

# Distributed Systems

## Master of Science in Engineering in Computer Science

AA 2019/2020

---

LECTURE 5: 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

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

# Leader Election Specification

---

---

**Module 2.7:** Interface and properties of leader election

---

**Module:**

**Name:** LeaderElection, **instance** *le*.

**Events:**

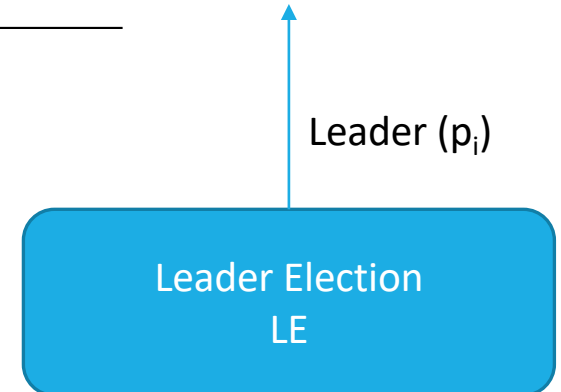
**Indication:**  $\langle le, Leader \mid p \rangle$ : Indicates that process  $p$  is elected as leader.

**Properties:**

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

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

---



# Leader Election Implementation

---

---

**Algorithm 2.6:** Monarchical Leader Election

---

**Implements:**

LeaderElection, **instance**  $le$ .

**Uses:**

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

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

$suspected := \emptyset$ ;

$leader := \perp$ ;

**upon event**  $\langle \mathcal{P}, Crash \mid p \rangle$  **do**

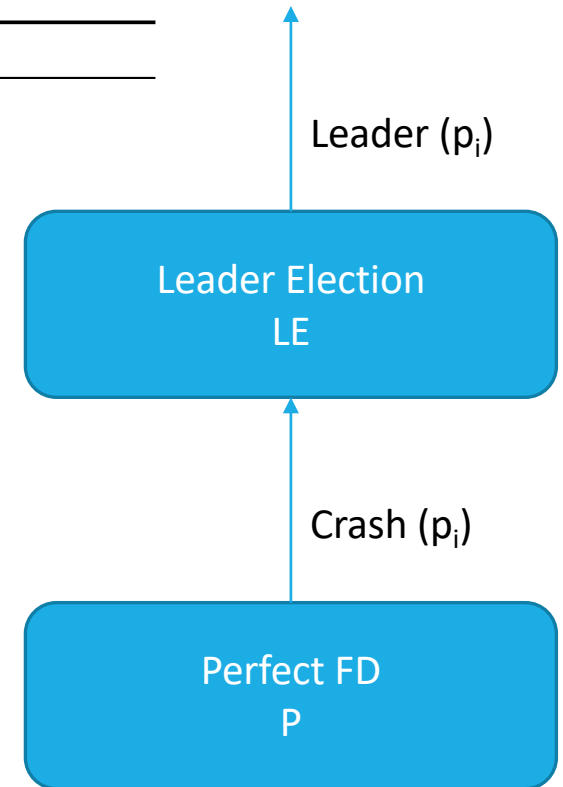
$suspected := suspected \cup \{p\}$ ;

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

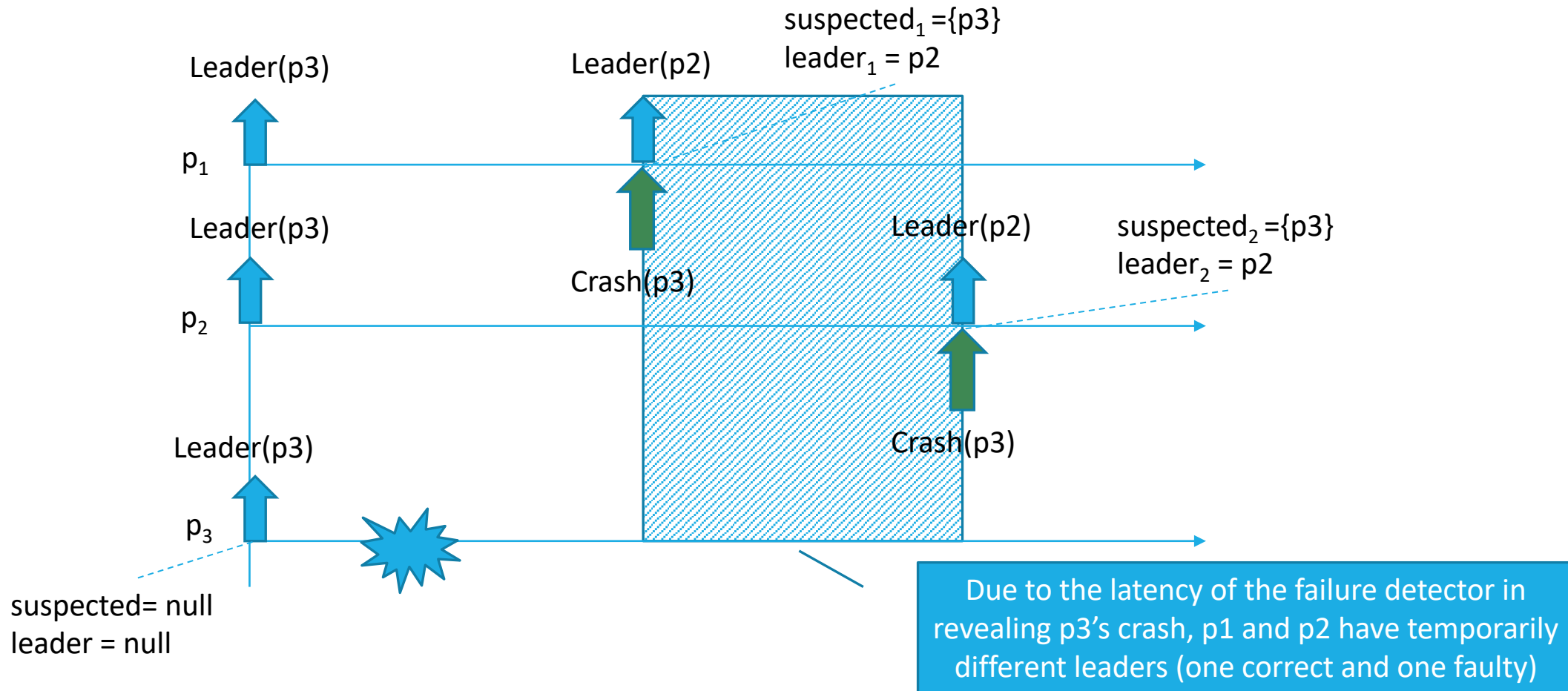
$leader := \text{maxrank}(\Pi \setminus suspected)$ ;

