

10|10|23

Dependable Distributed Systems

Master of Science in Engineering in Computer Science

AA 2023/2024

LECTURE 7: LEADER ELECTION

Recap on Timing Assumptions

Synchronous

- timing assumptions are explicit either on
 - Bounds on process executions and communication channels, or
 - Existence of a common global clock, or
 - Both

Asynchronous

- there are no timing assumptions

Recap on Timing Assumptions

Partial synchrony requires abstract timing assumptions (after an unknown time t the system becomes synchronous)

Two choices:

1. Put assumption on the system model (including links and processes)
2. Create a separate abstractions that encapsulates those timing assumptions

Note: manipulating time inside a protocol/algorithm is complex and the correctness proof may become very involved and sometimes prone to errors

An alternative

Sometimes, we may be interested in knowing one process that is alive instead of monitoring failures

- E.g., Need of a coordinator

monitor the process
and elect an alive process

We can use a different oracle (called leader election module) that reports a process that is alive

as master / slave approach

Leader Election Specification

Module 2.7: Interface and properties of leader election

Module:

Name: LeaderElection, **instance** le .

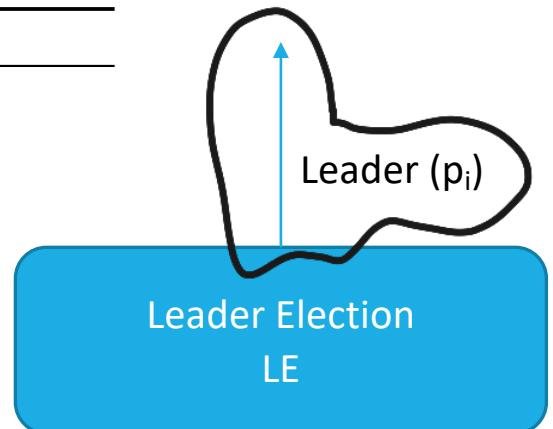
Events:

Indication: $\langle le, \text{Leader} | p \rangle$: Indicates that process p is elected as leader.

Properties:

LE1: *Eventual detection*: Either there is no correct process, or some correct process is eventually elected as the leader. **Liveness**

LE2: *Accuracy*: If a process is leader, then all previously elected leaders have crashed. **Safety**



→ he will remain the leader, until he will crash

Leader Election Implementation

Algorithm 2.6: Monarchical Leader Election

Implements:

LeaderElection, **instance** le .

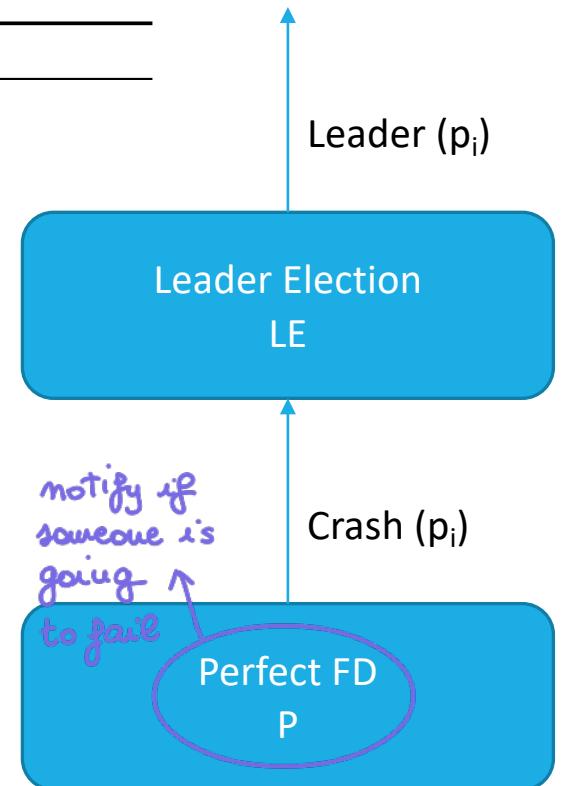
Uses:

PerfectFailureDetector, **instance** \mathcal{P} .

```
upon event < le, Init > do
    suspected := ∅;
    leader := ⊥; → no one is the leader

upon event <  $\mathcal{P}$ , Crash |  $p$  > do
    suspected := suspected ∪ { $p$ };

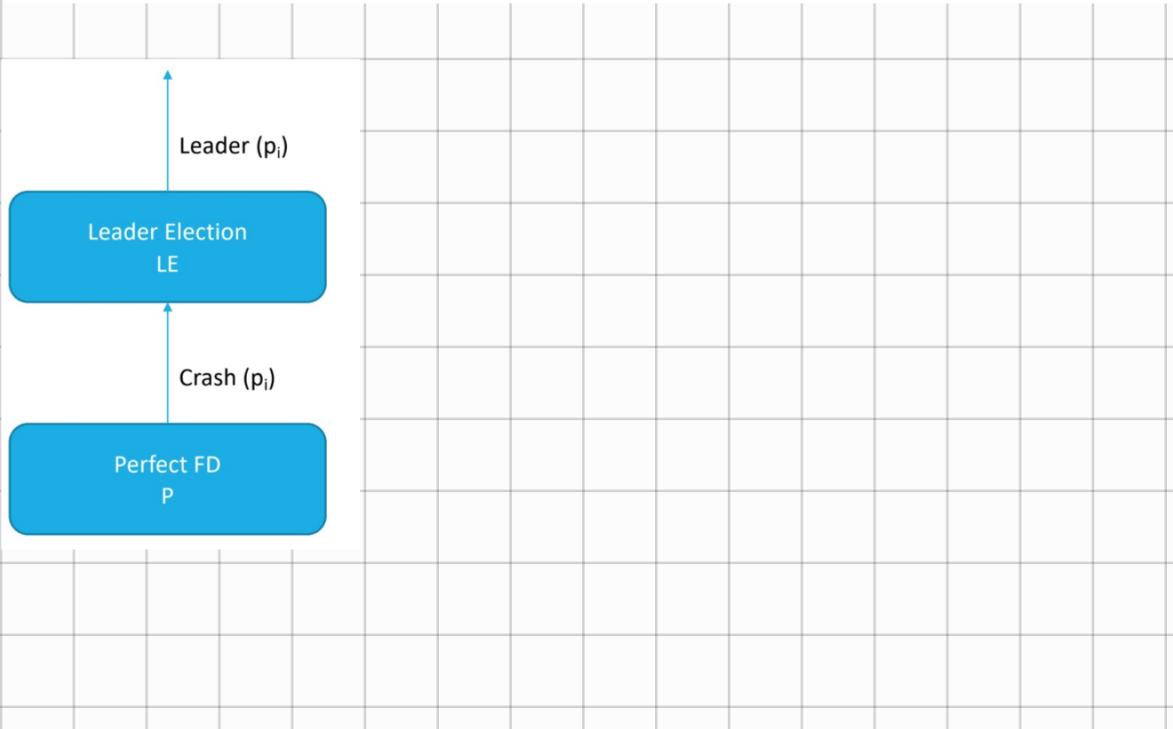
upon leader ≠ maxrank( $\Pi \setminus$  suspected) do
    leader := maxrank( $\Pi \setminus$  suspected);
    trigger < le, Leader | leader >;
    leader event
     $p_{10}$  CRASH
     $p_{10}$  is the leader
```



