >;
upon event < onar.self, WriteReturn > do
    trigger < nnar, WriteReturn >;

```

**CONSENSUS**: Given a set of initial values  $\in \{0,1\}$ , all processes shall decide the same value  $\in \{0,1\}$  based on the initial proposal.

IN SYNCHRONOUS:

**HIERARCHICAL CONSENSUS**: It uses a leader-based strategy. Every ROUND has a leader (which corresponds to the round number). When you are the LEADER you immediately decide your value and send a message to other processes that say "change your mind". When I receive a message from the leader I change my mind. I go to a NEW ROUND (and elect a new leader) in two cases: 1. The old leader died, and I elect a new one; 2. If I receive a message from the leader and I'm the new leader in the time. I become the leader, I decide and I send the msg to other people. The TRICK is that the first correct leader fix the value for everyone.

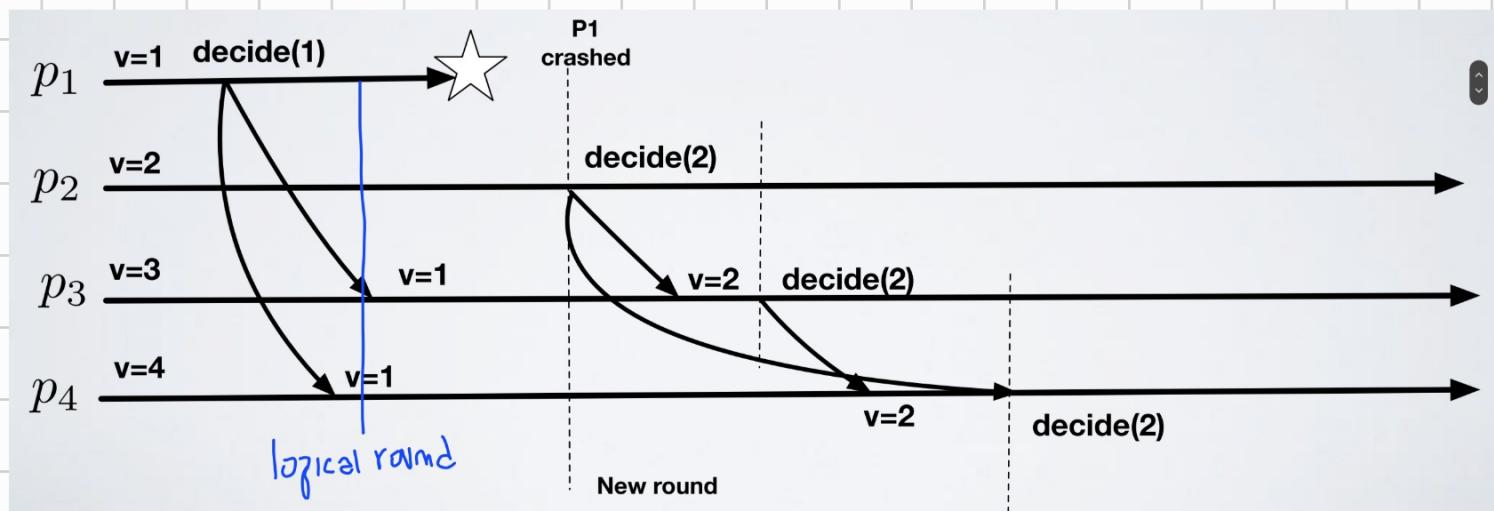
PROPERTIES:

**TERMINATION**: Every correct process eventually decides some value.

**VALIDITY**: If a process decides  $v$ , then  $v$  was proposed by some process.

**INTEGRITY**: No process decides twice.

**AGREEMENT**: No two correct processes decide differently.



# IMPLEMENTATION:

## Algorithm 5.2: Hierarchical Consensus

Implements:

Consensus, instance  $c$ .

Uses:

BestEffortBroadcast, instance  $beb$ ;

PerfectFailureDetector, instance  $\mathcal{P}$ .

$\nearrow$  Initialization Handler

upon event  $\langle c, \text{Init} \rangle$  do

$\text{detectedranks} := \emptyset; \rightarrow$  here you put ID of the dead processes  
 $\text{round} := 1; \rightarrow$  value you want to propose  
 $\text{proposal} := \perp; \text{proposer} := 0;$   
 $\text{delivered} := [\text{FALSE}]^N; \rightarrow$  ID of latest one that made you change your mind  
 $\text{broadcast} := \text{FALSE};$

upon event  $\langle \mathcal{P}, \text{Crash} | p \rangle$  do      return the ID of  $p$   
 $\text{detectedranks} := \text{detectedranks} \cup \{\text{rank}(p)\};$

1 Here You Propose a Value

upon event  $\langle c, \text{Propose} | v \rangle$  such that  $\text{proposal} = \perp$  do  
 $\text{proposal} := v;$

$\downarrow$  from the above  
every process can trigger this handler  
just once

Here You become the leader  
2 upon round = rank(self)  $\wedge \text{proposal} \neq \perp \wedge \text{broadcast} = \text{FALSE}$  do  
 $\text{broadcast} := \text{TRUE}; \rightarrow$  we are the leader

trigger  $\langle beb, \text{Broadcast} | [\text{DECIDED}, \text{proposal}] \rangle;$

3 trigger  $\langle c, \text{Decide} | \text{proposal} \rangle;$

leader crashed OR you received a msg from the leader  
5 upon round  $\in \text{detectedranks} \vee \text{delivered}[\text{round}] = \text{TRUE}$  do  
 $\text{round} := \text{round} + 1;$