**trigger**  $\langle le, Leader \mid leader \rangle$ ;

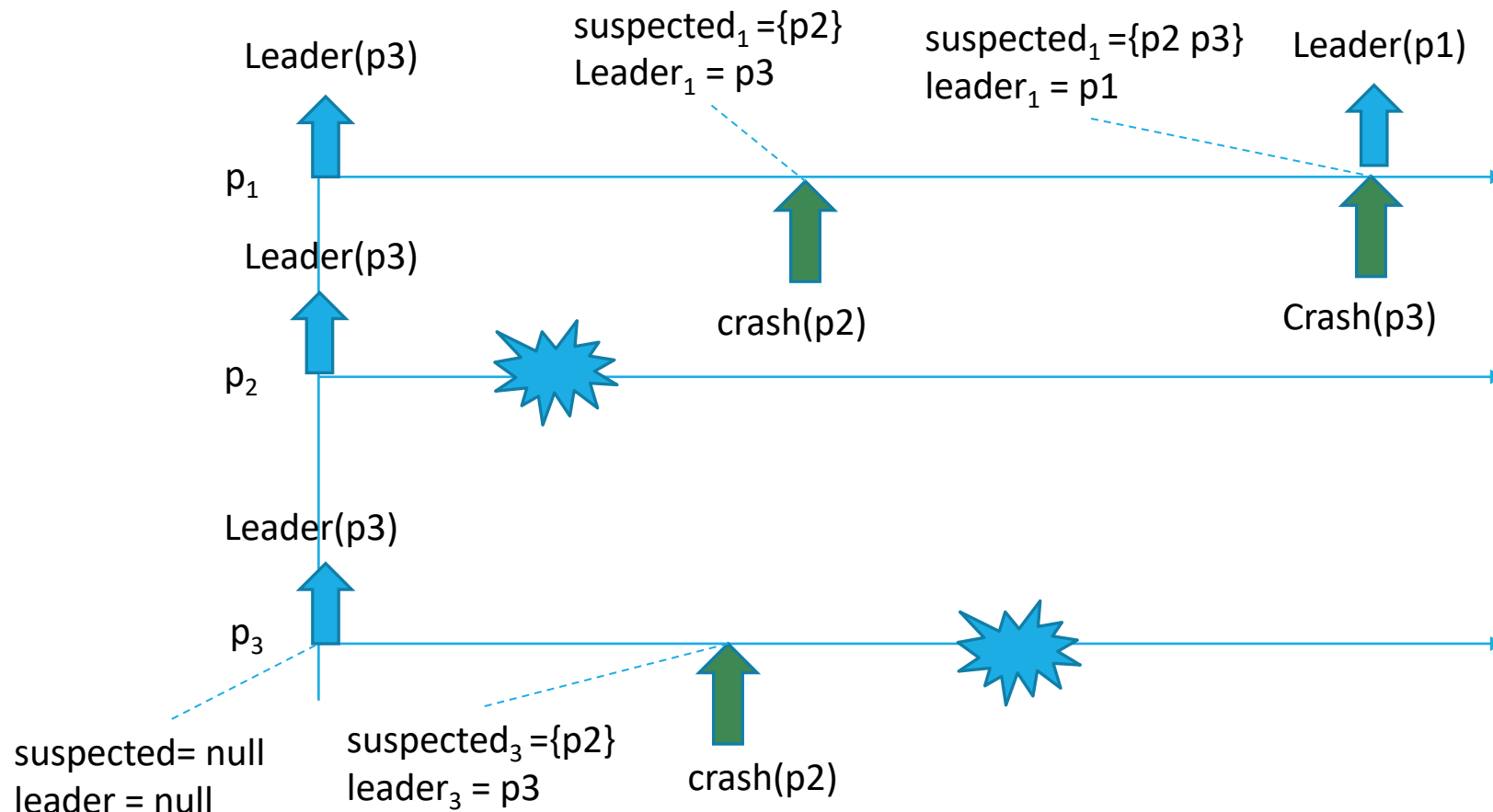
---



# Example

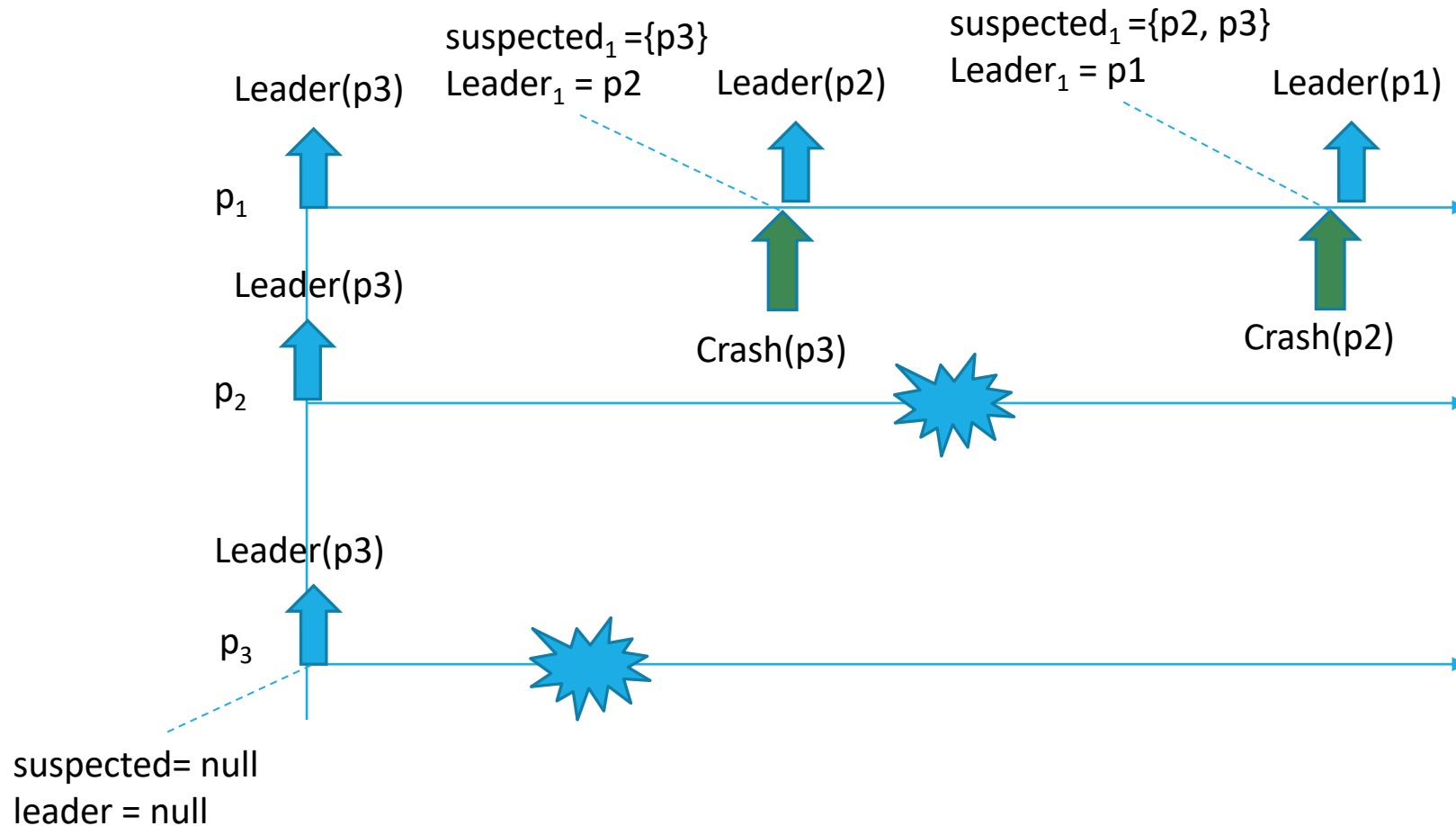


# Example 2

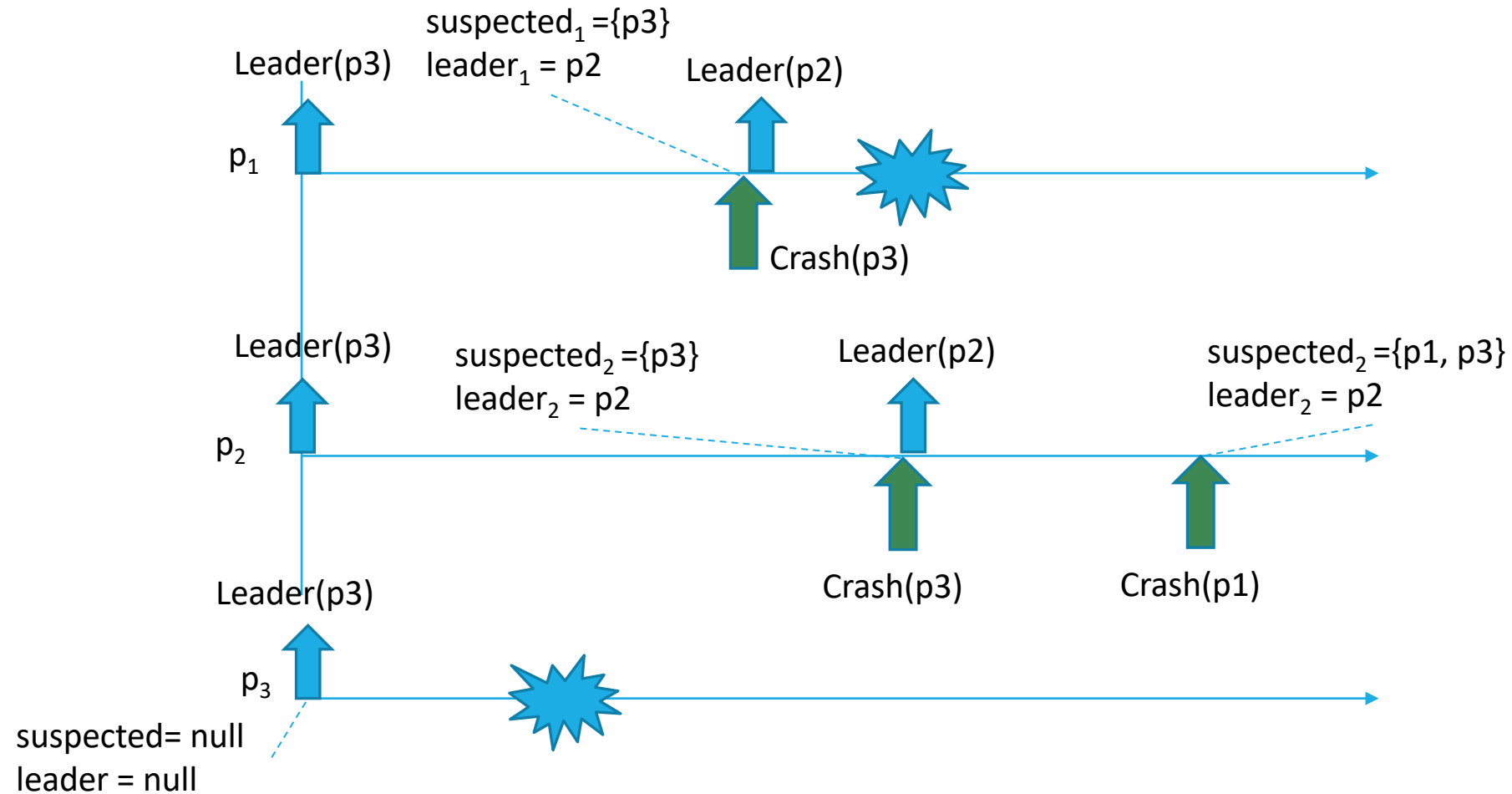




# Example 3



## Example 4



# Correctness

---

What if the Failure detector is not perfect?

# Eventual leader election ( $\Omega$ )

---

---

**Module 2.9:** Interface and properties of the eventual leader detector

---

**Module:**

**Name:** EventualLeaderDetector, **instance**  $\Omega$ .

**Events:**

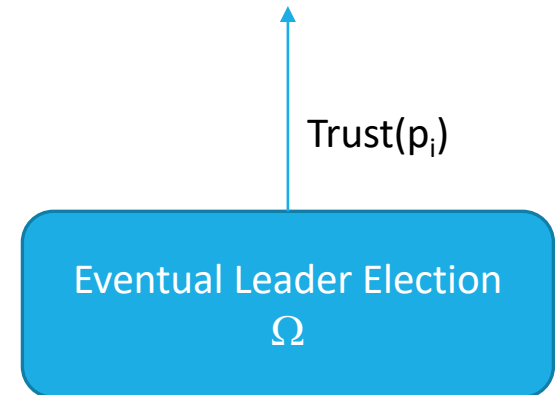
**Indication:**  $\langle \Omega, Trust \mid p \rangle$ : Indicates that process  $p$  is trusted to be leader.

**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.

---



# Observation on $\Omega$

---

$\Omega$  ensures that *eventually* correct processes will elect the same correct process as their leader

