

Fault-tolerant broadcasts

Part II: Algorithms

Fault-tolerant broadcasts

■ Reliable broadcast

– System model

- $\Pi = \{ p_1, p_2, \dots, p_n \}$
- Asynchronous system (i.e., no timing assumptions)
- Crash-stop failure model (at most f failures)
- Quasi-reliable channels
 - If p sends a message m to q and both processes are correct, then q eventually receives m

© Fernando Pedone

Fault-tolerant broadcasts

– Non-uniform reliable broadcast

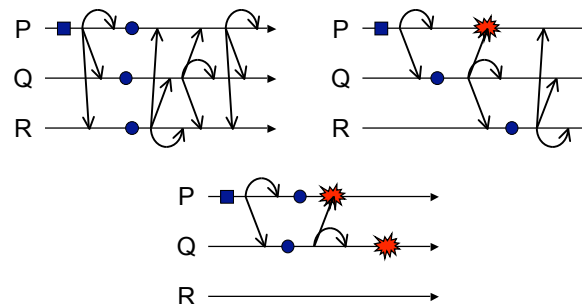
Agreement: if a correct process delivers m , then all correct processes eventually deliver m .

broadcast(R, m) works as follows:
 send(m) to all
 upon receive(m) for the first time do
 deliver(R, m)
 send(m) to all

© Fernando Pedone

Fault-tolerant broadcasts

– Non-uniform reliable broadcast (cont'd)



© Fernando Pedone

Fault-tolerant broadcasts

- Uniform reliable broadcast

Uniform agreement: if a process delivers m , then all correct processes eventually deliver m .

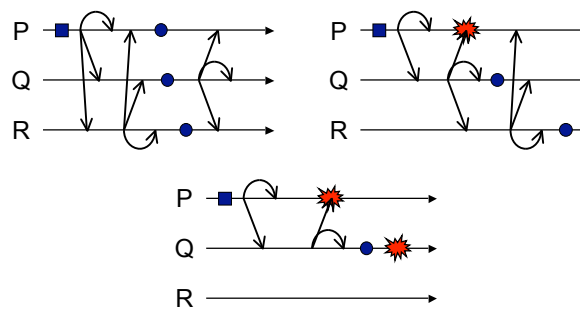
broadcast(UR, m) works as follows:
send(m) to all

upon receive(m)
if [for $f+1$ processes q : received(m) from q] then
deliver(UR, m)
if [didn't send m yet] then send(m) to all

© Fernando Pedone

Fault-tolerant broadcasts

- Uniform reliable broadcast (cont'd) - $f = 1$



© Fernando Pedone

Fault-tolerant broadcasts

- Cost of reliable broadcast algorithms

	Number of Messages	Latency
Reliable Broadcast (optimized)	n^2	δ
Uniform Reliable Broadcast	n^2	2δ

© Fernando Pedone

Fault-tolerant broadcasts

- Reducing the number of messages

- Asynchronous system + $\diamond S$

broadcast(R, m) works as follows:
send(m) to all

upon receive(m) for the first time do
deliver(R, m)
start Task Propagate(m , sender(m))

Task Propagate(m , q)
if suspect q then send(m) to all

© Fernando Pedone

Fault-tolerant broadcasts

■ FIFO reliable broadcast

– Transforms Reliable Broadcast into
FIFO Uniform Reliable Broadcast

– Assumptions

- If m is the i -th message broadcast by p , then m is tagged with $\text{sender}(m) = p$ and $\text{seq\#}(m) = i$
- Each process keeps a local vector of counters $\text{next}[1..n]$

© Fernando Pedone

Fault-tolerant broadcasts

Initialization

```
msgSet  $\leftarrow \emptyset$   
next[q]  $\leftarrow 1$ , for each process  $q$ 
```

To execute **broadcast**(F, m):
broadcast(R, m)

deliver($F, -$) occurs as follows:

```
upon deliver( $R, m$ ) do  
   $q \leftarrow \text{sender}(m)$   
   $\text{msgSet} \leftarrow \text{msgSet} \cup \{m\}$   
  while ( $\exists m' \in \text{msgSet}: \text{sender}(m') = q$  and  $\text{next}[q] = \text{seq\#}(m')$ )  
    deliver( $F, m'$ )  
     $\text{next}[q] \leftarrow \text{next}[q] + 1$ 
```