when you receive a decided message

4 upon event  $\langle beb, \text{Deliver} | p, [\text{DECIDED}, v] \rangle$  do  
 $r := \text{rank}(p);$

if  $r < \text{rank}(\text{self}) \wedge r > \text{proposer}$  then

$\text{proposal} := v;$

$\text{proposer} := r;$

$\text{delivered}[r] := \text{TRUE};$

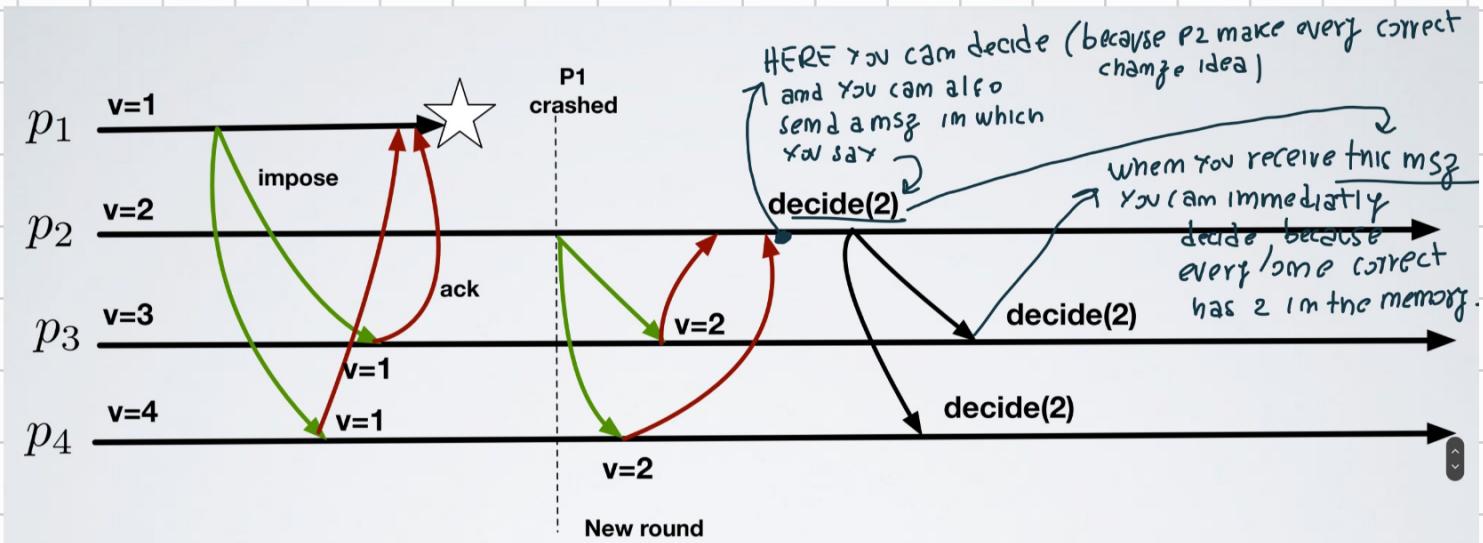
what happened if we remove this?  
we are taking messages from old leaders.

# PERFORMANCE:

COMPLEXITY:  $O(m^2)$  messages

MSG DELAY:  $O(m)$

- UNIFORM HIERARCHICAL CONSENSUS: The leader try to IMPOSE his value, then he waits for the ACKs. He DECIDE only when he receives the ACKs from ALL the CORRECT.



PROPERTIES: UC1-UC3: Same as above

UC4: UNIFORM AGREEMENT: No two processes decide differently

↳ means that must holds also for the crashed process

## IMPLEMENTATION:

Uses:

PerfectPointToPointLinks, instance  $pl$ ;  
BestEffortBroadcast, instance  $beb$ ;  
ReliableBroadcast, instance  $rb$ ;  
PerfectFailureDetector, instance  $P$ .

**upon event**  $\langle uc, \text{Init} \rangle$  do  
 $detectedranks := \emptyset;$   
 $ackranks := \emptyset;$   
 $round := 1;$  → you decide  
 $proposal := \perp; decision := \perp;$  → to set the value  
 $proposed := [\perp]^N;$  → proposed by previous leader

**upon event**  $\langle P, \text{Crash} \mid p \rangle$  do  
 $detectedranks := detectedranks \cup \{\text{rank}(p)\};$

**upon round**  $\in detectedranks$  do → you go to a new round only when the actual leader of the round is dead  
**if**  $proposed[round] \neq \perp$  **then**  
 $proposal := proposed[round];$   
 $round := round + 1;$  →

**4 upon event**  $\langle pl, \text{Deliver} \mid q, [\text{ACK}] \rangle$  do  
 $ackranks := ackranks \cup \{\text{rank}(q)\};$   
→ everyone is either dead or received my msg (and have my value)

**5 upon**  $detectedranks \cup ackranks = \{1, \dots, N\}$  do  
**trigger**  $\langle rb, \text{Broadcast} \mid [\text{DECIDED}, proposal] \rangle;$

**upon event**  $\langle rb, \text{Deliver} \mid p, [\text{DECIDED}, v] \rangle$  such that  $decision = \perp$  do  
 $decision := v;$   
**trigger**  $\langle uc, \text{Decide} \mid decision \rangle;$   
↓ not decided yet

**1 upon event**  $\langle uc, \text{Propose} \mid v \rangle$  such that  $proposal = \perp$  do

$proposal := v;$

you have not decided yet

**2 upon**  $round = \text{rank}(\text{self}) \wedge proposal \neq \perp \wedge decision = \perp$  do  
**trigger**  $\langle beb, \text{Broadcast} \mid [\text{PROPOSAL}, proposal] \rangle;$  → this is my value please change your mind

**3 upon event**  $\langle beb, \text{Deliver} \mid p, [\text{PROPOSAL}, v] \rangle$  do

$proposed[\text{rank}(p)] := v;$

**if**  $\text{rank}(p) \geq round$  **then**

**trigger**  $\langle pl, \text{Send} \mid p, [\text{ACK}] \rangle;$  →

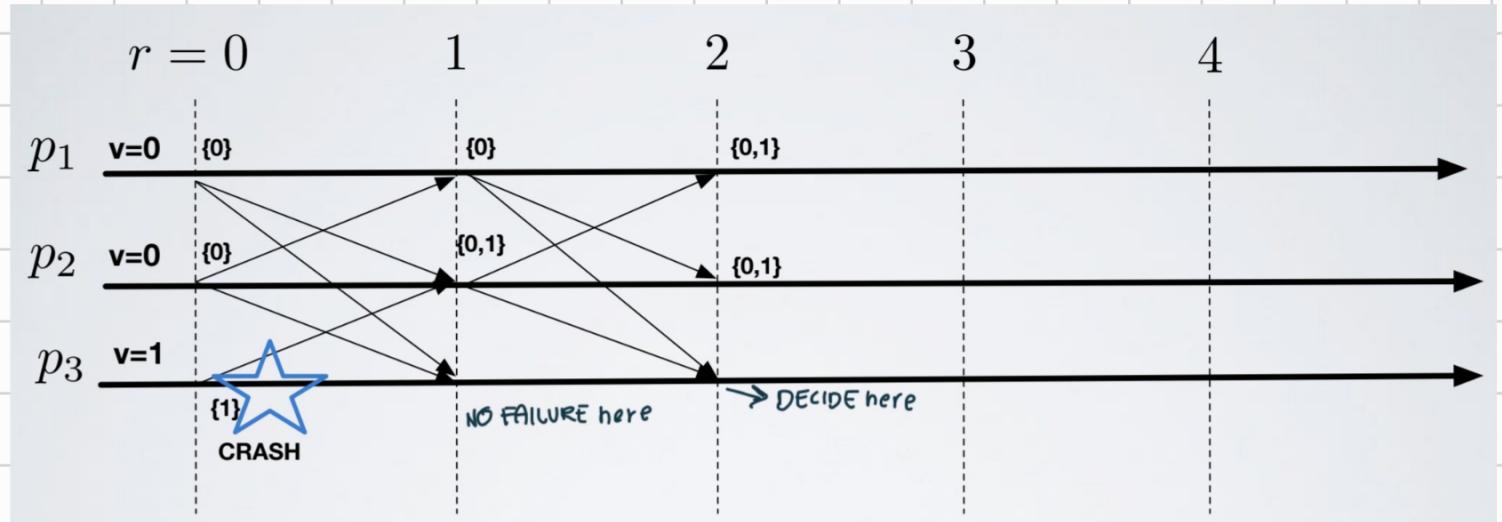