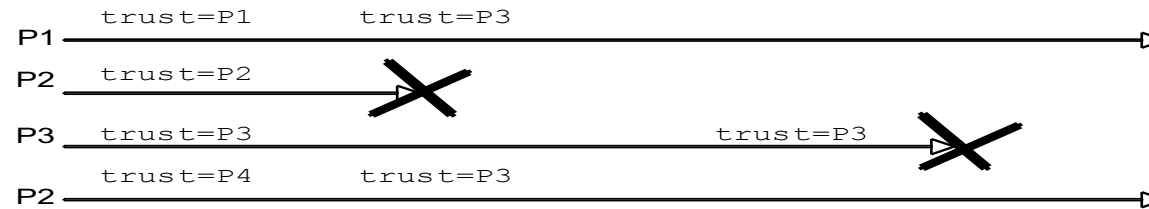
$\Omega$  does not guarantee that

- 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

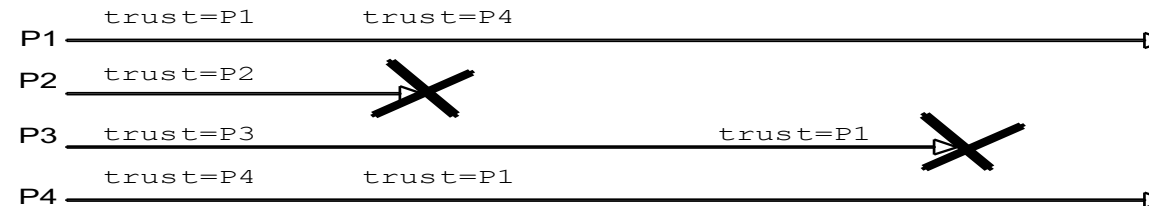
Once a unique leader is determined, and does not change again, we say that the leader has *stabilized*

# Study of Properties

Run 1



Run 2



	Run 1	Run 2
Eventual Accuracy	Not verified	Verified
Eventual Agreement	Verified	Not verified

# Eventual leader election ( $\Omega$ )

---

## Using Crash-stop process abstraction

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

Assume the existence of a correct process (otherwise  $\Omega$  cannot be built)

# $\Omega$ Implementation

---

**Algorithm 2.8:** Monarchical Eventual Leader Detection

---

**Implements:**

EventualLeaderDetector, **instance**  $\Omega$ .

**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**

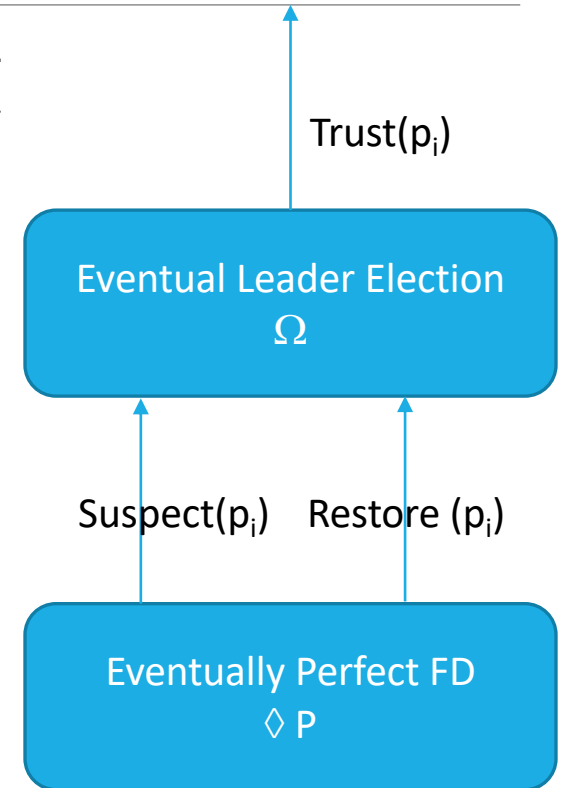
$\text{suspected} := \text{suspected} \setminus \{p\};$