© Fernando Pedone

Fault-tolerant broadcasts

■ Causal reliable broadcast

– Transforms Uniform Reliable Broadcast into
Causal Reliable Broadcast

– Notation

- Message sequences: $\langle m_1, m_2, \dots \rangle$
- \perp is the empty sequence
- \oplus is the concatenation operator

© Fernando Pedone

Fault-tolerant broadcasts

Initialization

```
 $\text{rcntDlvs} \leftarrow \perp$ 
```

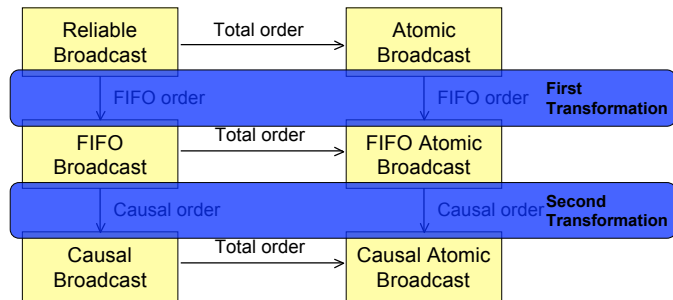
To execute **broadcast**(C, m):
broadcast($F, \text{rcntDlvs} \oplus \langle m \rangle$)
 $\text{rcntDlvs} \leftarrow \perp$

deliver($C, -$) occurs as follows:

```
upon deliver( $F, \langle m_1, m_2, \dots, m_l \rangle$ ) do  
  for  $i$  from 1 to  $l$  do  
    if  $p$  has not previously executed deliver( $C, m_i$ ) then  
      deliver( $C, m_i$ )  
       $\text{rcntDlvs} \leftarrow \text{rcntDlvs} \oplus \langle m_i \rangle$ 
```

© Fernando Pedone

Fault-tolerant broadcasts

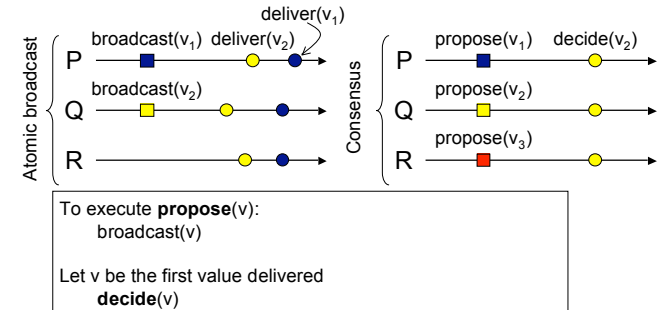


© Fernando Pedone

Fault-tolerant broadcasts

Atomic broadcast

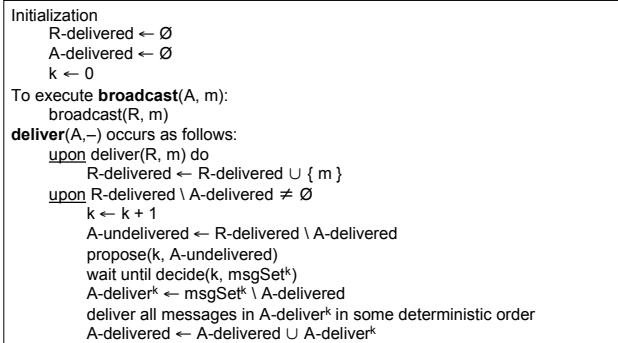
– From atomic broadcast to consensus



© Fernando Pedone

Fault-tolerant broadcasts

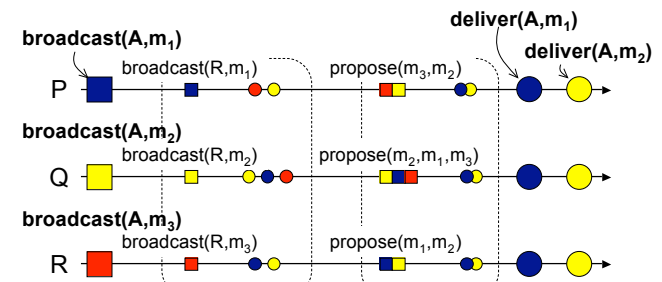
– From consensus to atomic broadcast



© Fernando Pedone

Fault-tolerant broadcasts

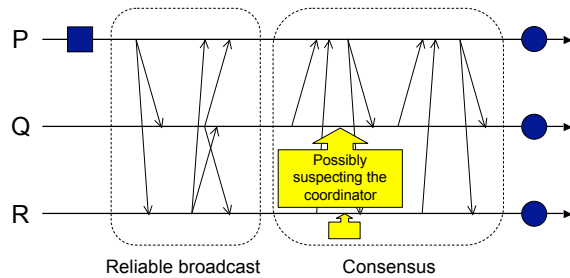
– From consensus to atomic broadcast (cont'd)