## PERFORMANCE:

COMPLEXITY:  $(1 \text{ BER} + 1 \text{ RB}) \cdot (f+1)$

MSG DELAY: 2 delays for each leader failed, and 3 for each leader that succeeds in sending and delivering one decided =  $O(f)$ .

- FULL-INFORMATION PROTOCOL (FIP)**: with FIP each process collects all possible information on the system, it builds a view of the system and then decides locally. At the beginning of R each process broadcast its set of proposed values (initially, just its value) to all. At the end of R collect all messages and update proposed value as union of the received message. At round R+k take the maximum value in the

proposed value set and decides,



IF a process CRASHES in a round, the algo restart from the first step, to be sure that everyone have the same SET of proposed values. The step continue until in a round nobody CRASHES, so  $K = f+1$ .

IN ASYNCHRONOUS : If  $f > 0$  then consensus is unsolvable.

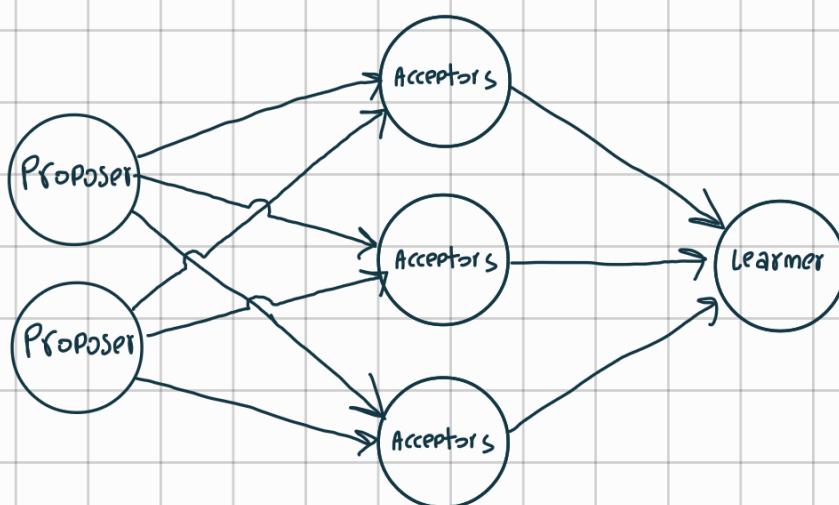
IN EVENTUALLY SYNCHRONOUS :

- **PAXOS** : In Paxos there are 3 actors :

**PROPOSERS** : Only Propose values .

**ACCEPTORS** : processes that must commit in a final decided value .

**LEARNERS** : Passively assist to the decision and obtain the final decided value .



The IDEAS behind the property and the algorithm are pretty complex. So I will show directly them. For more details consult the PROF's slide.

### PROPERTIES :

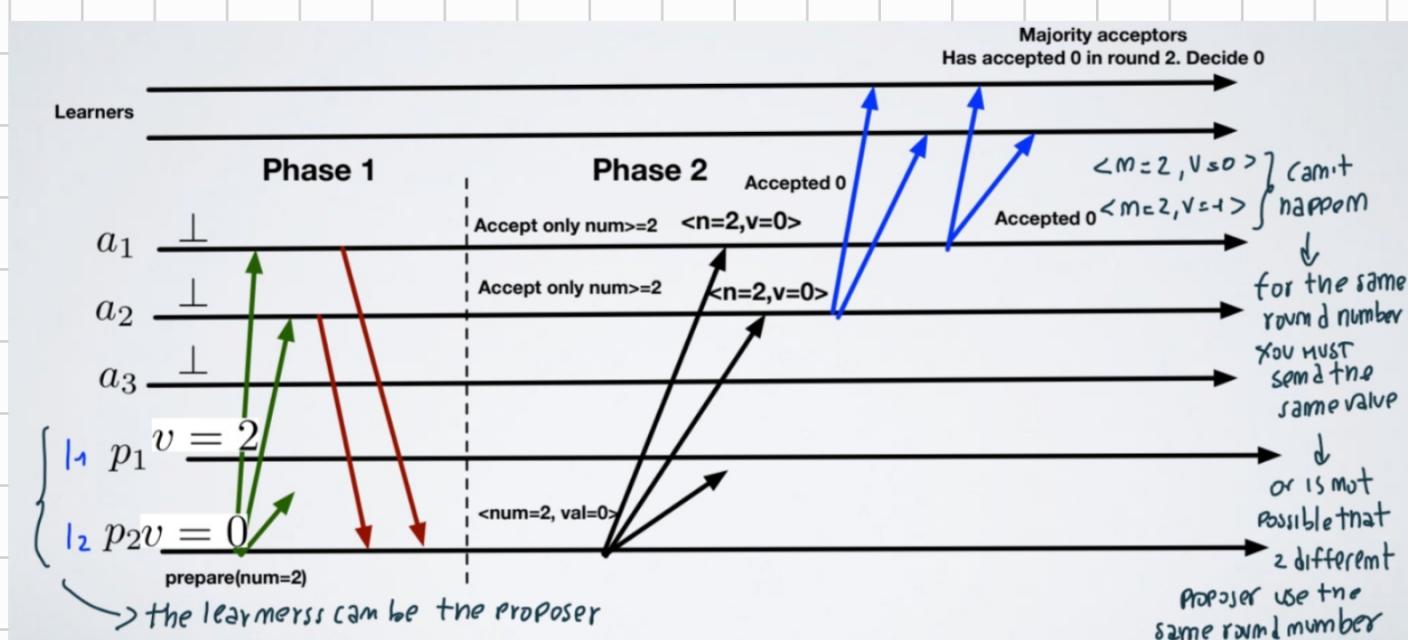
- P<sub>1</sub>: An acceptor must accept the first proposal it receives.
- P<sub>2</sub>: If a proposal with a value  $v$  is chosen, every higher-numbered proposal that is chosen has value  $v$ .
- P<sub>2a</sub>: If a proposal with a value  $v$  is chosen, every higher-numbered proposal that is accepted by any acceptor has value  $v$ .
- P<sub>2b</sub>: If a proposal with value  $v$  is chosen, every higher-numbered proposal issued by any proposer has value  $v$ .

P<sub>2c</sub>: For any  $v$  and  $n$ , if a proposal with value  $v$  and round number  $n$  is issued, then there is a set  $S$  consisting of a majority of acceptors such that either:

- (a) no acceptor in  $S$  has accepted any proposal round numbered less than  $n$ . You are the first in the majority (the majority is empty)
- (b)  $v$  is the value of the highest-numbered proposal among all proposals round numbered less than  $n$  accepted by the acceptors in  $S$ . the majority is not empty, you take all the proposal from the majority and you take the highest number

$$P_{2c} \rightarrow P_{2b} \rightarrow P_{2a} \rightarrow P_2$$

CAN'T BE IMPLEMENTED  
But the first in the implication can, and implies all the others



## IMPLEMENTATIONS :

### PROPOSERS:

- 1: **Constants:**  $\text{failure}$