7: LEADER ELECTION

LE1: *Eventual detection:* Either there is no correct process, or some correct process is eventually elected as the leader.

LE2: *Accuracy:* If a process is leader, then all previously elected leaders have crashed.



Algorithm 2.6: Monarchical Leader Election

Implements:

LeaderElection, **instance** *le*.

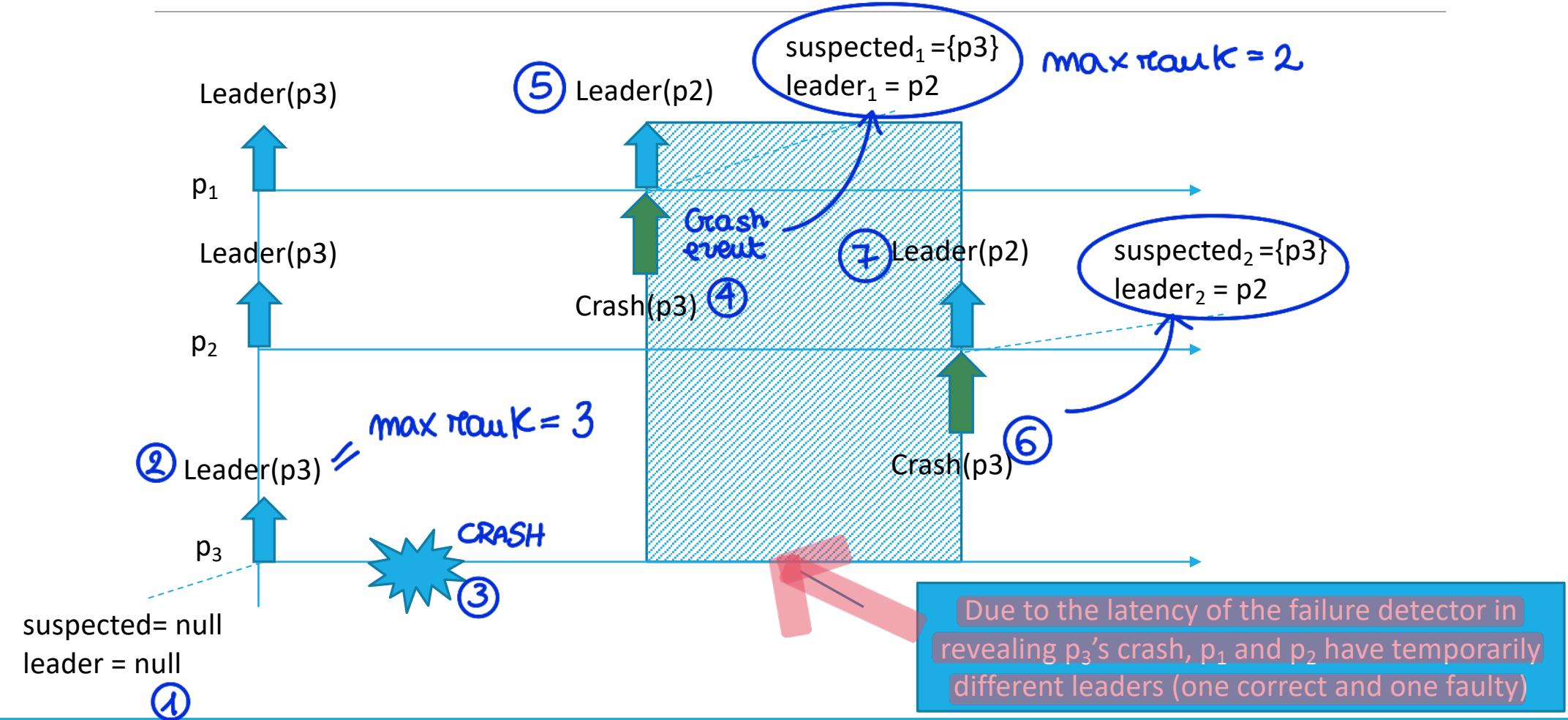
Uses:

PerfectFailureDetector, **instance** *P*.