© Fernando Pedone

Fault-tolerant broadcasts

– From consensus to atomic broadcast (cont'd)

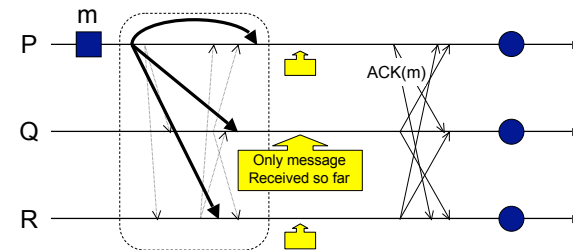


© Fernando Pedone

Fault-tolerant broadcasts

■ Generic broadcast

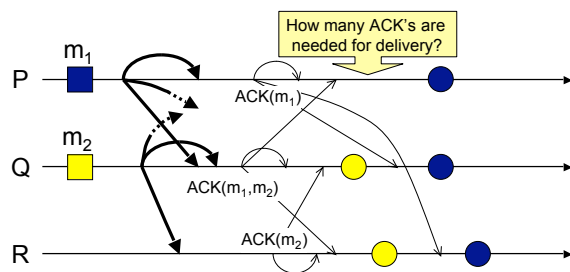
– Single-message case



© Fernando Pedone

Fault-tolerant broadcasts

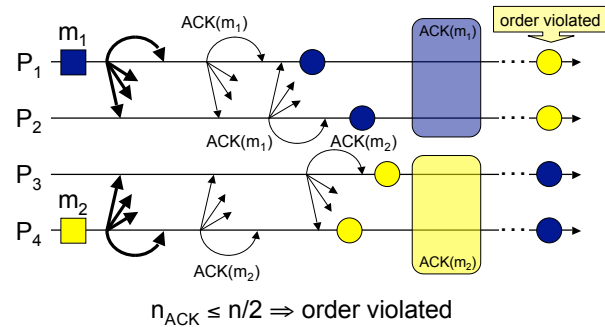
– Two non-conflicting messages



© Fernando Pedone

Fault-tolerant broadcasts

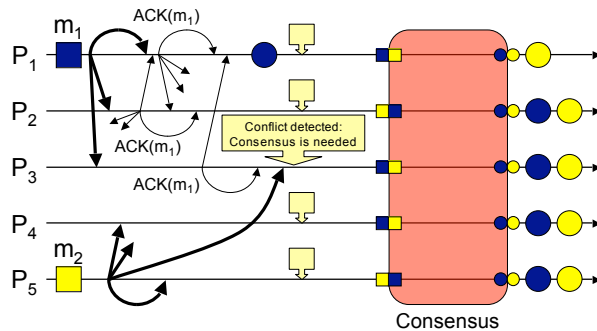
– The general case (m_1 and m_2 conflict)



© Fernando Pedone

Fault-tolerant broadcasts

– The general case (m_1 and m_2 conflict)



© Fernando Pedone

Fault-tolerant broadcasts

– The general case (m_1 and m_2 conflict)

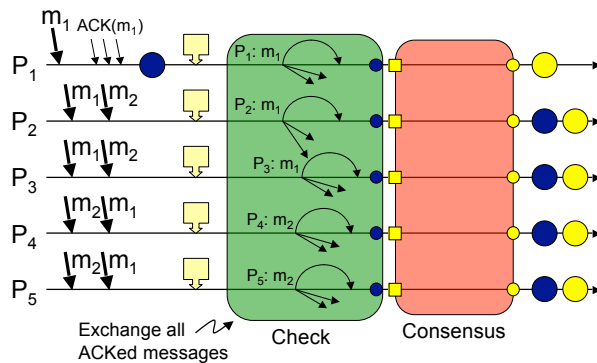
- Conflicts should be detected: $n_{ACK} > n / 2$
- For conflicting messages:
“If m_1 and m_2 conflict, and m_1 (or m_2) has been delivered before consensus, then consensus cannot contradict this order.”

How to ensure this???

