

Coordination

Part-2



**Distributed
Mutual Exclusion**



**Election
Mechanisms**

Maekawa Voting

Filosophy

- Drawback of Ricart-Agrawala : all processes need to agree —> linear scaling
- Maekawa voting : only SUBSET of processes involved
- Basic idea : “vote on behalf of others”
- Candidate process must collect sufficient votes before entering

Voting set for each process p_i : V_i

- $V_i \subseteq \{p_1, \dots, p_N\}$ (Required)
- $p_i \in V_i$ (Required)
- $V_i \cap V_j \neq \emptyset$ (Required)
- $|V_i| = K$ (fairness) (Bonus)
- Every p_j contained in exactly M voting sets (Bonus)

Basic idea

- process $q \in (V_i \cap V_j)$
- —> q makes sure p_i and p_j are not simultaneously executing critical section
- —> safety condition met
- q only votes for one process
- Additional state needed per process : **voted**

The algorithm

Process i requesting access

Initialization

state = Released
voted = False

enter()

state = Wanted
multicast **Request(pi)** to voting set **Vi**
wait until **K** Reply messages received
state = Held

leave()

state = Released
multicast **Release** to voting set **Vi**

Member of Vi receiving request

when receiving Request(pj) at pi

```
if (state == Held) or (voted == True) {  
    enqueue Request in Q  
    // do NOT reply  
} else {  
    send Reply to pj  
    voted = True  
}
```

when receiving Release(pj) at pi

```
if (Q != empty) {  
    dequeue pending Request m from Q  
    send Reply to sender(m)  
    voted = True  
} else {  
    voted = False  
}
```

Differences w.r.t. Ricart-Agrawala

- additional state variable per process needed (voted or not)
- multicast request to enter to voting set only
- explicit leave needed (so voting processes can vote for other process)

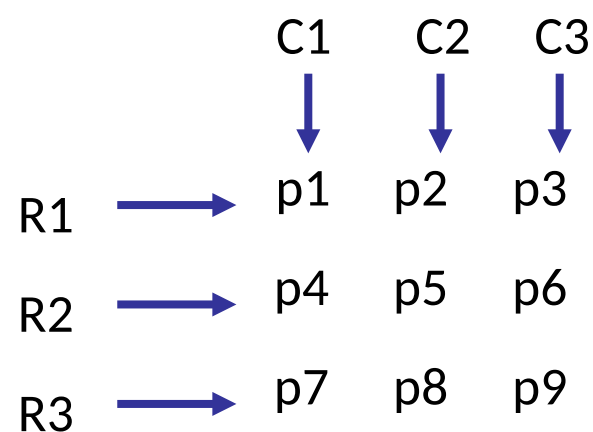
Constructing voting sets

Choosing parameters

- Theoretical result : optimal solution (minimal K)
 - $K \approx \sqrt{N}$
 - $M = K$
- In practice : difficult to calculate optimal V_i
Sub-optimal solution
 - $K \approx 2\sqrt{N}$
 - $M = K$

Practical algorithm (for $N=S^2$)

for $S=3$



V1 = R1 U C1
V2 = R1 U C2
V3 = R1 U C3
V4 = R2 U C1
V5 = R2 U C2
V6 = R2 U C3
V7 = R3 U C1
V8 = R3 U C2
V9 = R3 U C3

Theory:



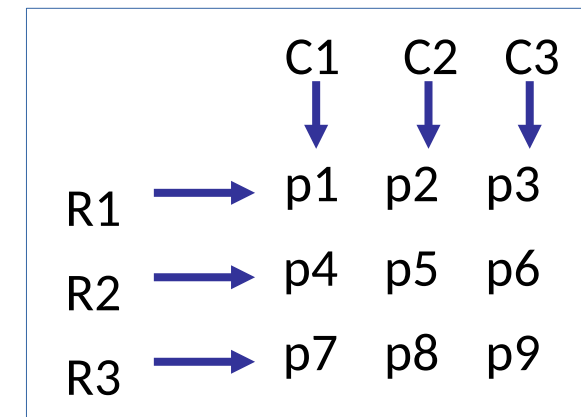
Constructing voting sets

In general

Construct $S \times S$ matrix A , consisting of all processes

- i is row where p is found
- j is column where p is found
- $R_i = i$ -th row of A
- $C_j = j$ -th column of A

Voting set $V_p = R_i \cup C_j$



Checking this V_p ...

$$p \in R_i, p \in C_j \Rightarrow p \in V_p$$

$$V_p = R_i \cup C_j$$

$$V_q = R_s \cup C_t$$

$$\begin{aligned} \Rightarrow V_p \cap V_q &= (R_i \cup C_j) \cap (R_s \cup C_t) \\ &= (R_i \cap R_s) \cup (R_i \cap C_t) \cup (C_j \cap R_s) \cup (C_j \cap C_t) \end{aligned}$$

$$\text{BUT } R_k \cap C_l \neq \Phi \text{ (for any } k, l) \text{ (} s_{kl} \in R_k \text{ and } s_{kl} \in C_l \text{)}$$