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

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

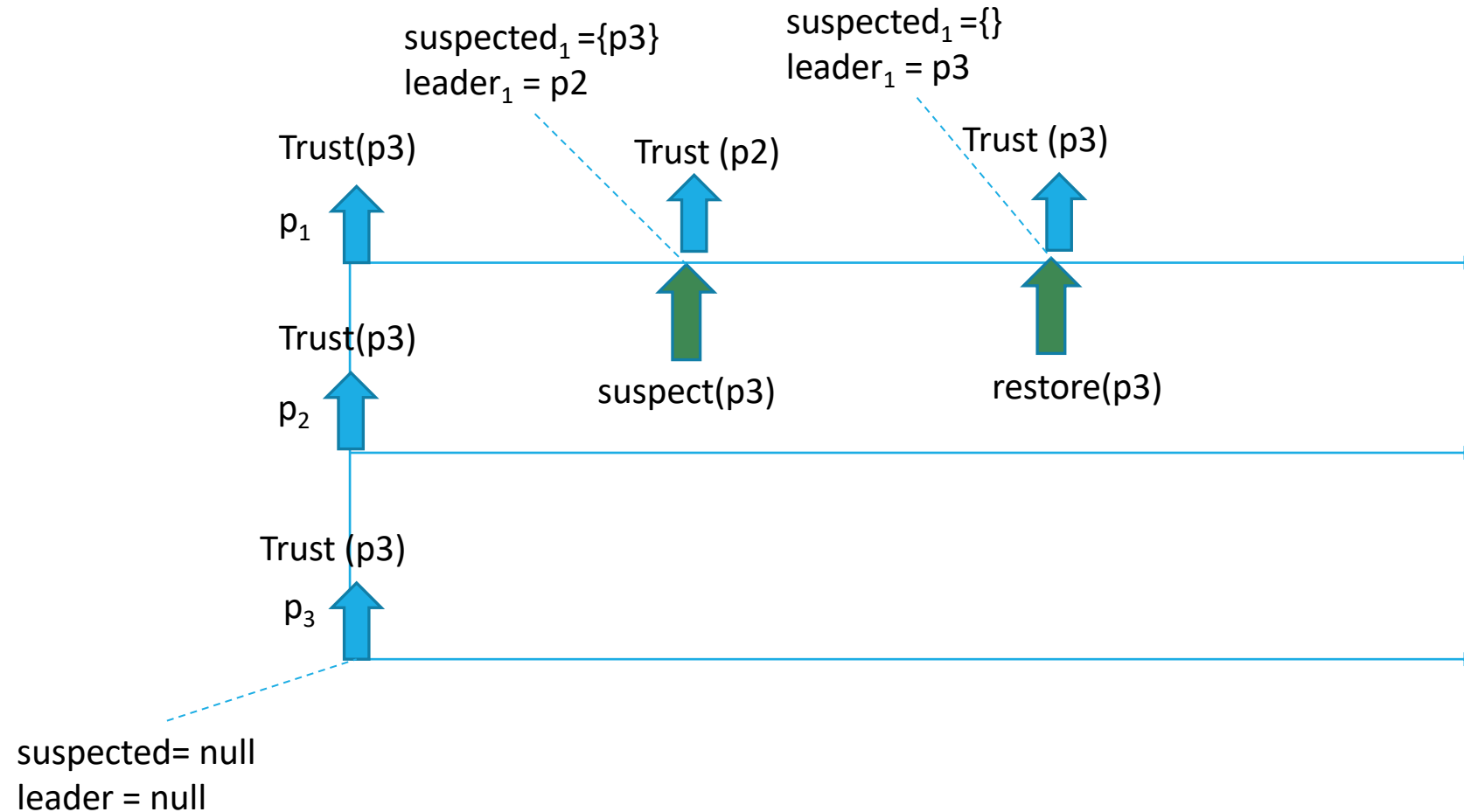
**trigger**  $\langle \Omega, \text{Trust} \mid \text{leader} \rangle;$