© Fernando Pedone

Fault-tolerant broadcasts

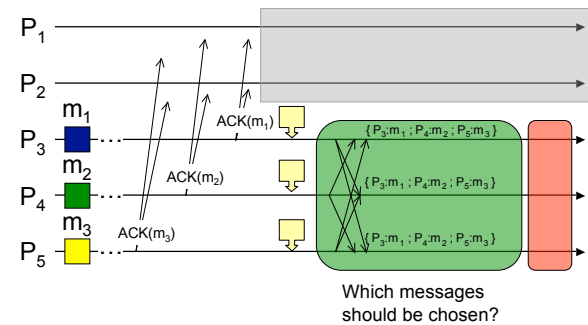
– The general case (m_1 and m_2 conflict)



© Fernando Pedone

Fault-tolerant broadcasts

– The general case (m_1 and m_2 conflict)



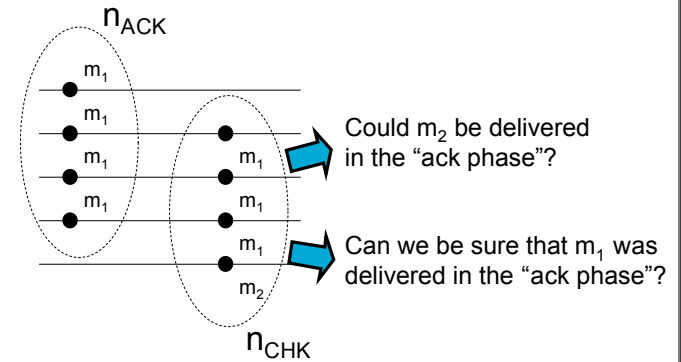
© Fernando Pedone

Fault-tolerant broadcasts

- If $n_{ACK} = \lceil (n+1) / 2 \rceil$ (i.e., majority), how much should n_{CHK} be?
- How to choose a message in the "check phase"?
- Problem: we can't tolerate failures
- What can be done?
- Increase n_{ACK} ... for example, if $n_{ACK} = 4$ and $n = 5$, could we have $n_{CHK} = 4$?

© Fernando Pedone

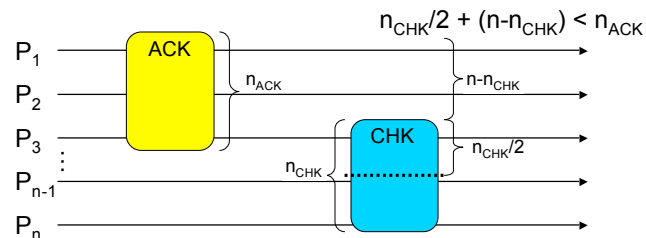
Fault-tolerant broadcasts



© Fernando Pedone

Fault-tolerant broadcasts

- What is the relation between n_{ACK} and n_{CHK} ?



© Fernando Pedone

Fault-tolerant broadcasts

- The general case (cont'd)

$$n_{CHK}/2 + (n - n_{CHK}) < n_{ACK}$$

$$n_{CHK} + 2n - 2n_{CHK} < 2n_{ACK}$$

$$-n_{CHK} + 2n < 2n_{ACK}$$

$$2n_{ACK} + n_{CHK} > 2n$$

Best case: $n_{ACK} = n_{CHK} = x$

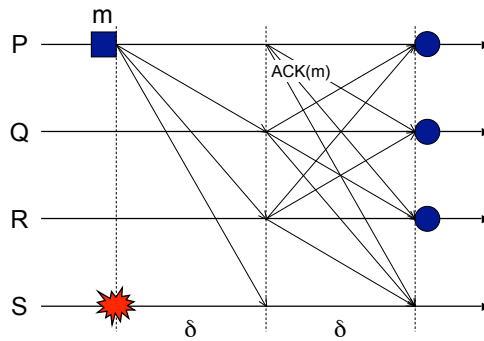
$$2x + x > 2n$$

$$x > 2n / 3 \text{ (Optimal for ACK with } 2\delta)$$

© Fernando Pedone

Fault-tolerant broadcasts

– Optimality (i.e., lower bound)



© Fernando Pedone

Fault-tolerant broadcasts

```

Initialization
R-delivered ← G-delivered ← ∅
pending1 ← g-Deliver1 ← ∅
k ← 1

To execute broadcast(m):
broadcast(R, m)

when deliver(R, m) do
R-delivered ← R-delivered ∪ { m }

when (R-delivered \ (G-delivered ∪ pendingk) ≠ ∅)
if [ for all m, m' ∈ (R-delivered \ G-delivered): m doesn't conflict with m' ] then
pendingk ← R-delivered \ G-delivered
send (k, pendingk, ACK) to all
else
Handle conflict (next slide)

when receive (k, pendingk, ACK) from pi
while ∃m such that [ for nACK pj: received (k, pendingk, ACK) from pj and
m ∈ (pendingk \ g-Deliverk) ] do
g-Deliverk ← g-Deliverk ∪ { m }
deliver(m)
    
```