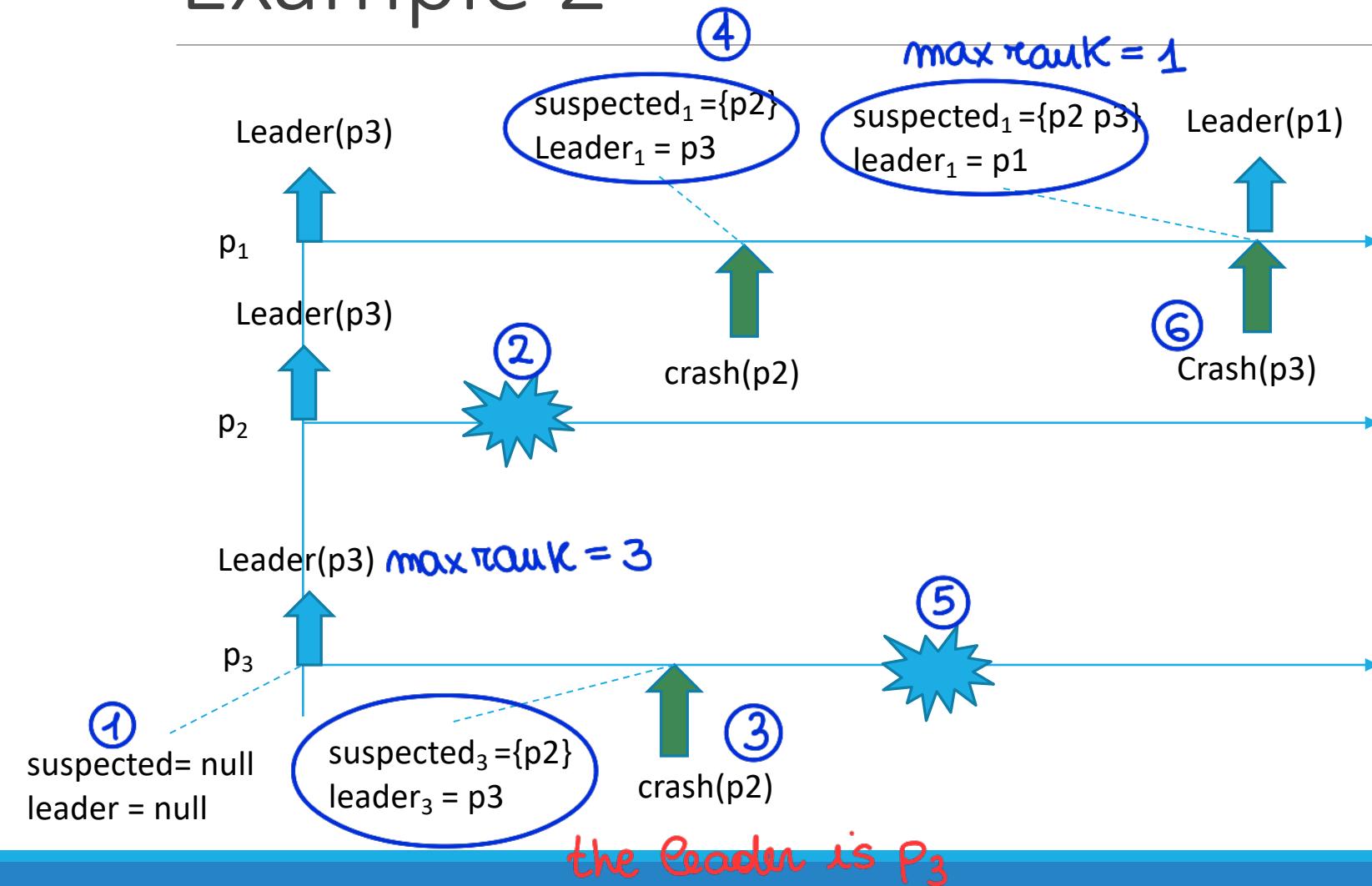
upon event $\langle le, \text{Init} \rangle$ **do**
① *suspected* := \emptyset ;
 leader := \perp ; *no one is the leader*

upon event $\langle P, \text{Crash} \mid p \rangle$ **do**
 suspected := *suspected* $\cup \{p\}$;
 sort the list *all the process* *the crashed processes*
 upon *leader* $\neq \text{maxrank}(P \setminus \text{suspected}) **do** \rightarrow WHO IS THE 1^o IN THE LIST?
 ② *leader* := *maxrank*(*P* \setminus *suspected*);
 trigger $\langle le, \text{Leader} \mid \text{leader} \rangle$;$