- 2:  $A, n,$  and  $f.$   $\{A\}$  is the set of acceptors.  $n = |A|$  and  $f = \lfloor (n - 1)/2 \rfloor.$   $\{\text{minority}\}$
- 3: **Init:**
- 4:  $crnd \leftarrow -1$   $\xrightarrow{\text{fictional round}}$  {Current round number}
- 5: **on**  $\langle \text{PROPOSE}, val \rangle$   $\xrightarrow{\text{create a different round number on each proposer}}$
- 6:  $crnd \leftarrow \text{pickNextRound}(crnd)$
- 7:  $cval \leftarrow val$
- 8:  $P \leftarrow \emptyset \rightarrow \text{set of reply from the acceptor}$
- 9: **send**  $\langle \text{PREPARE}, crnd \rangle$  **to**  $A$   $\xrightarrow{\text{FILTER the old REPLY}}$
- 10: **on**  $\langle \text{PROMISE}, \overset{\text{ACK}}{rnd}, vrnd, vval \rangle$  with  $rnd = crnd$  **from** acceptor  $a$   $\xrightarrow{\text{Highest round number and relative value that the acceptor has in memory}}$
- 11:  $P \leftarrow P \cup (vrnd, vval)$
- 12: **on event**  $|P| \geq n - f$   $\xrightarrow{\text{majority}}$
- 13:  $j = \max\{vrnd : (vrnd, vval) \in P\}$   $\xrightarrow{\text{Highest round number}}$
- 14: **if**  $j \geq 0$  **then** there is a value  $\neq \perp$  inside
- 15:  $V = \{vval : (j, vval) \in P\}$   $\xrightarrow{\text{Take the value associated with the highest round number}}$
- 16:  $cval \leftarrow \text{pick}(V)$  {Pick proposed value  $vval$  with largest  $vrnd\}$
- 17: **send**  $\langle \text{ACCEPT}, crnd, cval \rangle$  **to**  $A$   $\xrightarrow{\text{ }}$

### ACCEPTORS:

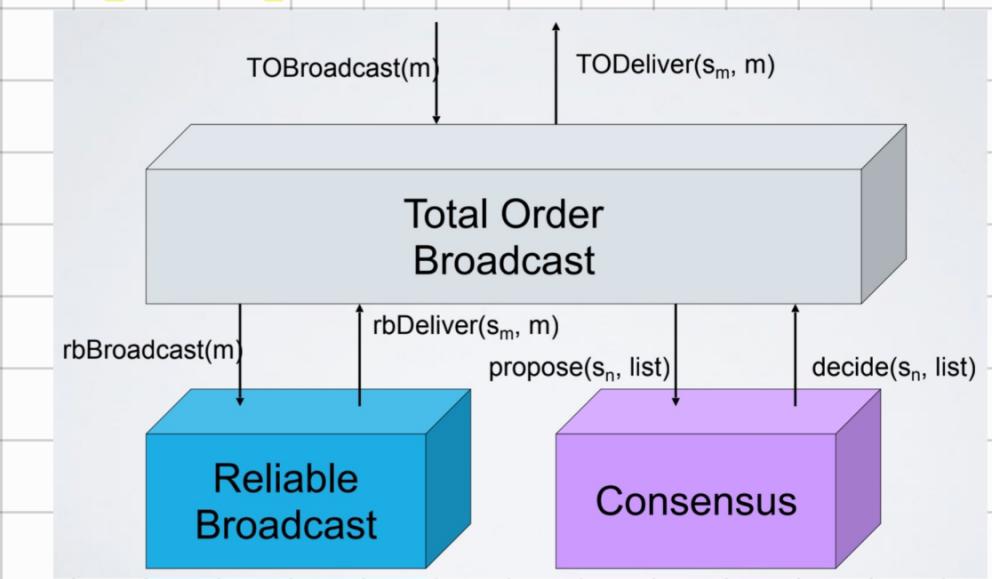
- 1: **Constants:**
- 2:  $L \leftarrow \text{have all the IP of learner}$  {Set of learners}
- 3: **Init:**
- 4:  $rnd \leftarrow -1 \rightarrow \text{send } -1 = \text{send } \perp$
- 5:  $vrnd \leftarrow -1$
- 6:  $vval \leftarrow -1$
- 7: **on**  $\langle \text{PREPARE}, prnd \rangle$  with  $prnd > rnd$  **from** proposer  $p$
- 8:  $rnd \leftarrow prnd$
- 9: **send**  $\langle \text{PROMISE}, \overset{\text{ACK}}{rnd}, vrnd, vval \rangle$  **to** proposer  $p$
- 10: **on**  $\langle \text{ACCEPT}, i, v \rangle$  with  $i \geq rnd$  **from** proposer  $p$
- 11:  $rnd \leftarrow i$   $\downarrow$  it is something NEW?
- 12:  $vrnd \leftarrow i$
- 13:  $vval \leftarrow v$
- 14: **send**  $\langle \text{LEARN}, i, v \rangle$  **to**  $L$

### LEARNERS:

- 1: **Init:**
- 2:  $V \leftarrow \emptyset$  All the msgc we received
- 3: **on**  $\langle \text{LEARN}, (i, v) \rangle$  **from** acceptor  $a$
- 4:  $V \leftarrow V \uplus (i, v)$   $\uparrow$  update the set  $\xrightarrow{\text{majority}}$
- 5: **on event**  $\exists i, v : |\{(i, v) : (i, v) \in V\}| \geq n - f$
- 6:  $v$  is chosen