$$\Rightarrow V_p \cap V_q \neq \Phi$$

$$K = 2S - 1$$

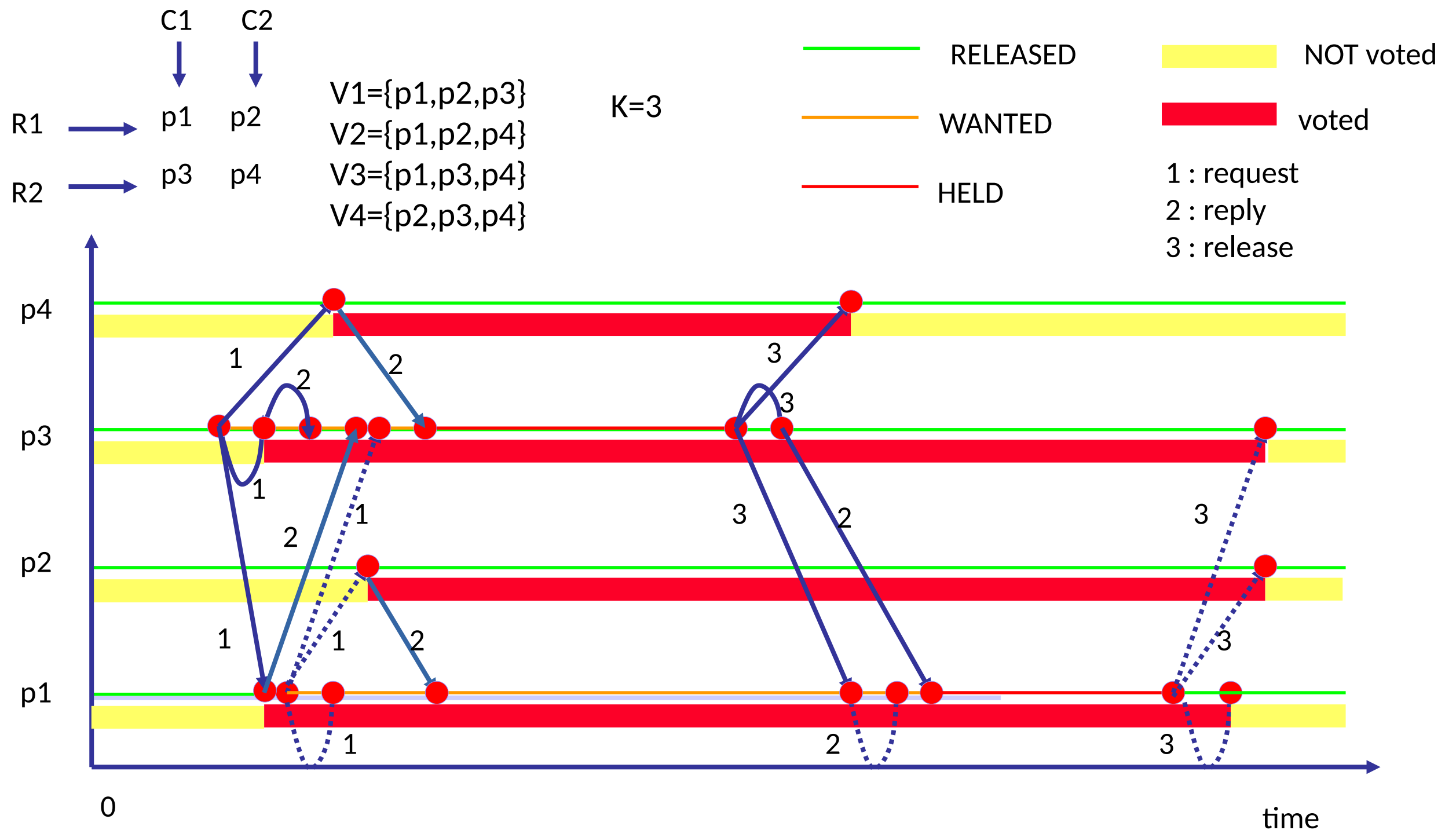
$$K = M$$

Voting set for each process $p_i : V_i$

- $V_i \subseteq \{p_1, \dots, p_N\}$, satisfying
- $p_i \in V_i$
- $V_i \cap V_j \neq \emptyset$
- $|V_i| = K$ (fairness)
- p_j contained in M voting sets

Example

Processes p1 and p3 request entry, p2 and p4 not



Algorithm OK ?

1

Safety

Suppose **p** and **q** simultaneously active

=> all processes in **V_p** and **V_q** voted for **p** AND **q**

=> process **t** in **V_p ∩ V_q ≠ ∅** voted for **p** AND **q**

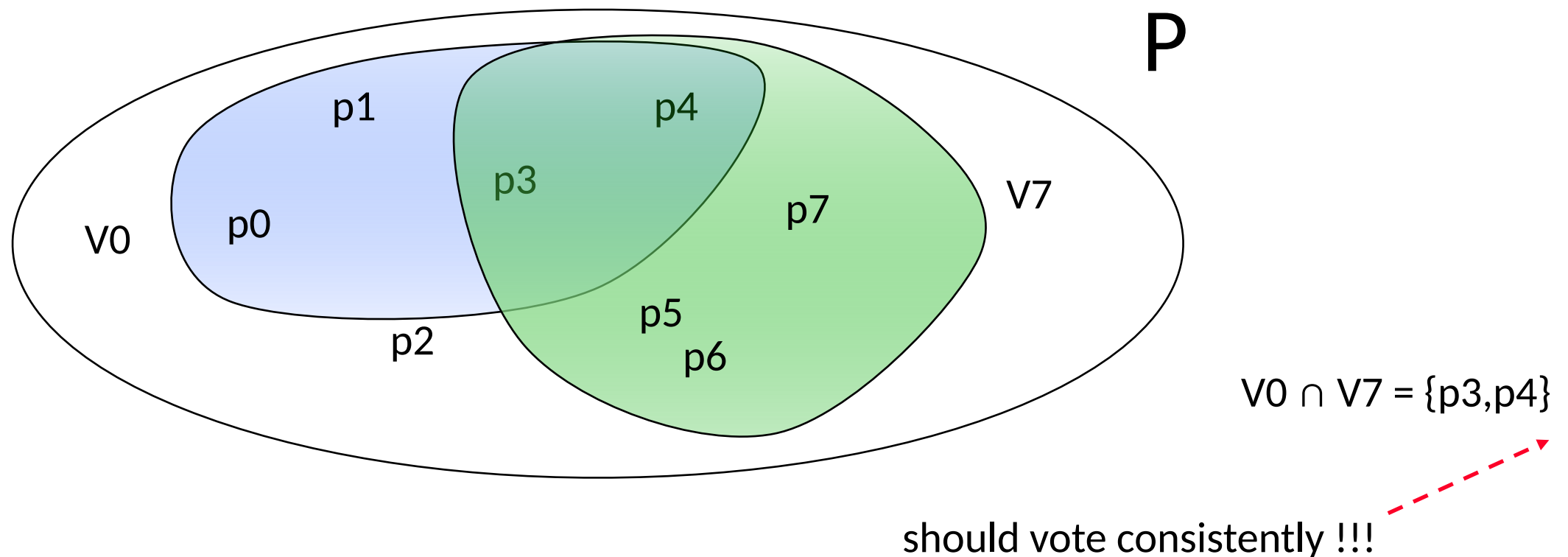
impossible :