---

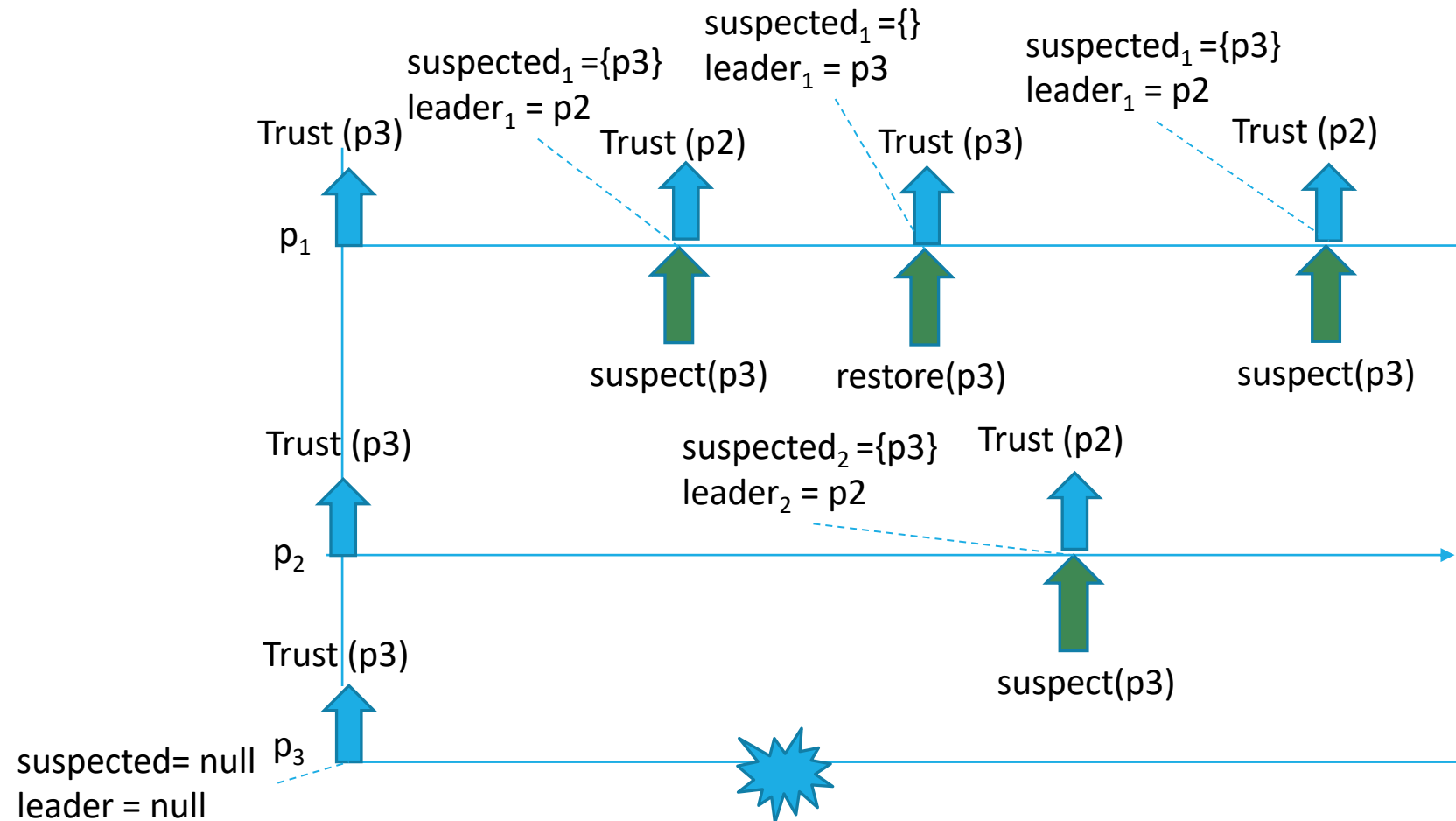




# Example



# Example 2



# Eventual leader election ( $\Omega$ )

---

## 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

# $\Omega$ With crash-recovery, fair lossy links and timeouts

---

**Algorithm 2.9:** Elect Lower Epoch

---

**Implements:**

EventualLeaderDetector, **instance**  $\Omega$ .

**Uses:**

FairLossPointToPointLinks, **instance**  $fll$ .

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

$epoch := 0$ ;

$store(epoch)$ ;

$candidates := \emptyset$ ;

**trigger**  $\langle \Omega, \text{Recovery} \rangle$ ;

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

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