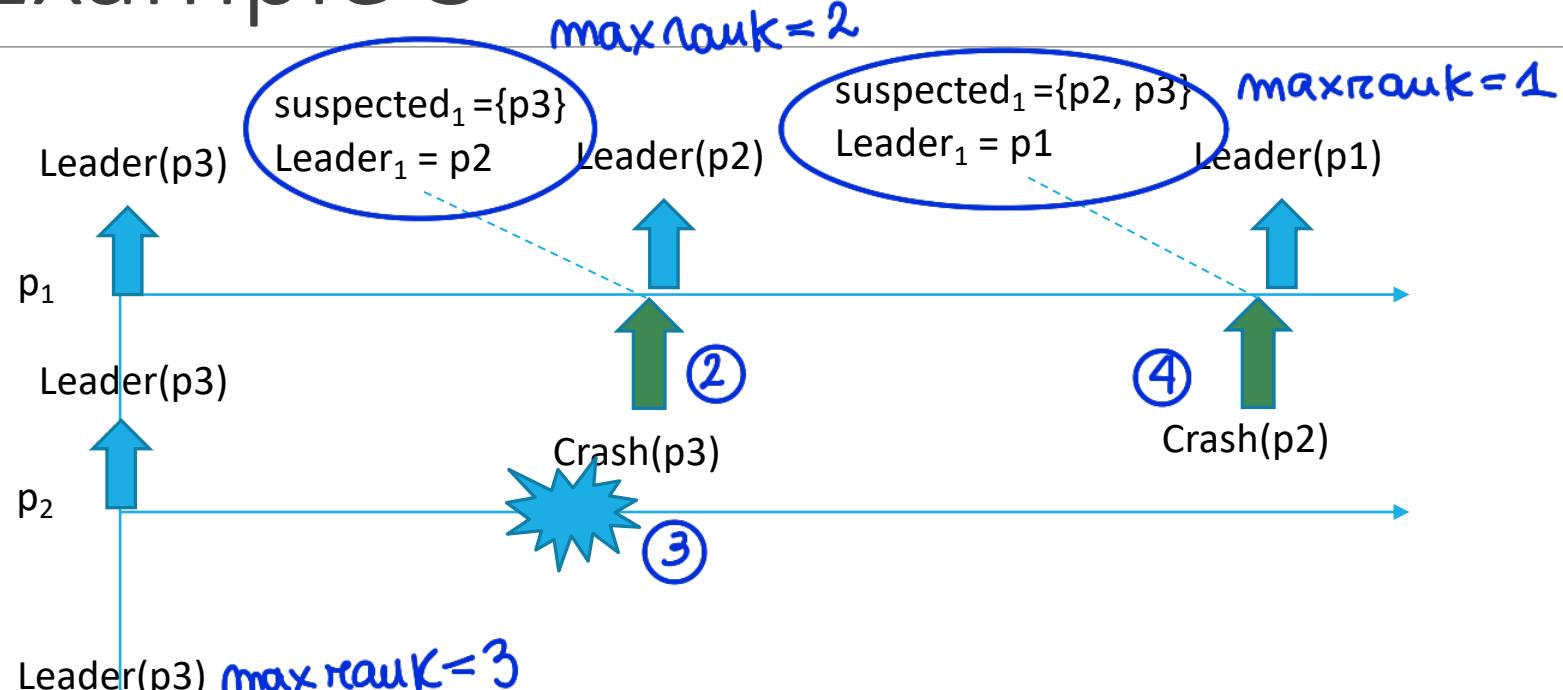
Example



Example 2



Example 3

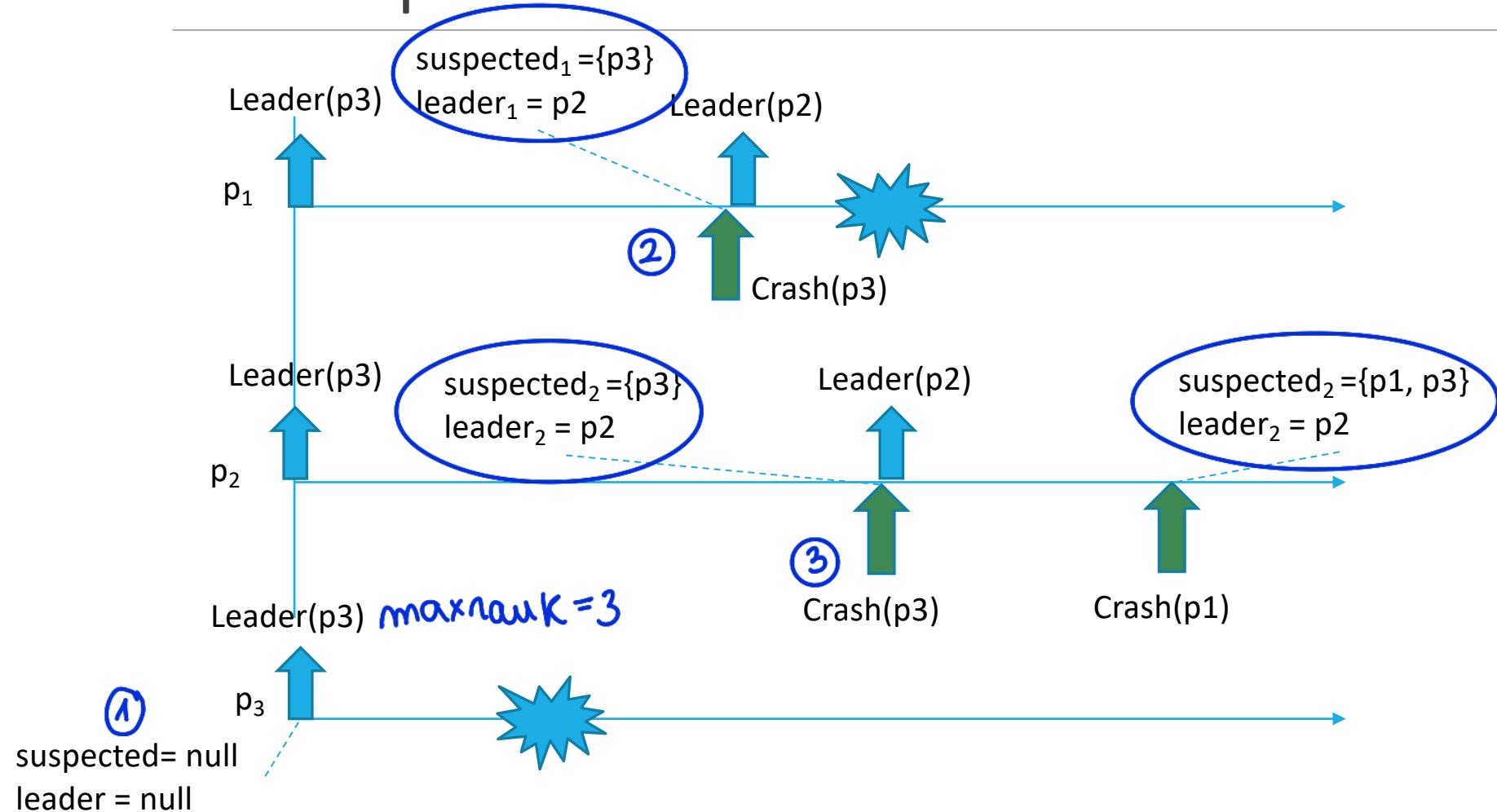


Leader(p3) maxrank = 3

①

suspected = null
leader = null

Example 4



Correctness

What if the Failure detector is not perfect?

Eventual leader election (Ω)

Module 2.9: Interface and properties of the eventual leader detector

Module:

Name: EventualLeaderDetector, **instance** Ω .

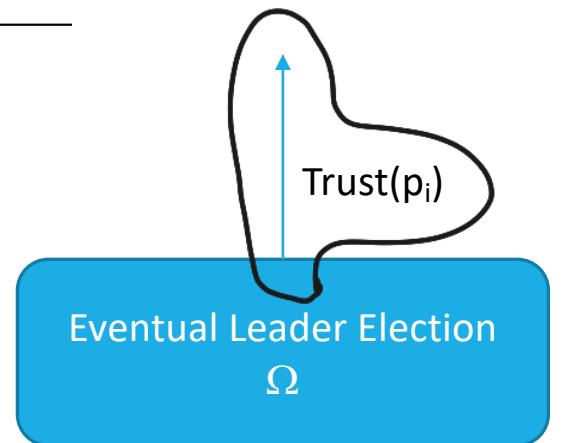
Events:

Indication: $\langle \Omega, Trust | p \rangle$: Indicates that process p is trusted to be leader.

Properties:

ELD1: *Eventual accuracy*: There is a time after which every correct process trusts some correct process.