voted=True immediately after voting for 1 process

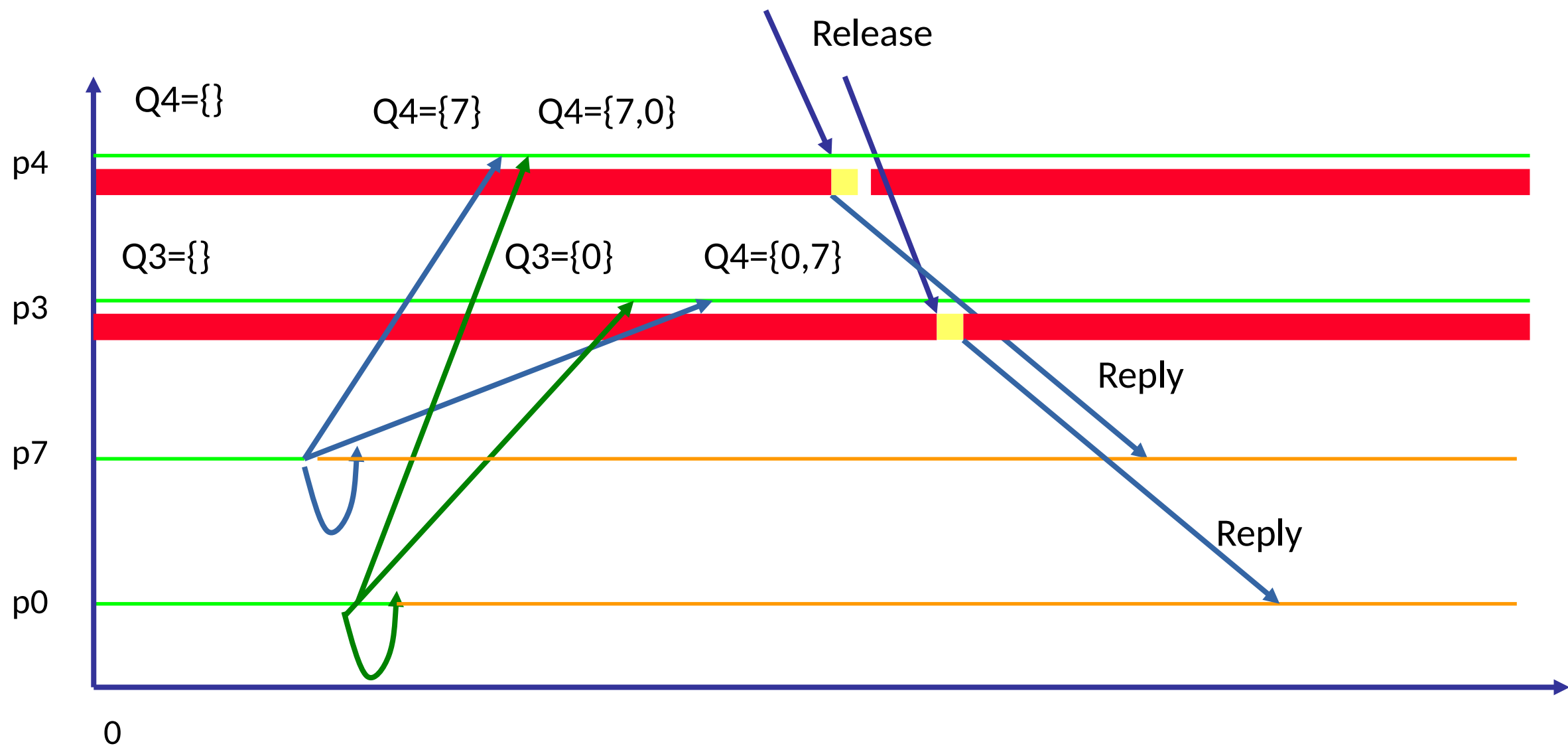
2

Liveness

deadlock prone !!!



Algorithm NOT OK !



p7 and p0 will NEVER
receive Reply from complete
voting set !

Algorithm OK ?

①

Safety

Suppose p and q simultaneously active
=> all processes in V_p and V_q voted for p AND q
=> process t in $V_p \cap V_q \neq \emptyset$ voted for p AND q
impossible :
voted=True immediately after voting for 1 process



②

Liveness

deadlock prone !!!
Solution: use totally ordered Lamport clocks



③

Fairness

Guaranteed if we use Lamport clocks



Algorithm efficient ?

Bandwidth usage

enter()

same as Ricart-Agrawala, but message sent to voting set only !

-> **K Request** messages

-> **K Reply** messages

leave()

explicit **Leave** message now needed

-> **K Release** messages

Client delay

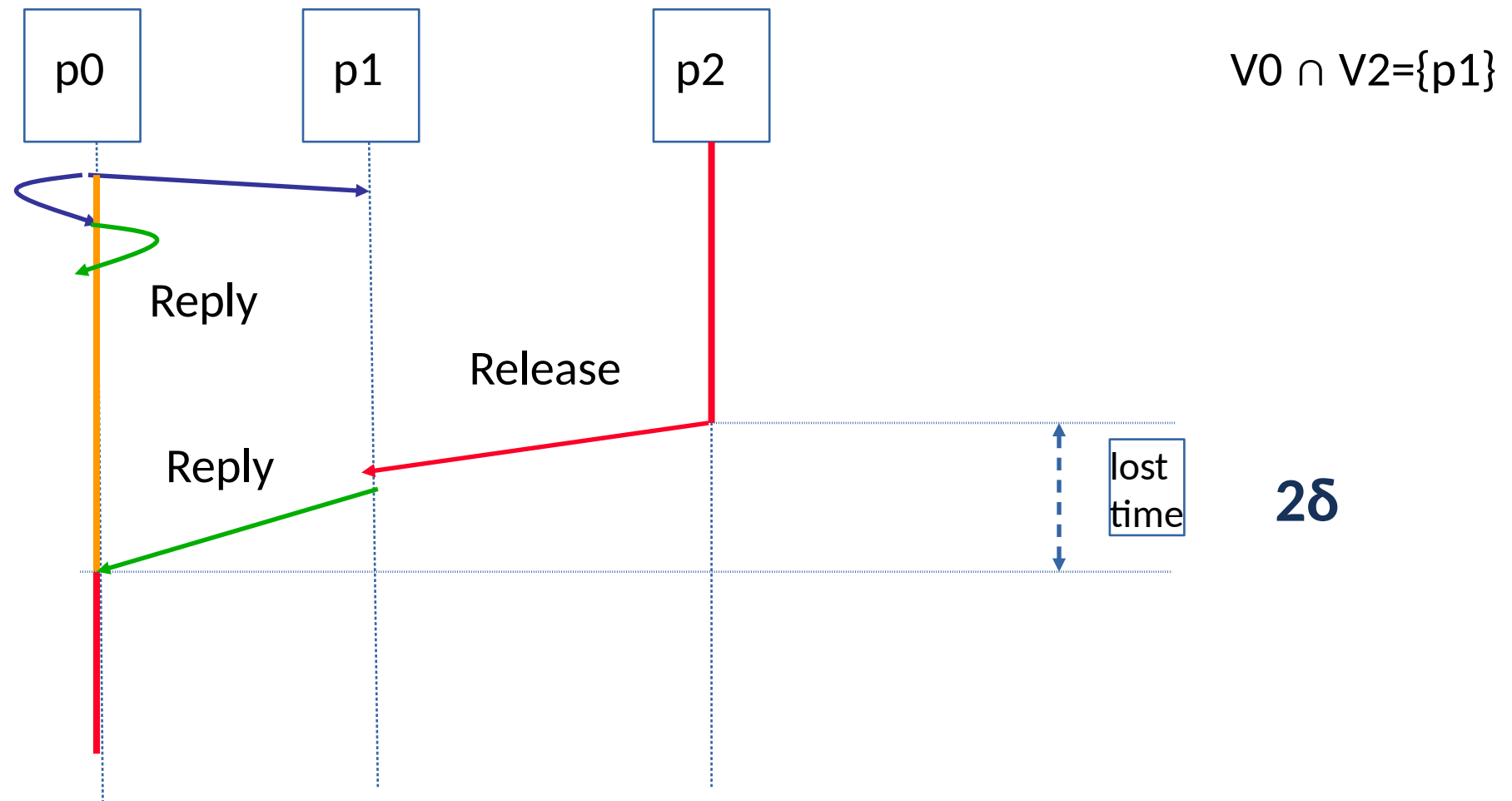
same as Ricart-Agrawala

enter() : 2δ

leave() : 0δ

Algorithm efficient ?

Synchronization delay



Summary on efficiency

	Bandwidth		Client delay		Synchronization delay
	enter()	leave()	enter()	leave()	
Central server	$2M$	$1M$	2δ	0δ	2δ
Ring algorithm	constant $(N+1)/2$		$N\delta/2$		$(N+1)\delta/2$
Ricart-Agrawala	$(2N-1)M$	$0M$	2δ	0δ	1δ
Maekawa voting	$2KM$	KM	2δ	0δ	2δ

N : number of processes
 M : number of voting sets that each process belongs to
 δ : cost of sending a message