**TOTAL ORDER (TO) BROADCAST**: It is used to solve the problem in which we want that several replicas of the same data structure are consistent, even if they are updated from different counter in a distributed way. All the replicas will see all the operations in the EXACT SAME ORDER (eventually all the replica will converge on the same value). [to have also a guarantee on time we need synchrony].

### • CONSENSUS-BASED TOTAL-ORDER BROADCAST:



### PROPERTIES :

TOB1      TOB2  
validity, No duplication,      TOB3  
No creation      SAME AS RB & BEB.

TOB4 : Agreement: If a message  $m$  is delivered by some correct process, then  $m$  is eventually delivered by every correct process.

TOB5 : Total order: Let  $m_1$  and  $m_2$  be any two messages and suppose  $p$  and  $q$  are any two correct processes that deliver  $m_1$  and  $m_2$ . If  $p$  delivers  $m_1$  before  $m_2$ , then  $q$  delivers  $m_1$  before  $m_2$ .

IMPLEMENTATION : In the following algorithm the rounds are DESYNCHRONIZED. One process could be at round 1 while another process could be already at round 2.

## Algorithm 6.1: Consensus-Based Total-Order Broadcast

Implements:

TotalOrderBroadcast, **instance** *tob*.

Uses:

ReliableBroadcast, **instance** *rb*;  
Consensus (multiple instances).

**upon event**  $\langle \text{tob}, \text{Init} \rangle$  **do**

*unordered* :=  $\emptyset$ ; *msg*, w/o order  
    *delivered* :=  $\emptyset$ ; To avoid duplication  
    *round* := 1; ROUND BASED  
    *wait* := FALSE;

triggered by the sender

**upon event**  $\langle \text{tob}, \text{Broadcast} \mid m \rangle$  **do**

    trigger  $\langle \text{rb}, \text{Broadcast} \mid m \rangle$ ;

**upon event**  $\langle \text{rb}, \text{Deliver} \mid p, m \rangle$  **do**

    if  $m \notin \text{delivered}$  then

*unordered* := *unordered*  $\cup \{(p, m)\}$ ;  
        ↑  
        waiting to  
        be ordered  
  
        ↑  
        original  
        source  
        YOU NOT DELIVER  
        HERE YET

**upon**  $\text{unordered} \neq \emptyset \wedge \text{wait} = \text{FALSE}$  **do**

*wait* := TRUE;

    Initialize a new instance *c.round* of consensus;

    trigger  $\langle c.\text{round}, \text{Propose} \mid \text{unordered} \rangle$ ;

a SET (can be codified as a bit)

// a process *P* is starting  
// round 1, then  
// round 2...

when the CONSENSUS decide (eventually, it will happen)

**upon event**  $\langle c.r, \text{Decide} \mid \text{decided} \rangle$  such that  $r = \text{round}$  **do**

**forall**  $(s, m) \in \text{sort}(\text{decided})$  **do** → sort locally // by the order in

    trigger  $\langle \text{tob}, \text{Deliver} \mid s, m \rangle$ ;

// in the resulting sorted list

*delivered* := *delivered*  $\cup \text{decided}$ ;

*unordered* := *unordered*  $\setminus \text{decided}$ ;

*round* := *round* + 1;

*wait* := FALSE;

↑  
decided set  
is the same for  
all the correct process