ELD2: *Eventual agreement*: There is a time after which no two correct processes trust different correct processes.



lost the **ACCURACY** → if a p is a leader, the previous leaders are crashed

Observation on Ω

Ω ensures that *eventually* correct processes will elect the same correct process as their leader

Ω does not guarantee that

- Leaders change in an arbitrary manner and for an arbitrary period of time
- many leaders might be elected during the same period of time without having crashed

based on the period of
ASYNCHRONOUS

Once a unique leader is determined, and does not change again, we say that the leader has *stabilized*

Eventual leader election (Ω)

Using Crash-stop process abstraction

- Obtained directly by $\text{<}P$ by using a deterministic rule on processes that are not suspected by $\text{<}P$
- trust the process with the highest identifier among all processes that are not suspected by $\text{<}P$

IMPORTANT ASSUMPTION

- There always exists at least one correct process (otherwise Ω cannot be built)

Ω Implementation

Algorithm 2.8: Monarchical Eventual Leader Detection

Implements:

EventualLeaderDetector, **instance** Ω .

Uses:

EventuallyPerfectFailureDetector, **instance** $\diamond \mathcal{P}$.

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

$\text{suspected} := \emptyset;$

$\text{leader} := \perp;$

upon event $\langle \diamond \mathcal{P}, \text{Suspect} \mid p \rangle$ **do**

$\text{suspected} := \text{suspected} \cup \{p\};$

upon event $\langle \diamond \mathcal{P}, \text{Restore} \mid p \rangle$ **do**

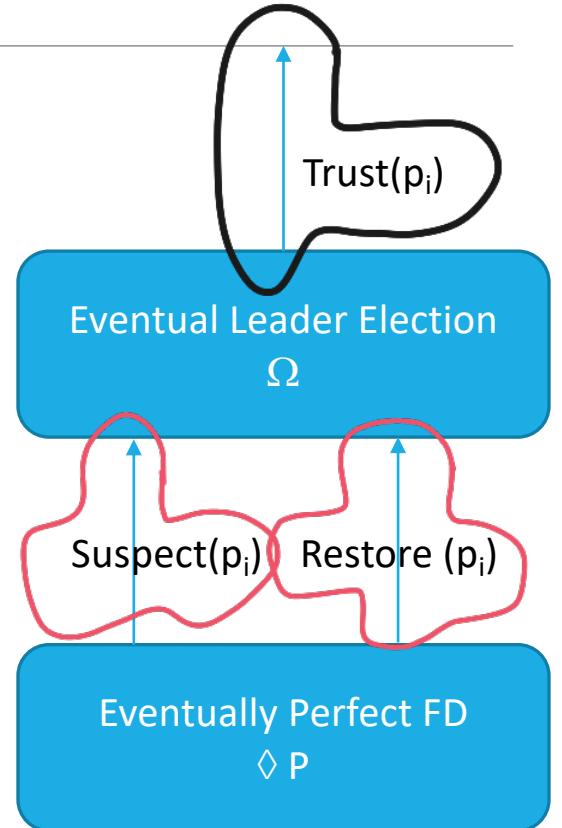
$\text{suspected} := \text{suspected} \setminus \{p\};$

upon $\text{leader} \neq \text{maxrank}(\Pi \setminus \text{suspected})$ **do**

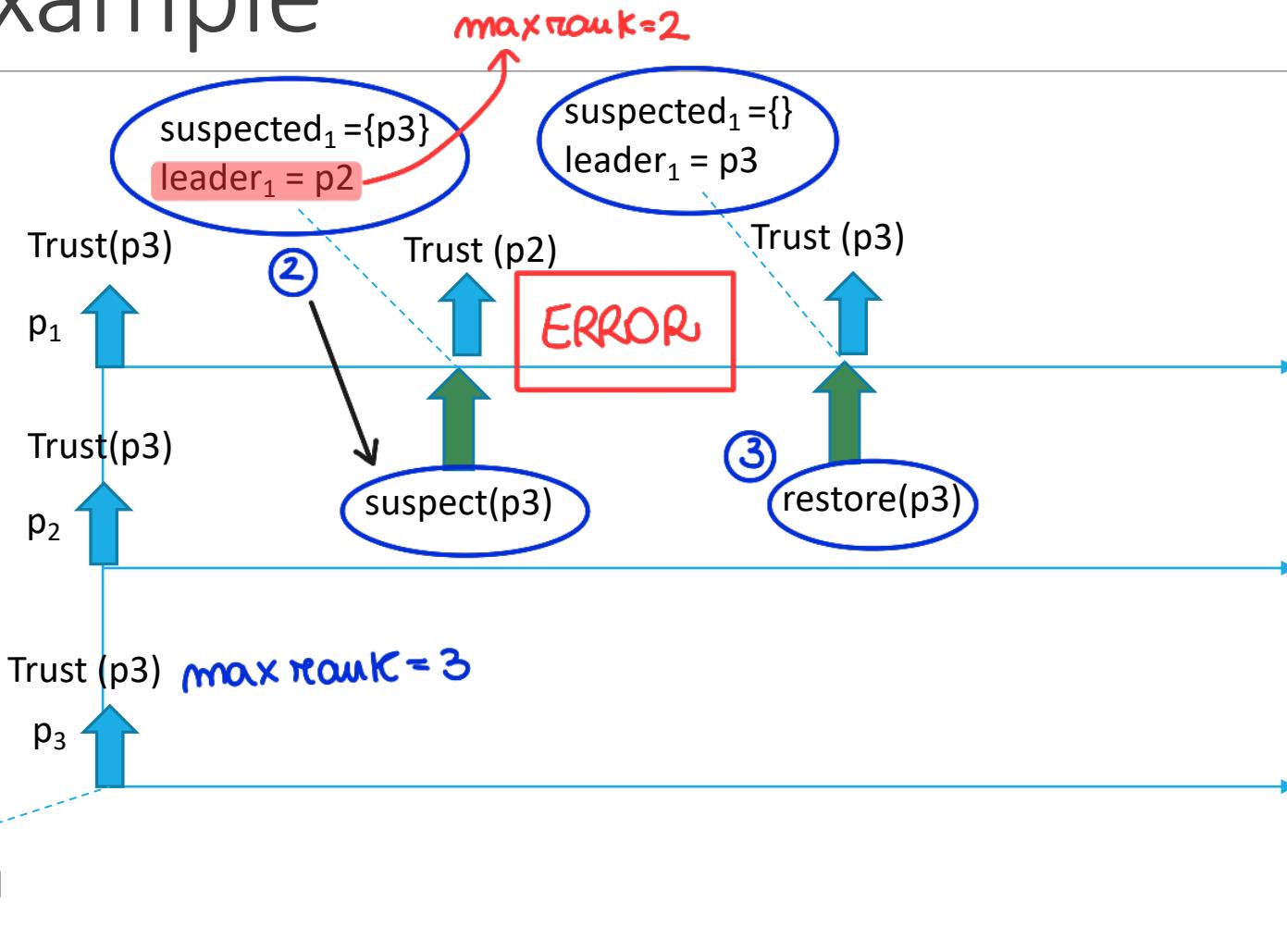
$\text{leader} := \text{maxrank}(\Pi \setminus \text{suspected});$