**Distributed
Mutual Exclusion**

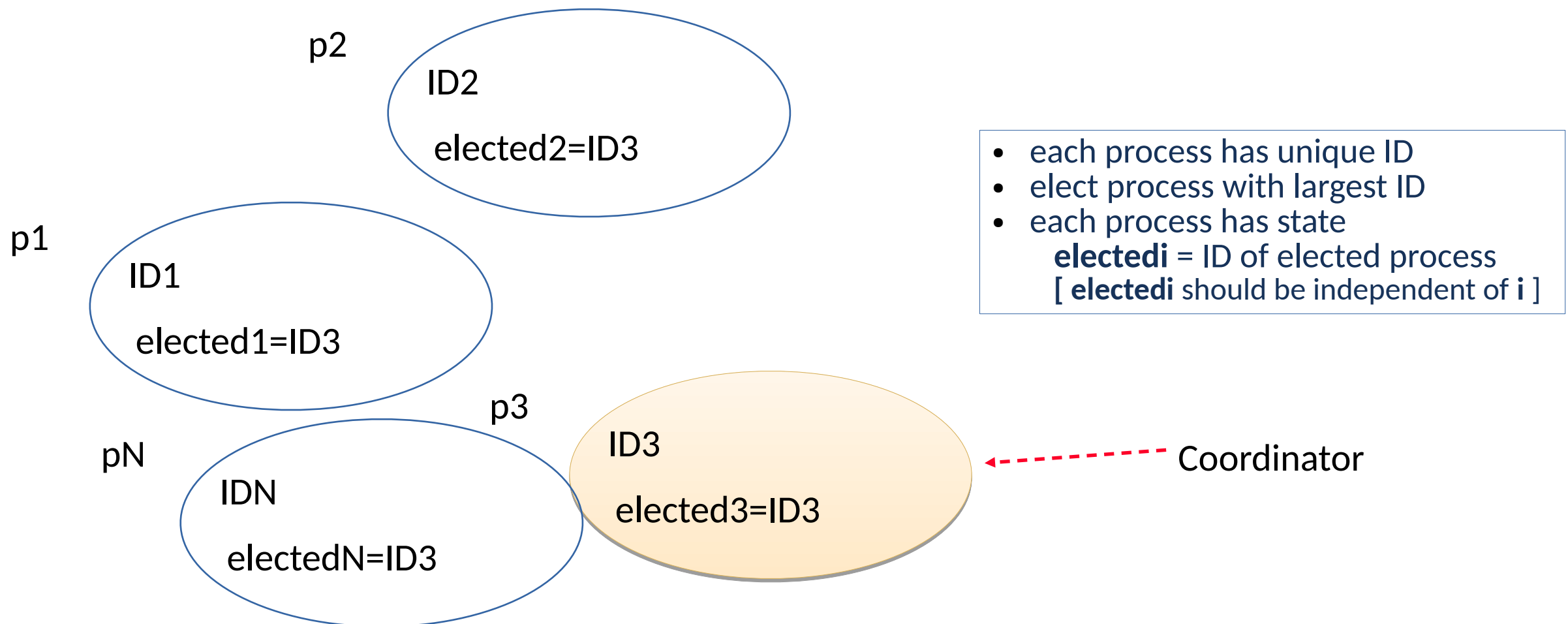


**Election
Mechanisms**

The Election Problem

Consider

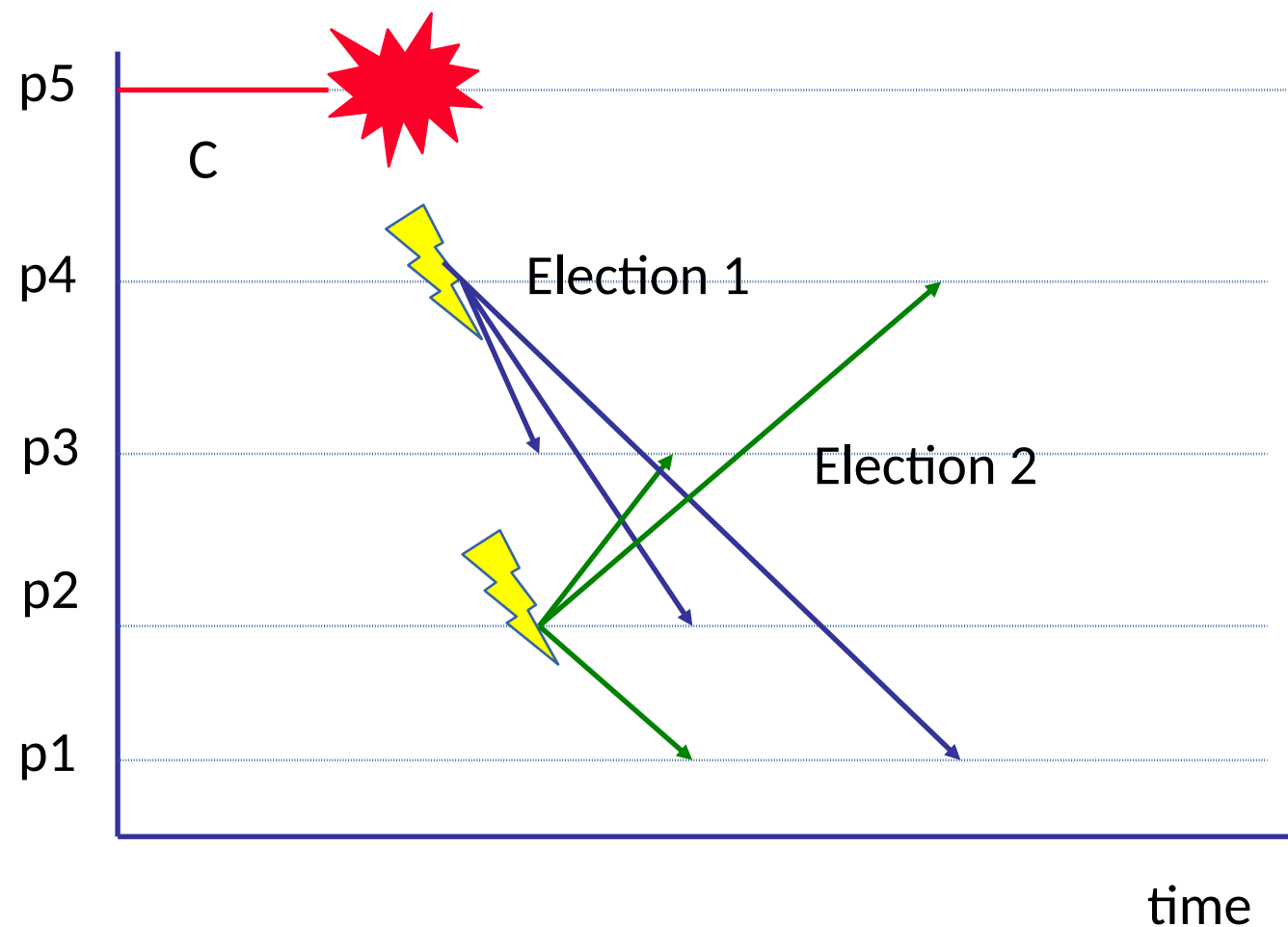
- N processes {p1, ... , pN}
 - NO shared variables
 - knowing each other (can communicate)
- select ONE process to play special role (e.g. coordinator)
- every process **pi** should have same coordinator
- if elected process fails : do new election round



Some terms

Assumed environment

- request election : a process “calls the election”
 - at most 1 election initiated per process
 - possibly **N** elections running simultaneously
- at any time, a process is either
 - *participant* : currently engaged in some election
 - *non-participant* : currently not engaged in any election



Good elections

Correctness requirements

①

Safety (REQUIRED)

each participant process p_i has :

elected $_i$ = ?