© Fernando Pedone

Fault-tolerant broadcasts

```

Handling conflicts:

when (R-delivered \ (G-delivered ∪ pendingk) ≠ ∅)
if [ for all m, m' ∈ (R-delivered \ G-delivered): m doesn't conflict with m' ] then
...
else
send (k, pendingk, CHK) to all
wait until [ for nCHK pj: received (k, pendingk, CHK) from pj ]
msgSet = { m | for [(nCHK+1)/2] pj: received (k, { ..., m, ... }, CHK) from pj }

propose(k, msgSet, (R-delivered \ (G-delivered ∪ pendingk)))
wait until decide (k, NCset, Cset)

for each m ∈ NCset \ (G-delivered ∪ g-Deliverk) do deliver(m)
in ID order: for each m ∈ Cset \ (G-delivered ∪ g-Deliverk) do deliver(m)
G-delivered ← G-delivered ∪ NCset ∪ Cset

k ← k + 1
pendingk ← g-Deliverk ← ∅
    
```

© Fernando Pedone

Paxos

- The part-time parliament
 - Leslie Lamport
 - Consensus algorithm
 - SRC Research Report 49, Sep. 1, 1989
 - ACM TOCS, May 1998
 - Submitted in 1990

© Fernando Pedone

Paxos

The protocol

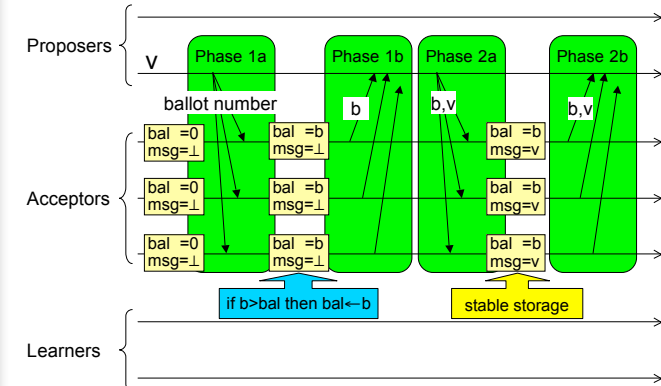
– System model

- $\Pi = \{ p_1, p_2, \dots, p_n \}$
- Asynchronous system, plus...
- Leader-election oracle (Ω)
- Crash-recovery failure model
- Unreliable channels
- Stable storage