trigger $\langle \Omega, \text{Trust} \mid \text{leader} \rangle;$

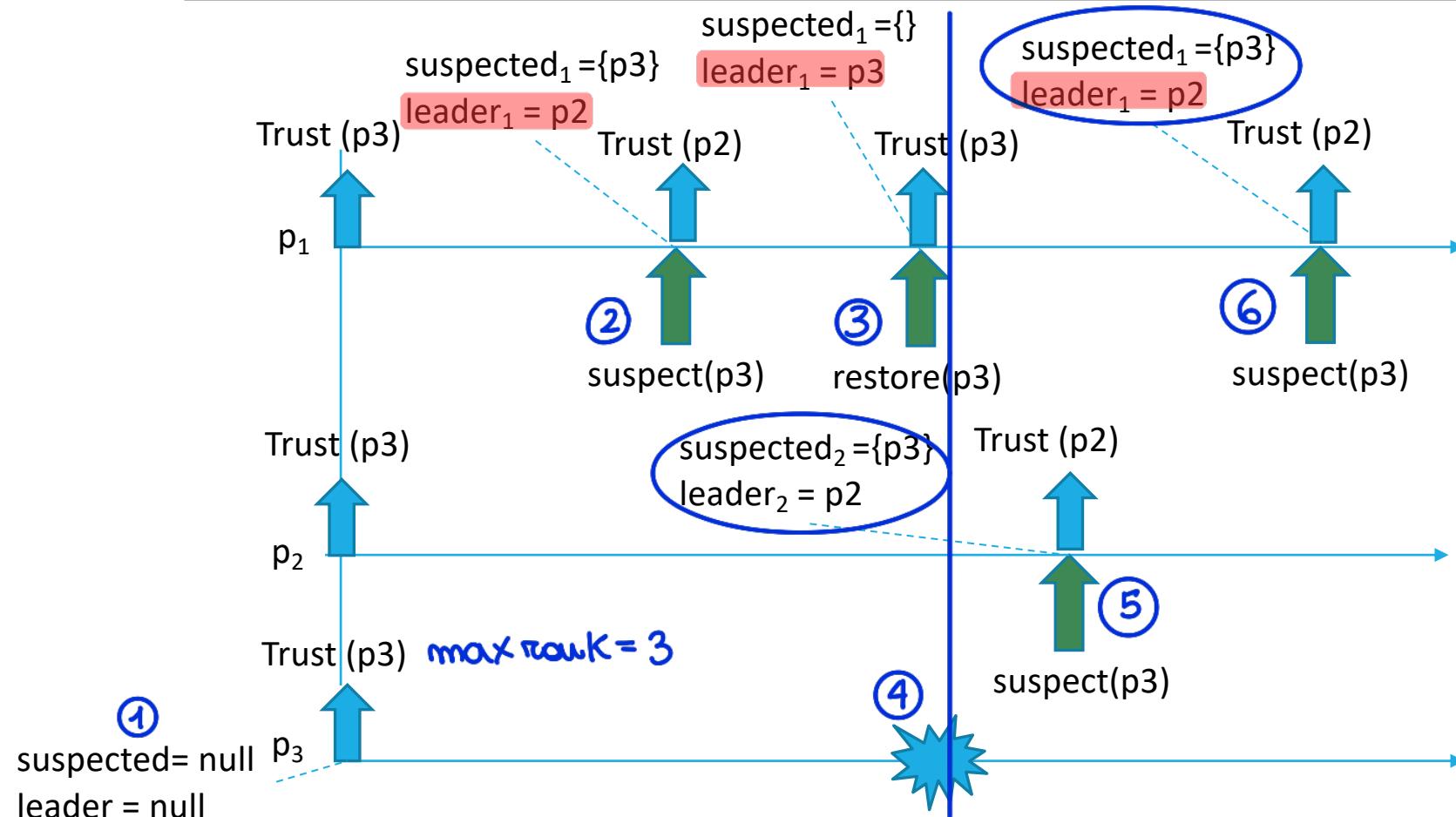
} deterministic
rule



Example



Example 2



Eventual leader election (Ω)

System model

- Crash-Recovery
- Partial synchrony

Under this assumption, a **correct process** means:

1. A process that does not crash or
2. A process that crashes, eventually recovers and never crashes again

Ω With crash-recovery, fair lossy links and timeouts

Algorithm 2.9: Elect Lower Epoch

Implements:

EventualLeaderDetector, instance Ω .