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

$delay := \Delta$ ;

$retrieve(epoch)$ ;

$epoch := epoch + 1$ ;

$store(epoch)$ ;

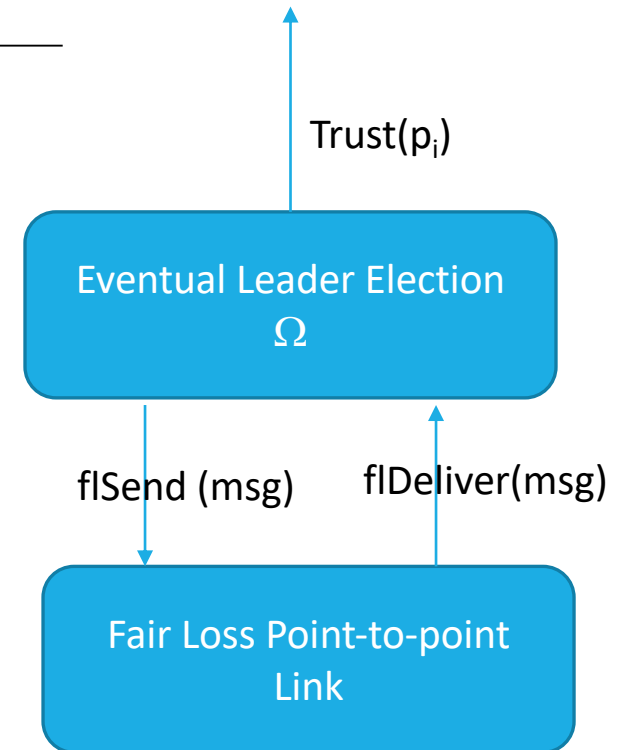
**forall**  $p \in \Pi$  **do**

**trigger**  $\langle fll, \text{Send} \mid p, [\text{HEARTBEAT}, epoch] \rangle$ ;

$candidates := \emptyset$ ;

$starttimer(delay)$ ;

keeps track of how many  
times the process  
crashed and recovered



# $\Omega$ With crash-recovery, fair lossy links and timeouts

---

**Algorithm 2.9:** Elect Lower Epoch

---

**Implements:**

EventualLeaderDetector, **instance**  $\Omega$ .

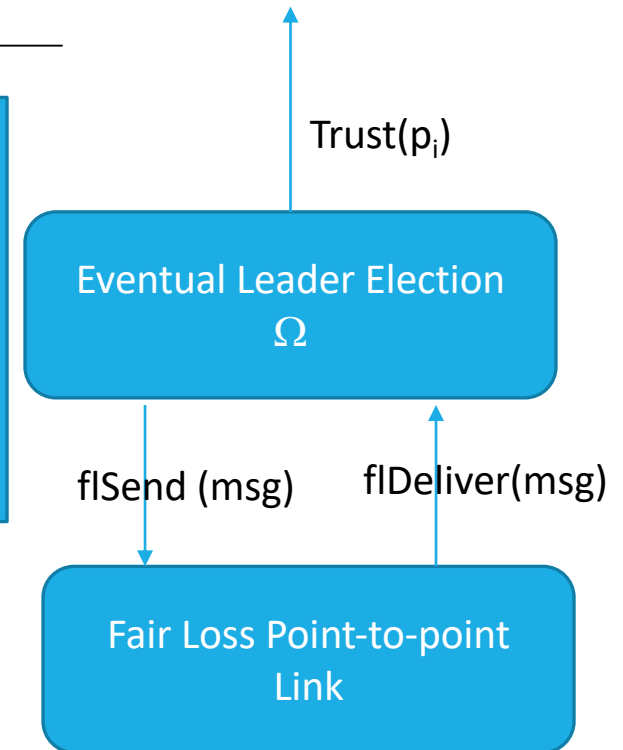
**Uses:**

FairLossPointToPointLinks, **instance**  $fll$ .