OR **elected $_i$** = P

(P is elected process, non-crashed with largest ID)

②

Liveness (REQUIRED)

all processes p_i participate

eventually set **elected $_i$** \neq ? or crash



Evaluation metrics

- Bandwidth

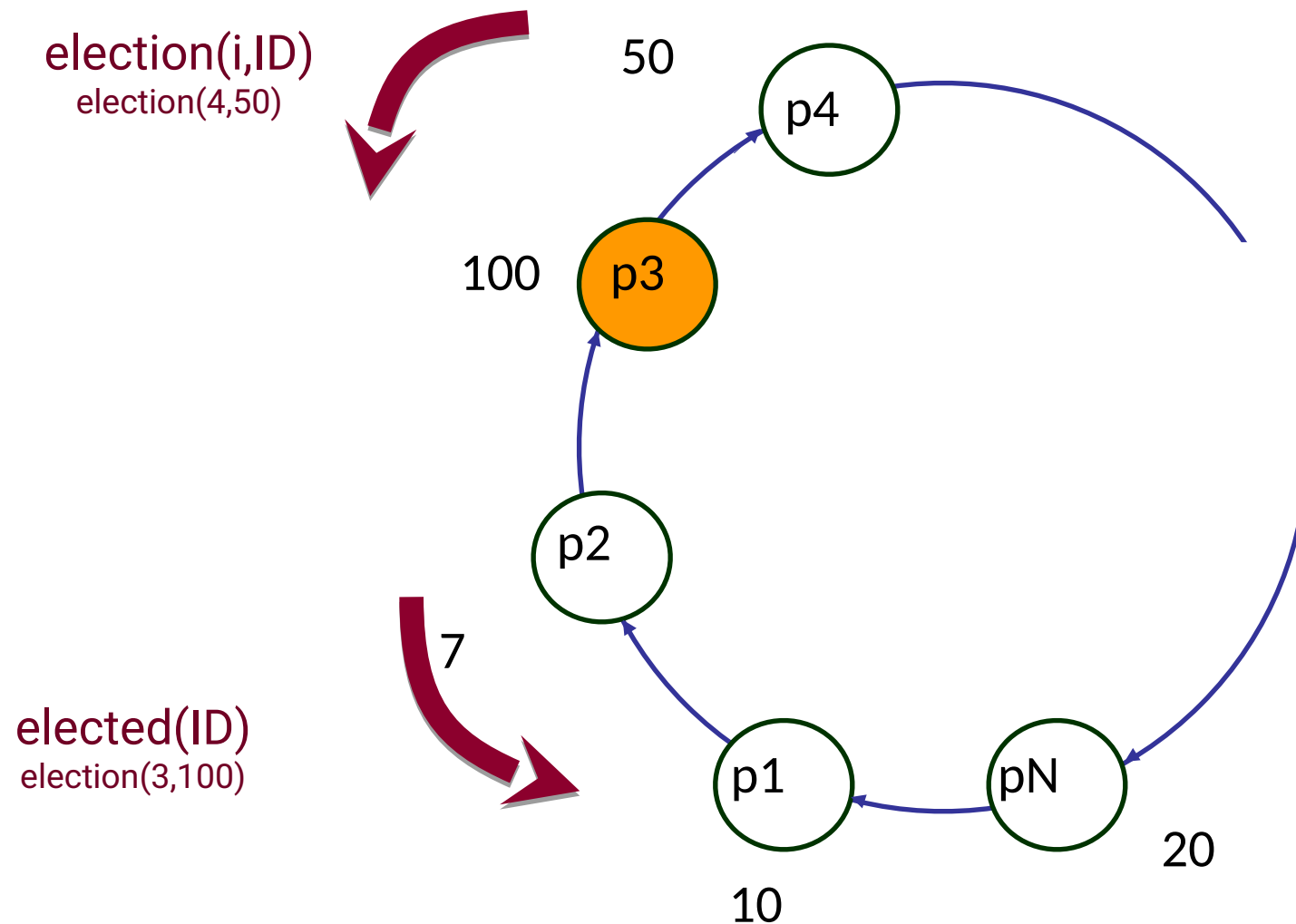
messages needed to do election process

- turnaround time

time needed for election round

Ring algorithm: Chang – Roberts

- processes organized in ring
- non-identical IDs (how to make IDs unique ?)
- processes know how to communicate



Failure modes

No failures :

- reliable channels
- no process crashes

Asynchronous system

Algorithm

Each process p

Initialization

participant_i = FALSE for all i

Start election process p_i

participant_i = TRUE

send message Election(i, ID_i)

Receipt of Elected(i)-message at p_j

if($i \neq j$) {

 participant_j = FALSE

 elected_j = i

 forward Elected(i)

}

Receipt of Election(i, ID)-message at p_j

if($ID > ID_j$) {

 forward Election(i, ID)

 participant_j = TRUE

}

if($(ID \leq ID_j) \text{ and } (i \neq j)$) {

 if(participant_j == FALSE) {

 send Election(j, ID_j)

 participant_j = TRUE

 }

}

if($i == j$) {

 participant_j = FALSE

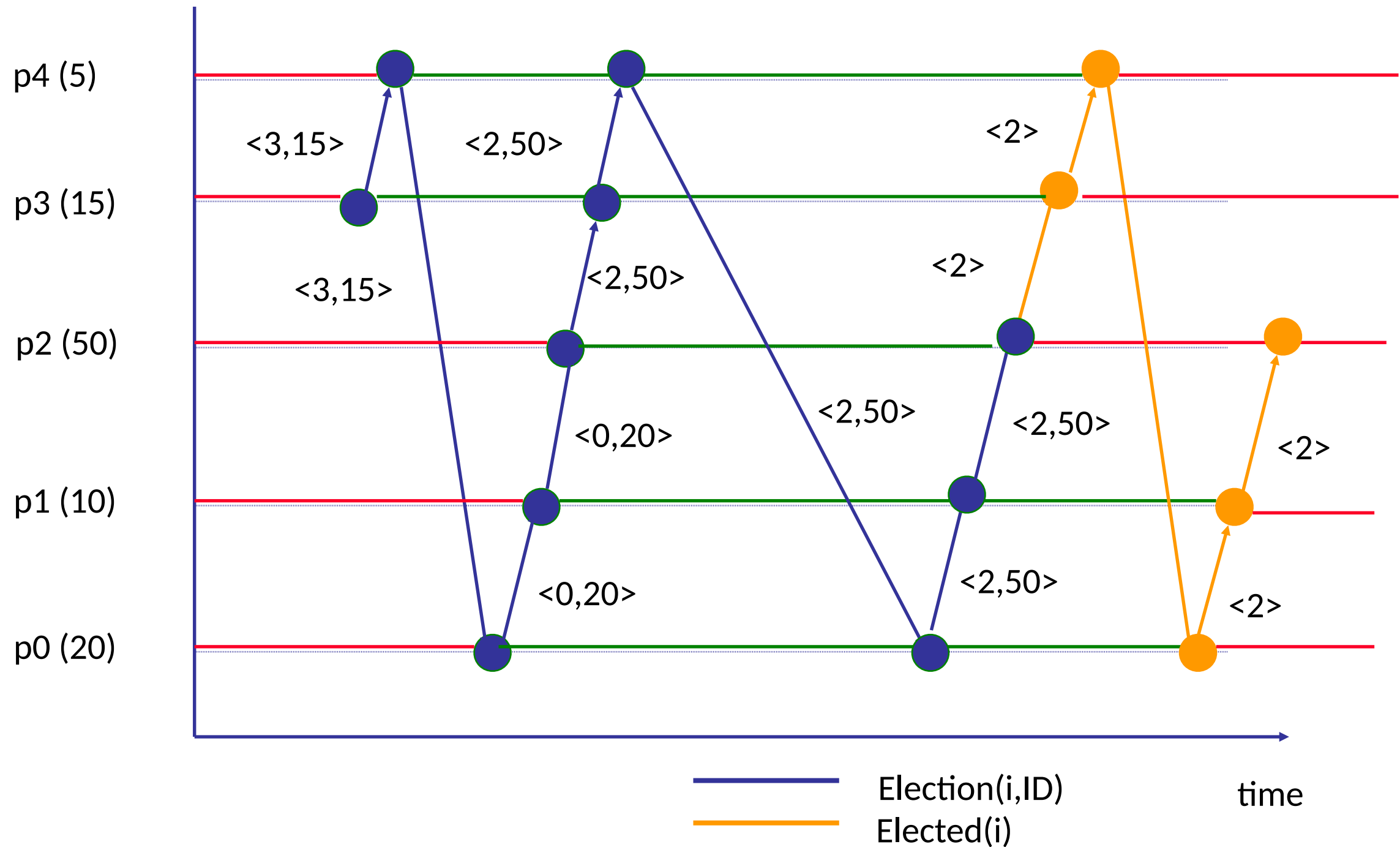
 elected_j = j

 send Elected(j)

}

Example

p3 calls the election



Algorithm OK ?