© Fernando Pedone

Paxos

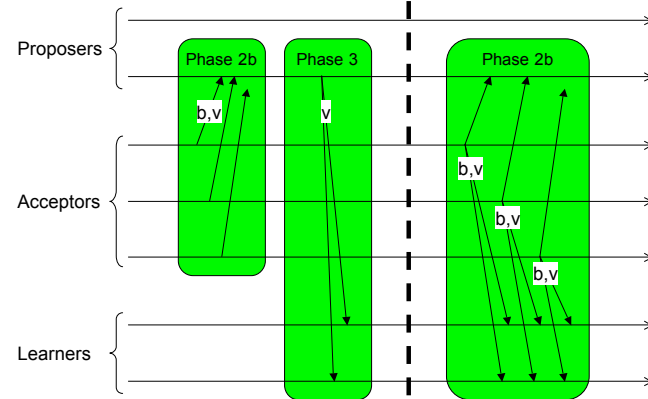
Choosing a value



© Fernando Pedone

Paxos

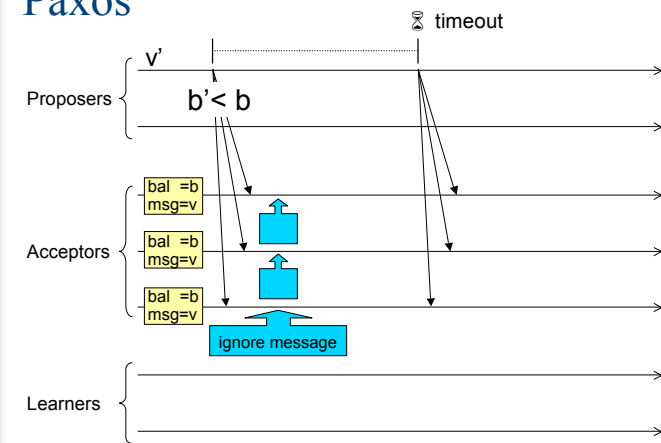
Learning a chosen value



© Fernando Pedone

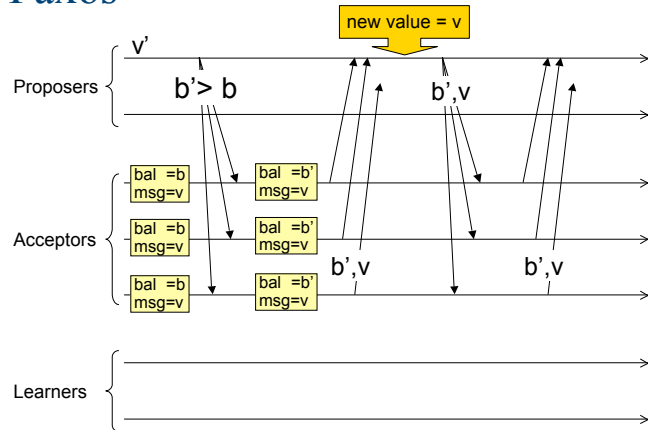
Paxos

Handling old ballot numbers



© Fernando Pedone

Paxos

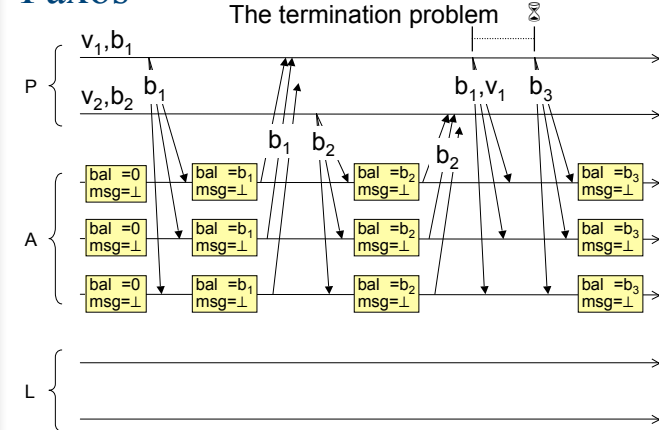


Handling new ballot numbers

© Fernando Pedone

Paxos

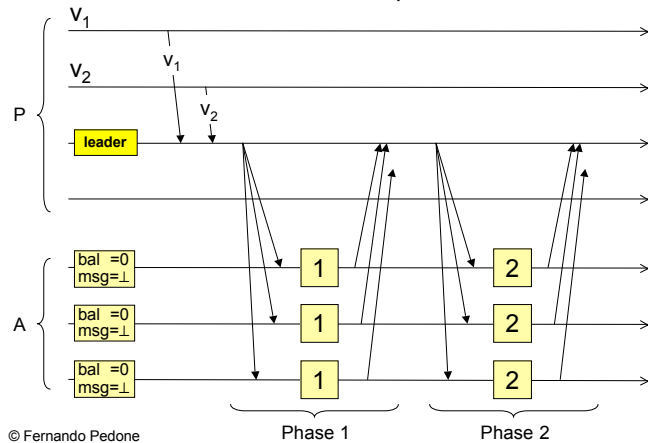
The termination problem



© Fernando Pedone

Paxos

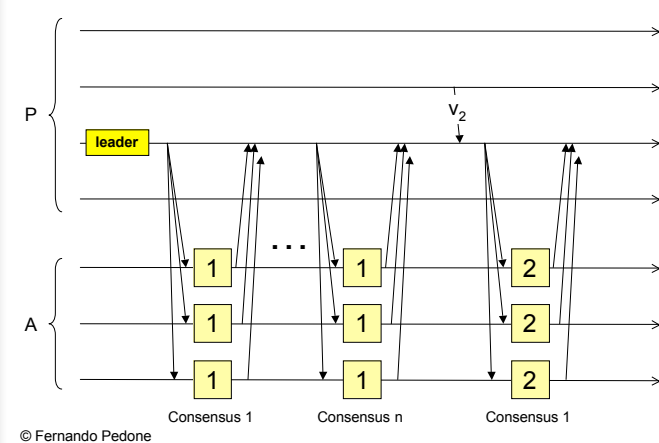
A leader-based protocol



© Fernando Pedone

Paxos

Ballot reservation



© Fernando Pedone