\* We need rb Broadcast to do the dissemination:  
to ensure that will be a time (eventually)  
in which everyone will start with a specific message

## CONSENSUS / RELIABLE BROADCAST RELATION :

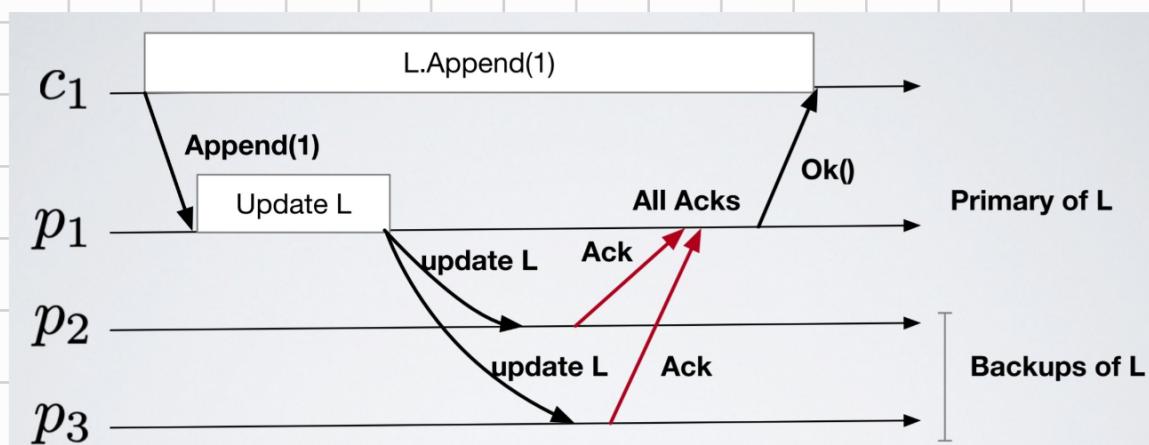
Consensus	Uniform	Non Uniform
Reliable Broadcast	Uniform	Non-Uniform
Uniform	TO	Non-uniform TO
Non Uniform	<b>Uniform</b> <b>TO</b>	Non-uniform TO

**SOFTWARE REPLICATION** : Suppose we want to implement a remote object (e.g. a LIST) in a server, and the clients interact with it to do some operation on the LIST. REMOTE because the object is shared between the clients (if client 1 does an op. client 2 will see it). This is a CENTRALIZED structure, so the PROBLEM is that if the server crashes the object is destroyed forever. So to increase the probability that our system works, we need to have an high available object distributed in several servers → we need REPLICATION.

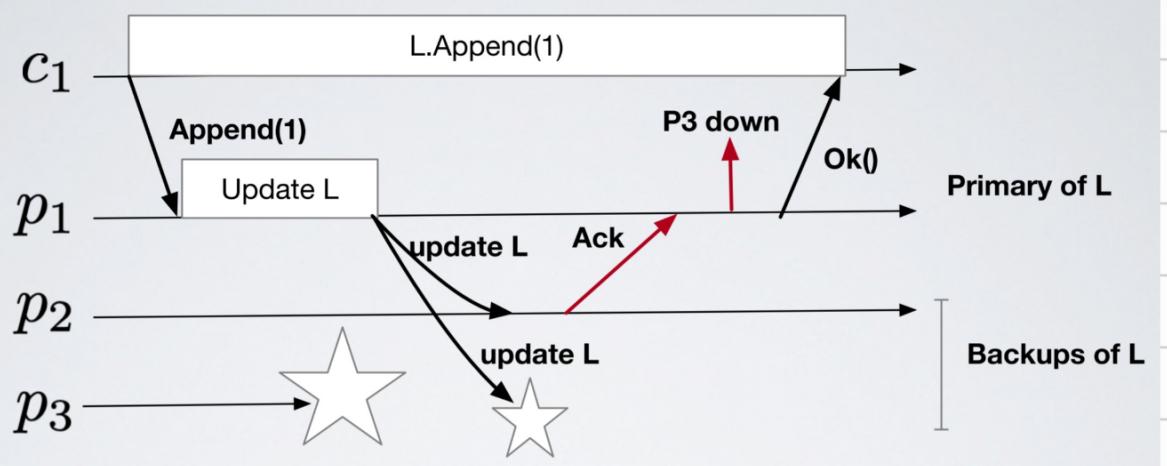
## SYNCHRONOUS SYSTEM

- **PRIMARY BACKUP**: Replicas are partitioned in two set: PRIMARY and BACKUP. The former contains only one element (the LEADER). It is the only that receives invocation from the clients and sends back the answers. The BACKUP set interact with PRIMARY only. It is used to guarantee fault tolerance by replacing a Primary when crashes.

CASE 1 - NO CRASH : CLIENT ask action to primary ; PRIMARY send UPDATE request to BACKUP, them he wait for all the ACKS. Only then PRIMARY send the RES to the client.

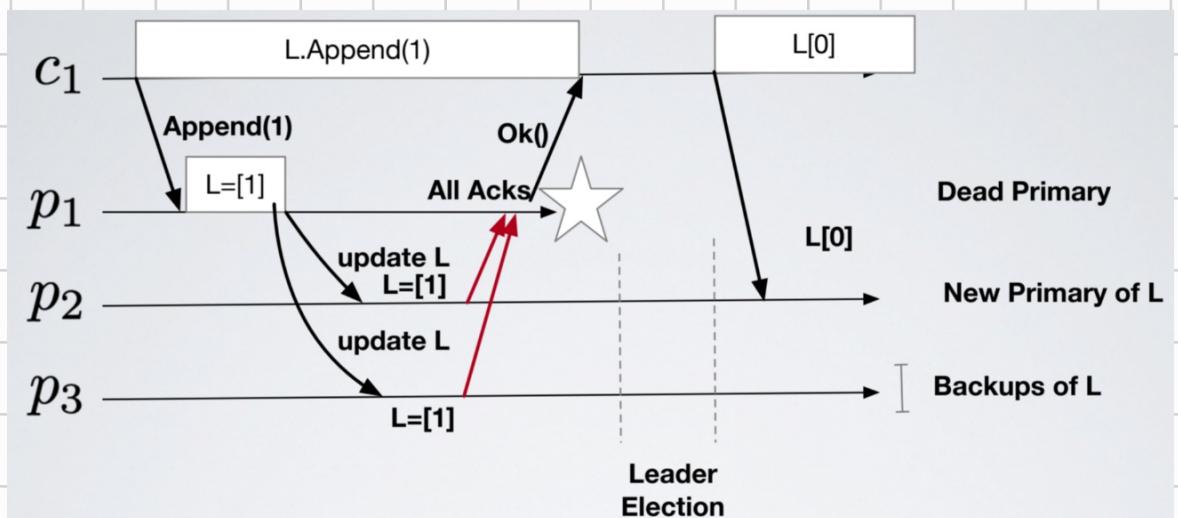


CASE 2 - BACKUP CRASH : The crash is detected by the FD. PRIMARY doesn't wait ACK from crashed process (wait until crashed ACKS = T).

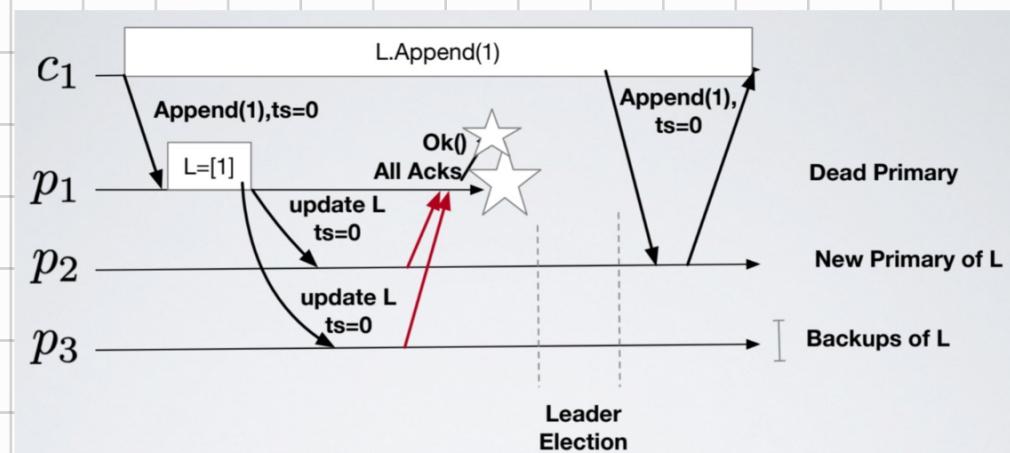


CASE 3 - PRIMARY CRASH: There are 3 scenarios,

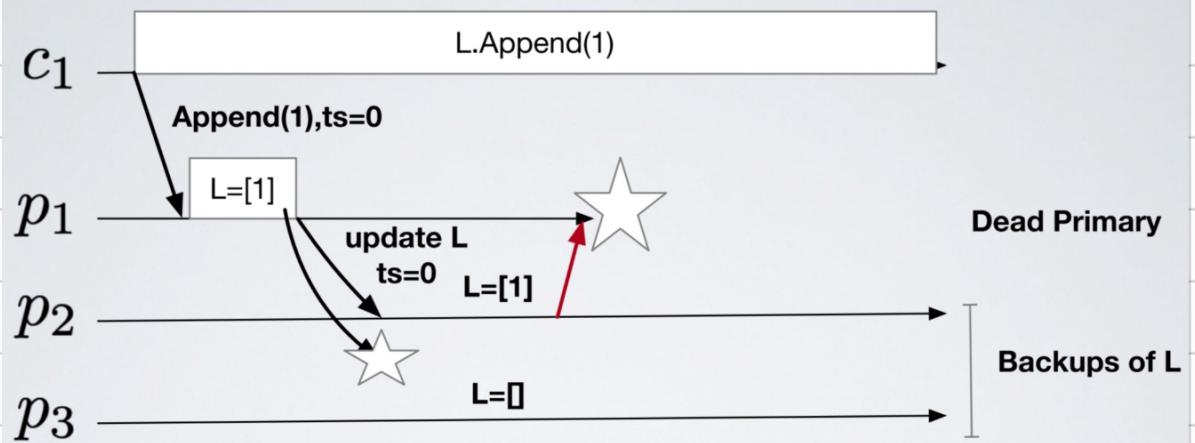
SCENARIO 1: PRIMARY FAILS after the operation is completed and when is not handling any operation (after the response to the client). We elect a new LEADER which has the latest COPY of the object and can answer to new requests.



SCENARIO 2: PRIMARY FAILS after it has received all the ACKS and before the client receives the response. the CLIENT TIMEOUTS and send again the request to the new primary. the NEW PRIMARY will recognize the request re-issued by the client as already processed and sends back result without updating the replicas. (REQUEST have to be UNIQUE  $\langle id, ts \rangle$ )



SCENARIO 3 : PRIMARY FAILS while sending update messages and before receiving all the ACK . If no backups receive the op is fine. We elect someone now, that will execute again the operation upon requests. The problematic case is if the update is received by only some of the backups. We have to take the backup with the most updated state as PRIMARY.



### - ELECTING THE MOST UPDATED BACKUP:

upon event INIT