1

Safety

Elected message only sent if Election-message with own ID received

Suppose **p** and **q** both elected

=> **p** received **Elected(p)**

q received **Elected(q)**

BUT ID's are unique

$(ID_p < ID_q) \Rightarrow q$ will NOT forward **Elected(p, ID_p)**

$(ID_p > ID_q) \Rightarrow p$ will NOT forward **Elected(q, ID_q)**

=> impossible for BOTH messages to visit complete ring

=> impossible **p** AND **q** to be elected

2

Liveness

No failures

=> messages allowed to circulate

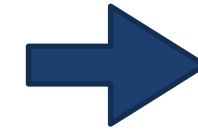
=> circulation stops (through participant state variable)

Algorithm efficient ?

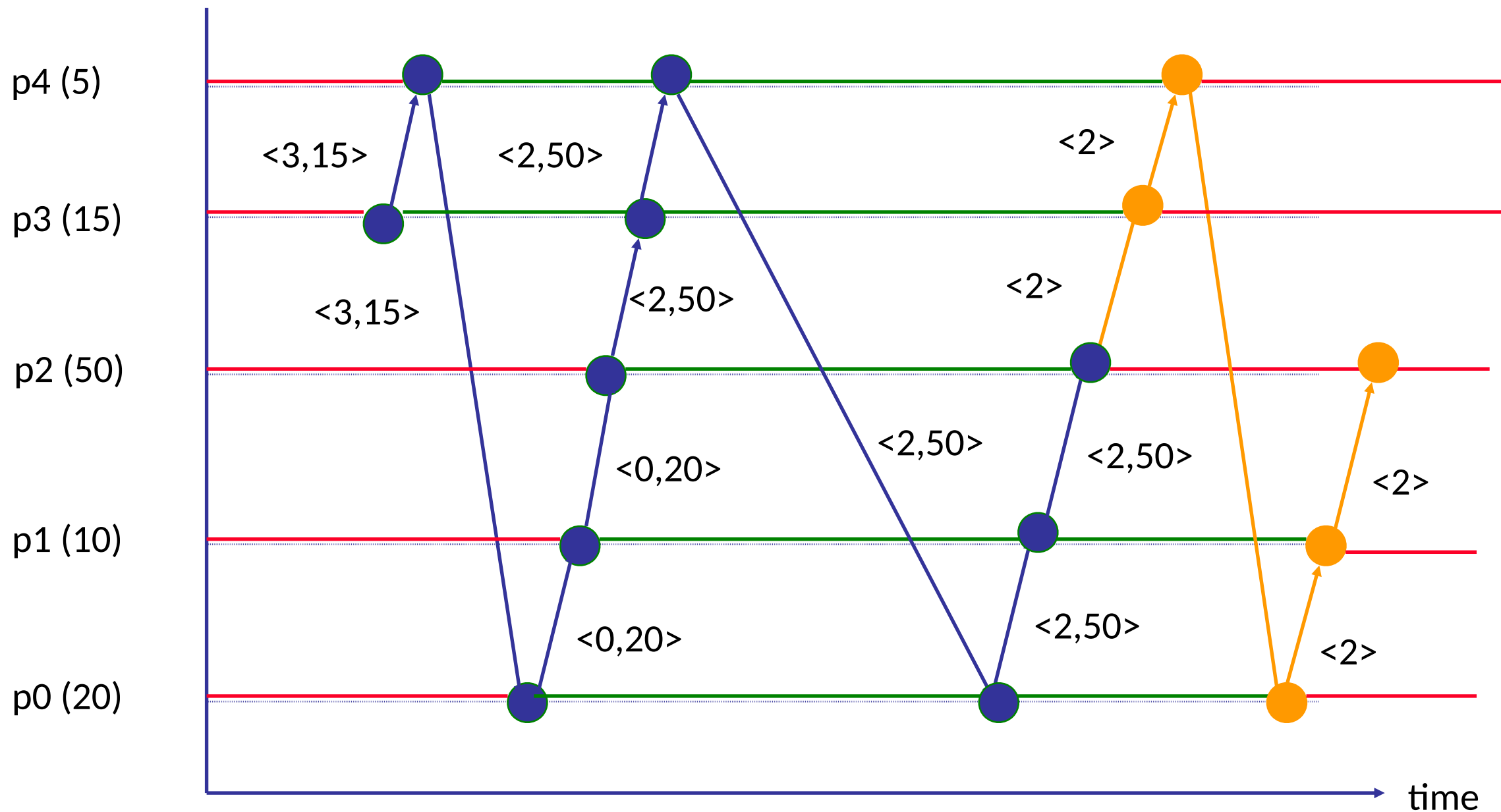
3 phases

1. ID in Election message grows
2. do complete round with constant ID
3. let the Elected message circulate

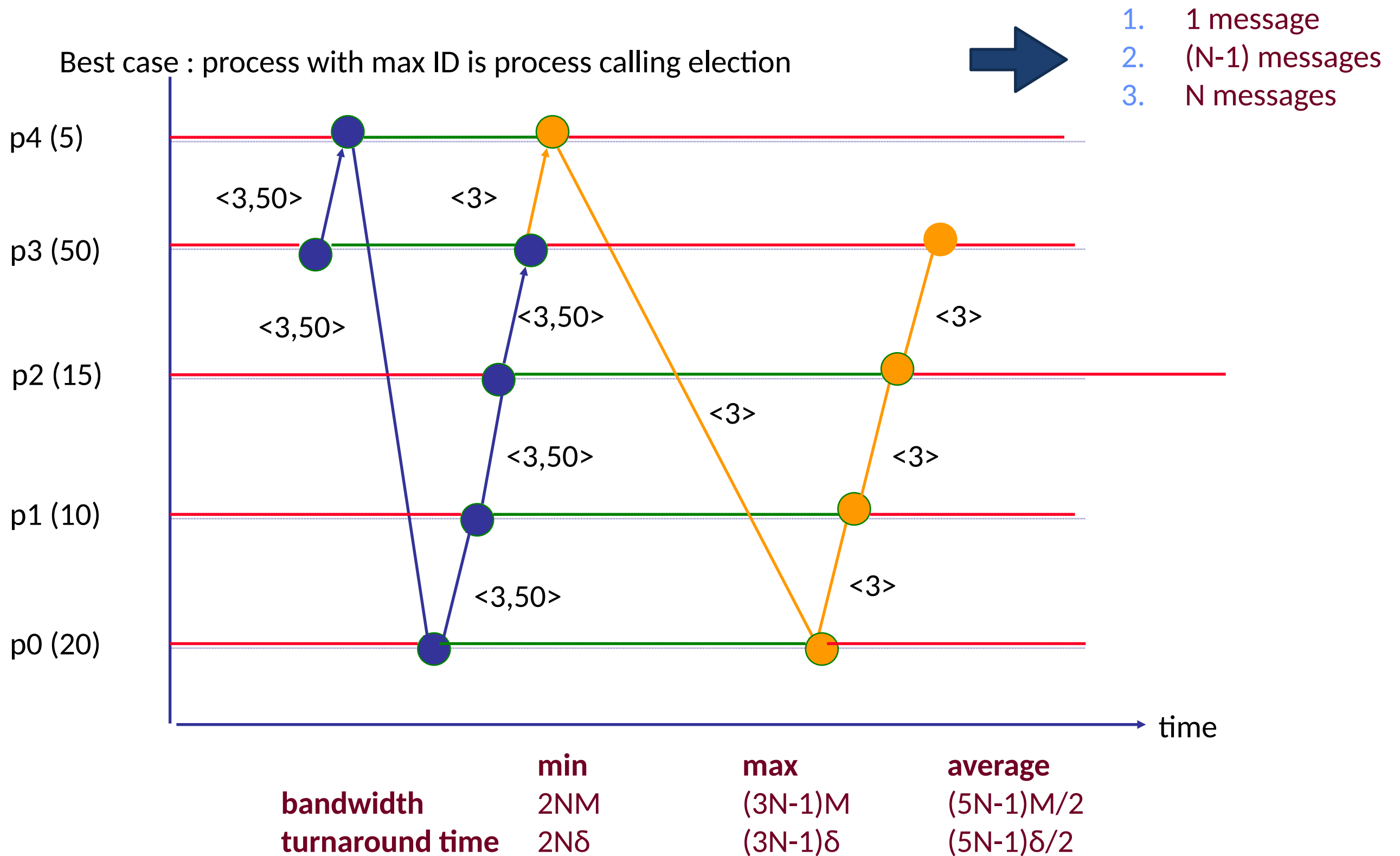
Worst case : process with max ID is last process visited