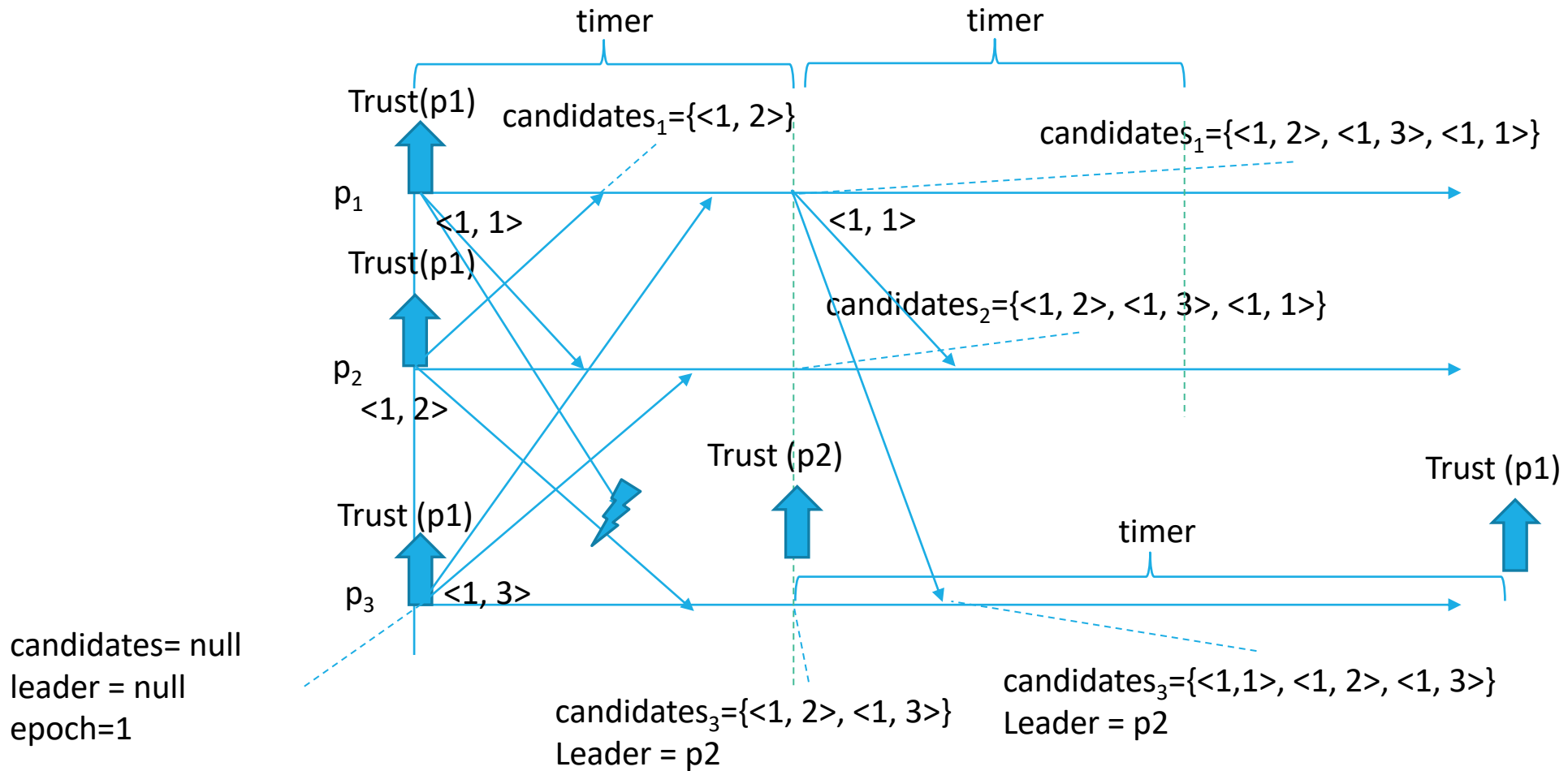
```
upon event  $\langle \text{Timeout} \rangle$  do
     $\text{newleader} := \text{select}(\text{candidates});$ 
    if  $\text{newleader} \neq \text{leader}$  then
         $\text{delay} := \text{delay} + \Delta;$ 
         $\text{leader} := \text{newleader};$ 
        trigger  $\langle \Omega, \text{Trust} \mid \text{leader} \rangle;$ 
    forall  $p \in \Pi$  do
        trigger  $\langle fll, \text{Send} \mid p, [\text{HEARTBEAT}, \text{epoch}] \rangle;$ 
     $\text{candidates} := \emptyset;$ 
     $\text{starttimer}(\text{delay});$ 

upon event  $\langle fll, \text{Deliver} \mid q, [\text{HEARTBEAT}, ep] \rangle$  do
    if exists  $(s, e) \in \text{candidates}$  such that  $s = q \wedge e < ep$  then
         $\text{candidates} := \text{candidates} \setminus \{(q, e)\};$ 
     $\text{candidates} := \text{candidates} \cup (q, ep);$ 
```

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



# 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