```

primary = ⊥
deadprimary = ⊥
Crashed = ∅
update_counter = 0      Array or set
ACKS = [∅, ∅, ..., ∅]   ACKS[p] is the set of acks associated with p

```

▷ size equal to  $|\Pi|$

upon event CRASH( $p$ )

$Crashed = Crashed \cup \{p\}$       dead primary

start the election  
upon event  $primary \in Crashed$

BebCast(< ELECTION, primary, myID, update\_counter >)  
 $primary = \perp$

sending process

upon event BEB DELIVER(< ELECTION, IDDEADPRIMARY, ID, CNT > from  $p$ )

if  $primary == IDDEADPRIMARY$  then

$Crashed = Crashed \cup \{primary\}$

BebCast(< ELECTION, primary, myID, update\_counter >) → I start the election also myself  
 $primary = \perp$

$ACKS[IDDEADPRIMARY] = ACKS \cup \{< ELECTION, IDDEADPRIMARY, ID, CNT >\}$

upon event EXISTS  $j$  SUCH THAT  $ACKS[j] \cup Crashed = \Pi$

$primary =$  take the lowest ID with the highest counter in  $ACKS[j]$

$ACKS[i] = \emptyset$

- ACTIVE REPUTATION : We use TOTALORDER broadcast, including the clients! It doesn't need recovery actions upon the failure of a replica.

## EVENTUALLY SYNCHRONOUS SYSTEM

- PRIMARY BACKUPS: RAFT : It's quite long and complex. Look at the original slide.

BYZANTINE FAILURES: A Byzantine process may :

- deviate arbitrarily from the instructions of the algorithm:
  - creating fake messages
  - dropping messages
  - delay the deliveries ...

It can send messages faking an arbitrary sender meaning that it can impersonate other processes. So AUTHENTICATION IS NECESSARY.

- AUTHENTICATED PERFECT LINK

### PROPERTIES

AL1: REUABLE DELIVERY : If a correct process sends a msg  $m$  to a correct process  $q$ , then  $q$  eventually delivers  $m$ .

AL2: NO DUPLICATION: No message is delivered by a correct process more than once.

AL3: AUTHENTICITY : If some correct process  $q$  delivers a message  $m$  with sender  $P$  and process  $p$  is correct, then  $m$  was previously sent to  $q$  by  $P$ .

## ALGORITHM ON STUBBORN

### Algorithm 2.4: Authenticate and Filter

#### Implements:

AuthPerfectPointToPointLinks, instance  $al$ .

#### Uses:

StubbornPointToPointLinks, instance  $sl$ .

**upon event**  $\langle al, Init \rangle$  **do**  
 $delivered := \emptyset;$

**upon event**  $\langle al, Send | q, m \rangle$  **do**  
 $a := \text{authenticate}(self, q, m);$  → process  $P$  invoke the oracle auth with msg  $m$  and dest  $q$  and  
trigger  $\langle sl, Send | q, [m, a] \rangle$ ; obtain object  $a$

**upon event**  $\langle sl, Deliver | p, [m, a] \rangle$  **do**  
**if** verifyauth( $self, p, m, a$ )  $\wedge m \notin delivered$  **then**  
 $delivered := delivered \cup \{m\};$  →  $q$  invoke the oracle call with msg  $m$ , sender  $p$  and object  $a$   
trigger  $\langle al, Deliver | p, m \rangle;$  → oracle answer yes iff



WE DON'T define any "UNIFORM" variant of primitive in Byz.

WE CAN'T build FAILURE DETECTOR for Byz.

**BYZANTINE BROADCAST**: Broadcast is a fundamental primitive. If the sender is Byz. it can send different messages with the same broadcast!



So we need a broadcast primitive that prevents the Byz. to make HONEST processes deliver different messages.

- **BYZANTINE CONSISTENT BROADCAST**: This is a ONE-SHOT algorithm. It means that we execute a single instance of the algo x time.

PROPERTIES:

BCB1: VAULTY: If a correct process  $P$  broadcast a message  $m$ , then every correct process eventually delivers  $m$ .

BCB2: NO DUPLICATION: Every correct process delivers at most one message.

BCB3: INTEGRITY: If some correct process delivers a message  $m$  with sender  $P$  and process  $P$  is correct, then  $m$  was previously broadcast by  $P$ .

BCB4: CONSISTENCY: If some correct process delivers a message  $m$  and another correct process delivers message  $m'$ , then  $m = m'$ .

IS POSSIBLE THAT  
SOMEONE DOESN'T  
DELIVER.

An ECHO is a relay of the original message.

# IMPLEMENTATION:

## Algorithm 3.16: Authenticated Echo Broadcast

Implements:

ByzantineConsistentBroadcast, instance  $bcb$ , with sender  $s$ .

Uses:

AuthPerfectPointToPointLinks, instance  $al$ .

is known by everyone

**upon event**  $\langle bcb, Init \rangle$  **do**

$sentecho := \text{FALSE}$ ;  
     $delivered := \text{FALSE}$ ;  
     $echos := [\perp]^N$ ;

**upon event**  $\langle bcb, Broadcast | m \rangle$  **do**

**forall**  $q \in \Pi$  **do**  
        **trigger**  $\langle al, Send | q, [\text{SEND}, m] \rangle$ ;

// only process  $s$

**upon event**  $\langle al, Deliver | p, [\text{SEND}, m] \rangle$  **such that**  $p = s$  **and**  $sentecho = \text{FALSE}$  **do**

$sentecho := \text{TRUE}$ ;  
    **forall**  $q \in \Pi$  **do**  
        **trigger**  $\langle al, Send | q, [\text{ECHO}, m] \rangle$ ;

RELAY of the original message

**upon event**  $\langle al, Deliver | p, [\text{ECHO}, m] \rangle$  **do**

**if**  $echos[p] = \perp$  **then**  
         $echos[p] := m$ ;

**upon exists**  $m \neq \perp$  **such that**  $\#(\{p \in \Pi \mid echos[p] = m\}) > \frac{N+f}{2}$   
    **and**  $delivered = \text{FALSE}$  **do**

$delivered := \text{TRUE}$ ;  
        **trigger**  $\langle bcb, Deliver | s, m \rangle$ ;

Correctness is guaranteed  
as long as  $N > 3f$

# PERFORMANCE:

COMPLEXITY: At most  $m^2$  messages

MSG DELAYS: 2 STEP.

## • BYZANTINE RELIABLE BROADCAST:

### PROPERTIES:

BRB1 - BRB6: Same as above.

BRB7; TOTALITY: If some message is delivered by any correct process, every correct process eventually delivers a message.

↑  
ALL OR NOTHING

# IMPLEMENTATION:

## Algorithm 3.18: Authenticated Double-Echo Broadcast

### Implements:

ByzantineReliableBroadcast, instance  $brb$ , with sender  $s$ .

### Uses:

AuthPerfectPointToPointLinks, instance  $al$ .

**upon event**  $\langle brb, Init \rangle$  **do**