Uses:

FairLossPointToPointLinks, instance fl .

```

upon event <  $\Omega$ , Init > do
    epoch := 0;
    store(epoch);
    candidates :=  $\emptyset$ ;
    trigger <  $\Omega$ , Recovery >;

```

period of time
in which you
are suspected crash
and then recover

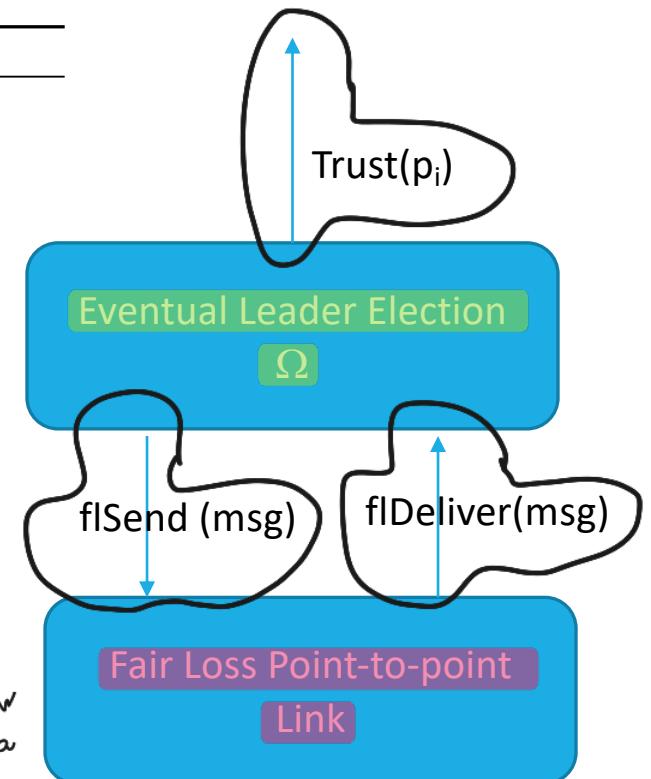
count how many
time you have CRASH-RECOVER

keeps track of how many
times the process
crashed and recovered

```

upon event <  $\Omega$ , Recovery > do
    leader := maxrank( $\Pi$ );
    trigger <  $\Omega$ , Trust | leader >;
    delay :=  $\Delta$ ;
    retrieve(epoch);
    epoch := epoch + 1;
    store(epoch);
    forall  $p \in \Pi$  do
        trigger <  $fl$ , Send |  $p$ , [HEARTBEAT, epoch] >;
    candidates :=  $\emptyset$ ;
    starttimer(delay);

```



Ω With crash-recovery, fair lossy links and timeouts

Algorithm 2.9: Elect Lower Epoch

Implements:

EventualLeaderDetector, instance Ω .

Uses:

FairLossPointToPointLinks, instance fl .

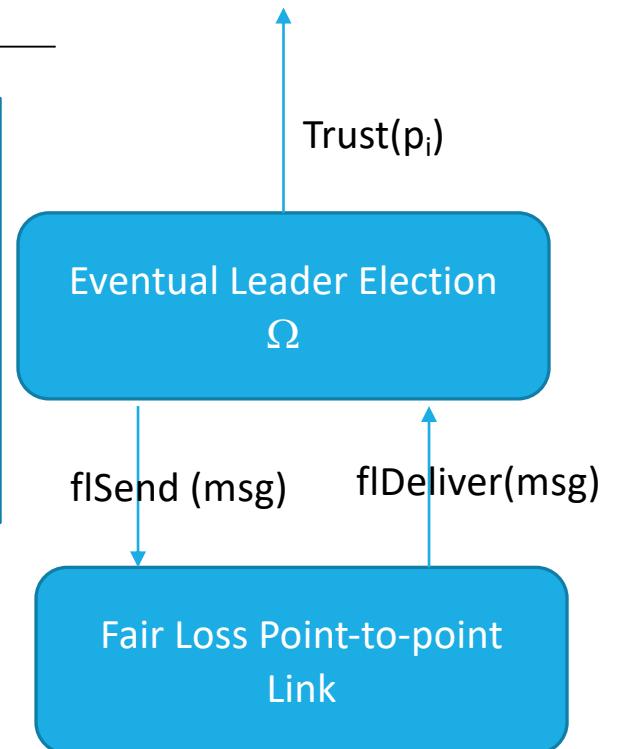
```

upon event < Timeout > do
    newleader := select(candidates);
    if newleader ≠ leader then → may there were  
a c-R
        delay := delay + Δ;
        leader := newleader;
        trigger <  $\Omega$ , Trust | leader >;
    forall  $p \in \Pi$  do
        trigger <  $fl$ , Send |  $p$ , [HEARTBEAT, epoch] >;
        candidates := ∅;
        starttimer(delay);

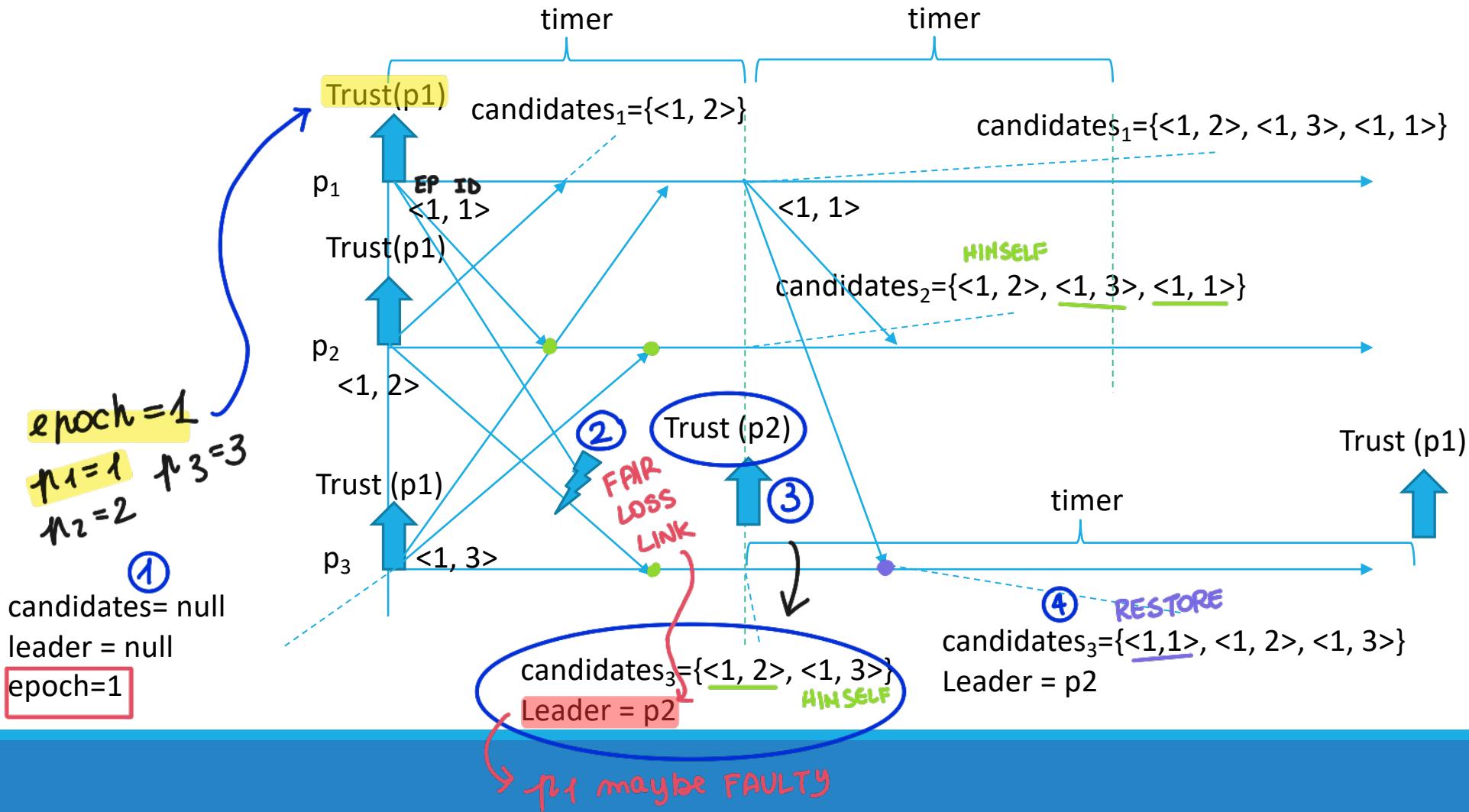
upon event <  $fl$ , Deliver |  $q$ , [HEARTBEAT, ep] > do
    if exists  $(s, e) \in candidates$  such that  $s = q \wedge e < ep$  then
        candidates := candidates \  $\{(q, e)\}$ ;
        candidates := candidates ∪  $(q, ep)$ ;