1. N messages
2. (N-1) messages
3. N messages



Algorithm efficient ?



Bully algorithm (Garcia – Molina)

Context

- Failure mode: process crashes dealt with
- System model: Synchronous system (uses time-outs to detect failure)
- A-priori knowledge: process knows all processes with larger ID

Filosophy

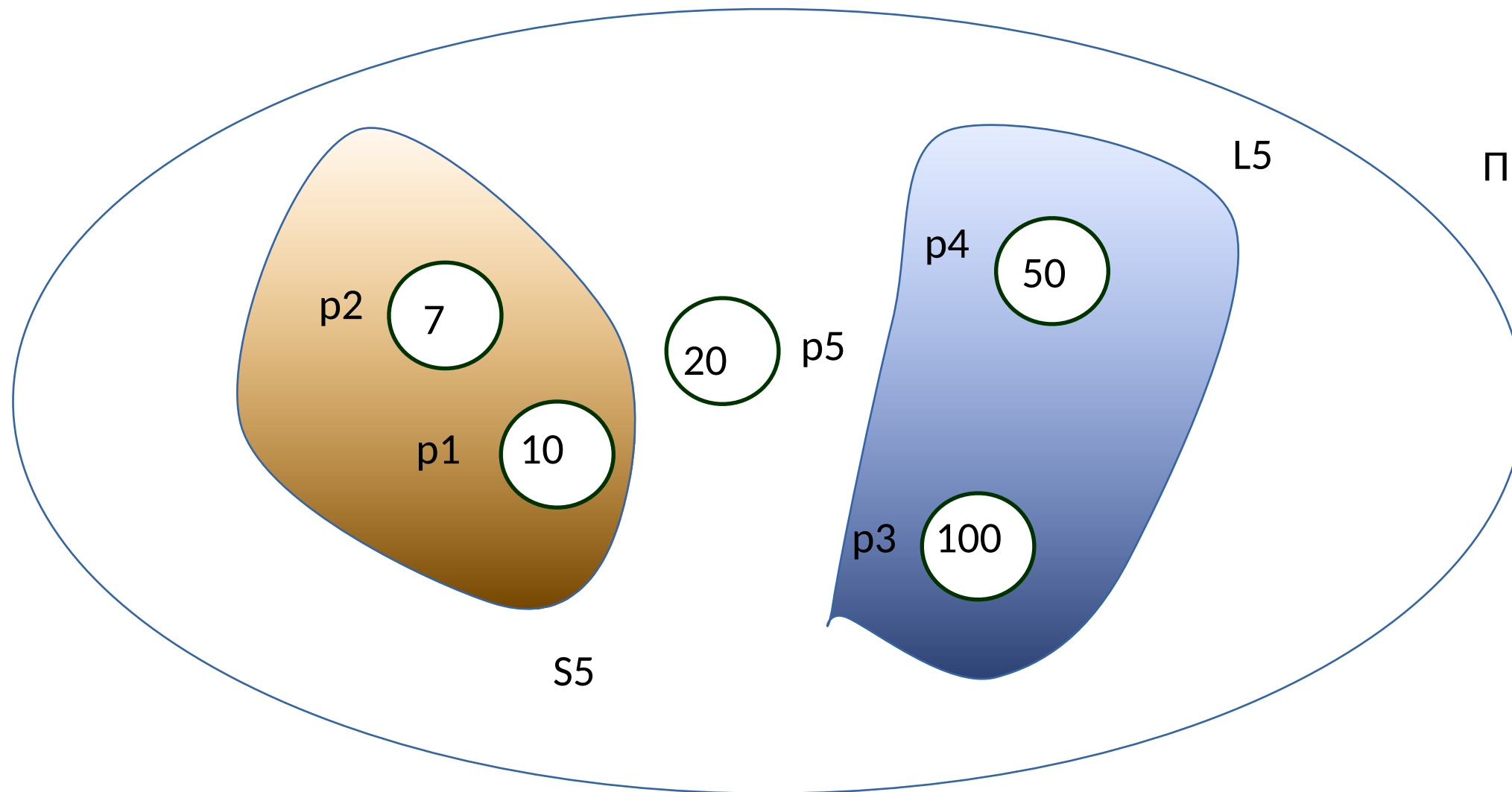
- Election starts when current coordinator fails
- Failure discovery :
 - by timeouts
 - election possibly by several processes
- Each process has
 - set **L** of candidate coordinators (set of processes with larger ID)
 - set **S** of other processes (smaller IDs)
- Upper bound for answering : **T**



Communication

Messages involved

- **election** announce election round
- **answer** reply to election message
- **coordinator** announce ID of elected coordinator



Algorithm

Call the election process p_i

```
if { $L == \{ \}$ } then {  
    elected =  $i$   
    send coordinator( $i$ ) to  $S$   
} else {  
    send election( $i$ ) to all processes in  $L$   
    if no answer-message in period  $T$  then {  
        elected =  $i$ ;  
        send coordinator( $i$ ) to  $S$   
    } else {  
        if no coordinator-message in  $T'$  then call election again  
    }  
}
```