```

    sentecho := FALSE;
    sentready := FALSE;
    delivered := FALSE;
    echos :=  $[\perp]^N$ ;
    readyss :=  $[\perp]^N$ ;
  
```

**upon event**  $\langle brb, Broadcast | m \rangle$  **do**

```

    forall q in  $\Pi$  do
      trigger  $\langle al, Send | q, [SEND, m] \rangle$ ;
  
```

**upon event**  $\langle al, Deliver | p, [SEND, m] \rangle$  **such that**  $p = s$  **and**  $sentecho = FALSE$  **do**

```

    sentecho := TRUE;
    forall q in  $\Pi$  do
      trigger  $\langle al, Send | q, [ECHO, m] \rangle$ ;
  
```

**upon event**  $\langle al, Deliver | p, [ECHO, m] \rangle$  **do**

```

    if echos[p] =  $\perp$  then
      echos[p] := m;
  
```

```

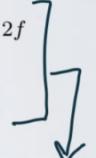
upon exists  $m \neq \perp$  such that  $\#\{p \in \Pi \mid echos[p] = m\} > \frac{N+f}{2}$ 
and  $sentready = FALSE$  do
  sentready := TRUE;
  forall q in  $\Pi$  do
    trigger  $\langle al, Send | q, [READY, m] \rangle$ ;

upon event  $\langle al, Deliver | p, [READY, m] \rangle$  do
  if readyss[p] =  $\perp$  then
    readyss[p] := m;

upon exists  $m \neq \perp$  such that  $\#\{p \in \Pi \mid readyss[p] = m\} > f$ 
and  $sentready = FALSE$  do
  sentready := TRUE;
  forall q in  $\Pi$  do
    trigger  $\langle al, Send | q, [READY, m] \rangle$ ;

upon exists  $m \neq \perp$  such that  $\#\{p \in \Pi \mid readyss[p] = m\} > 2f$ 
and  $delivered = FALSE$  do
  delivered := TRUE;
  trigger  $\langle brb, Deliver | s, m \rangle$ ;
  
```

NO DUPLICATION  
with the FLAG



**BYZANTINE AGREEMENT**: In this setting we assume authenticated channels (MAC) but no signatures. We assume also that  $f < N/3$  (if some third or more of the processes are byzantine, there is no algorithm solving consensus).  $\rightarrow$  number of correct

**KING ALGORITHM**: Because  $f < N/3 \Rightarrow N > 3f \Rightarrow N \geq 3f+1$ . It can be seen as an adaption of Hierarchical Consensus for Byt. This algo works in synchronous with no signatures and with MAC. The algo runs for  $f+1$  phases. In each phase  $j$  there is a KING (e.g. the KING of phase  $j$  is the process with  $ID=j$ ). Each phase is divided in 3 rounds, the VOTE ROUND where each correct process broadcast a variable  $x$ , PROPOSE ROUND where if a process received at least  $n-f$  votes for a value  $y$ , it broadcast  $y$ , otherwise do nothing. At the end of this round, if a value  $z$  reaches  $f+1$  proposals, a correct updates its value  $x$  to  $z$ . The KING ROUND is where the KING is the only which broadcast  $x$ .

At the end of this round, a process drops the king's value if he has seen a value  $m-f$  times (doesn't trust the king); choose king's value otherwise.

## PROPERTIES :

TERMINATION: Every correct eventually decides.

ALL-SAME VAULTY (WEAK VAULTY): If all correct propose  $v$ , the only decision is  $v$ .

INTEGRITY: No correct decides twice.

AGREEMENT: All correct decides the same value.

## IMPLEMENTATION :

```
1:  $x = \text{my input value}$ 
2: for phase = 1 to  $f + 1$  do
    if some value( $y$ ) received at least  $n - f$  times then
        Broadcast propose( $y$ )
    end if
    if some propose( $z$ ) received more than  $f$  times then
         $x = z$  ( $f+1$ )
    end if
    King
    Let node  $v_i$  be the predefined king of this phase  $i$ 
    The king  $v_i$  broadcasts its current value  $w$ 
    if received strictly less than  $n - f$  propose( $y$ ) then
         $x = w \rightarrow$  update to the value of the king
    end if
end for
```

Normal Process

## PERFORMANCE :

COMPLEXITY :  $O((f+1) \cdot m^2)$

MSG DELAY :  $3(f+1)$