```

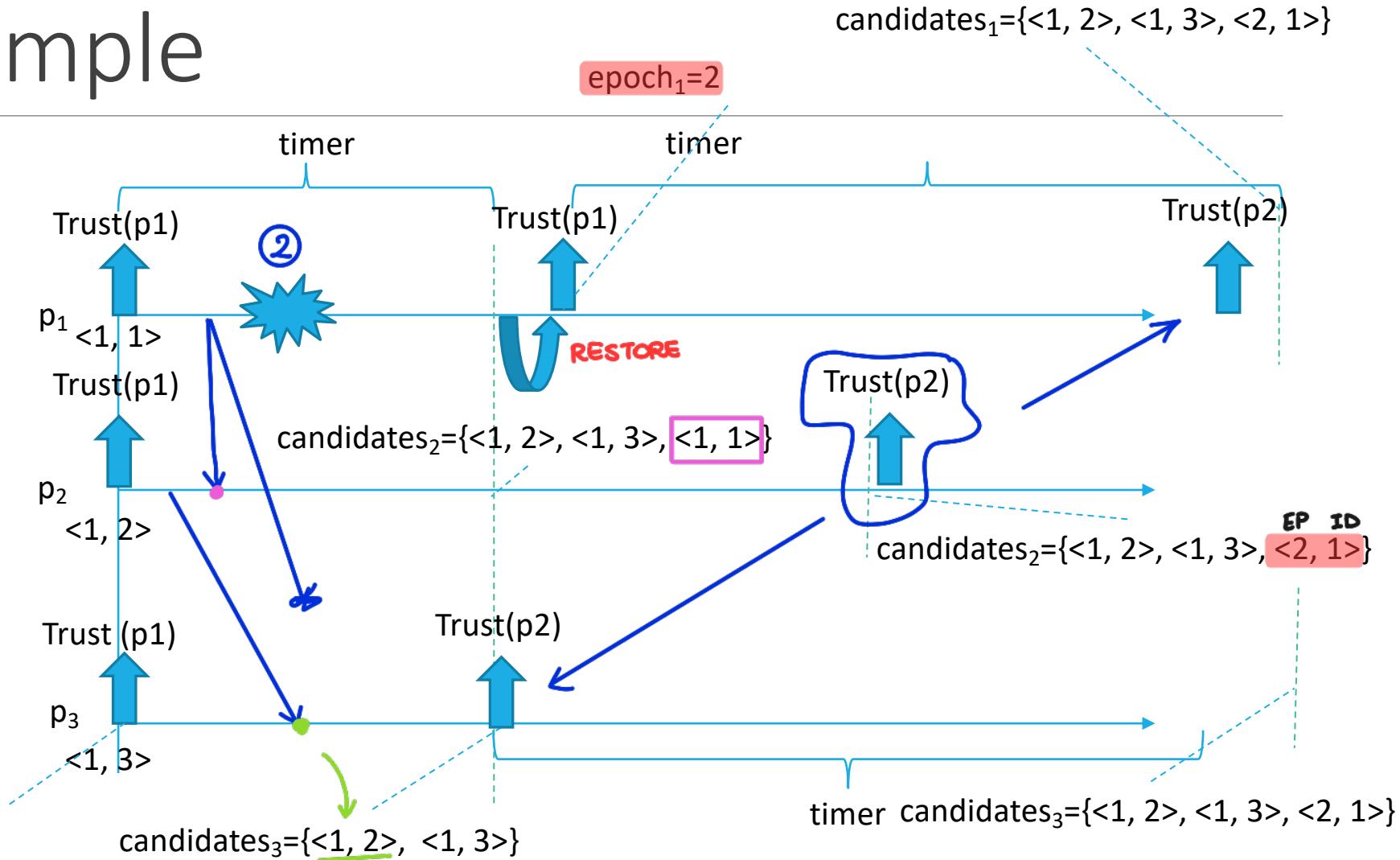
deterministic function returning one process among all candidates (i.e., process with the lowest epoch number and among the ones with the same epoch number the one with the lowest identifier)



Example



Example



References

C. Cachin, R. Guerraoui and L. Rodrigues. Introduction to Reliable and Secure Distributed Programming, Springer, 2011

- Chapter 2 – from Section 2.6.1 to Section 2.6.5