Receipt **coordinator(j)** at p_i

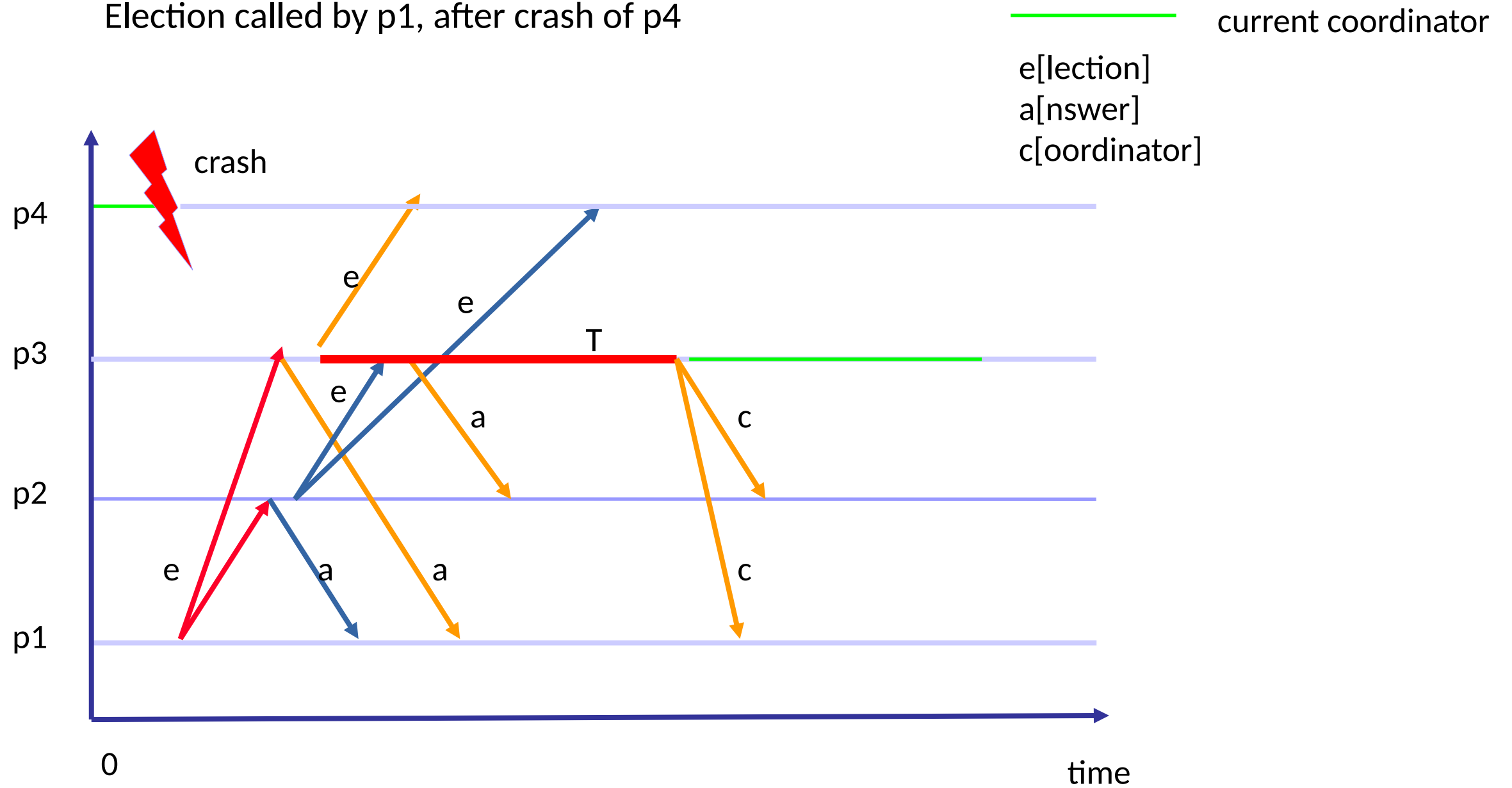
elected _{i} = j

Receipt of **election(j)-message** at p_i

```
if (no elections initiated by  $p_i$ ) {  
    send answer-message to  $p_j$   
     $p_i$  calls election  
}
```

Example

Election called by p1, after crash of p4



Algorithm OK ?

1

Safety

- Ok, if no process replacement
- If process replacement occurs & new process has highest ID
 - Duplicate coordinator messages
 - One from the replacing process
 - One from the largest-but-one ID

2

Liveness

messages delivered reliably (no communication faults) either

- answer from **L**
- process is coordinator itself
- in any case coordinator identified !

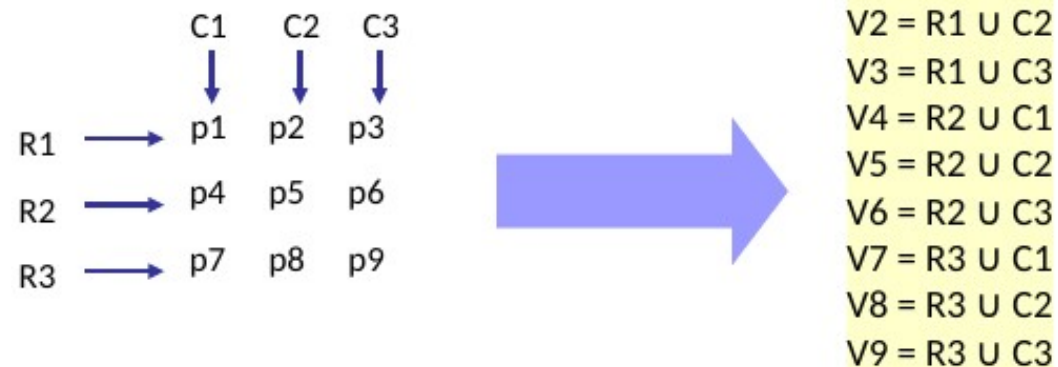


**Distributed
Mutual Exclusion**



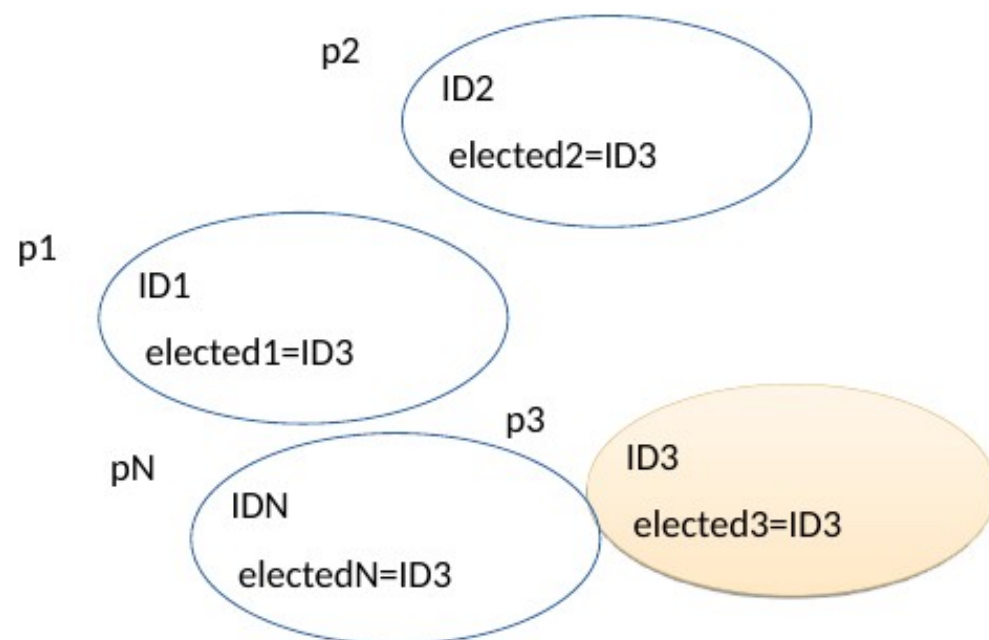
**Election
Mechanisms**

Summarizing



Coordination

- Maekawa Voting Algorithm
- Improvement over Ricard-Agrawala



Election

- Select a node for a special role
- Desirable properties & evaluation metrics
- Ring Algorithm
- Bully Algorithm



Lots of energy and success in your exams!

Questions?

Coordination

Part-2

Coordination