

04|10|23

Dependable Distributed Systems

Master of Science in Engineering in Computer Science

AA 2023/2024

LECTURE 6: FAILURE DETECTION

Recap on Timing Assumptions

Synchronous

timing assumptions are explicit and known

- Upper bounds on process execution time
- Upper bounds on communication time
- Existence of a common global clock (apart from synchronization error)
abstraction of physical time

Asynchronous

there are no timing assumptions no physical clock

- the only ordering is given by the happen-before relationship
- it is possible to build scalar or vector clocks
 \downarrow causal relat.

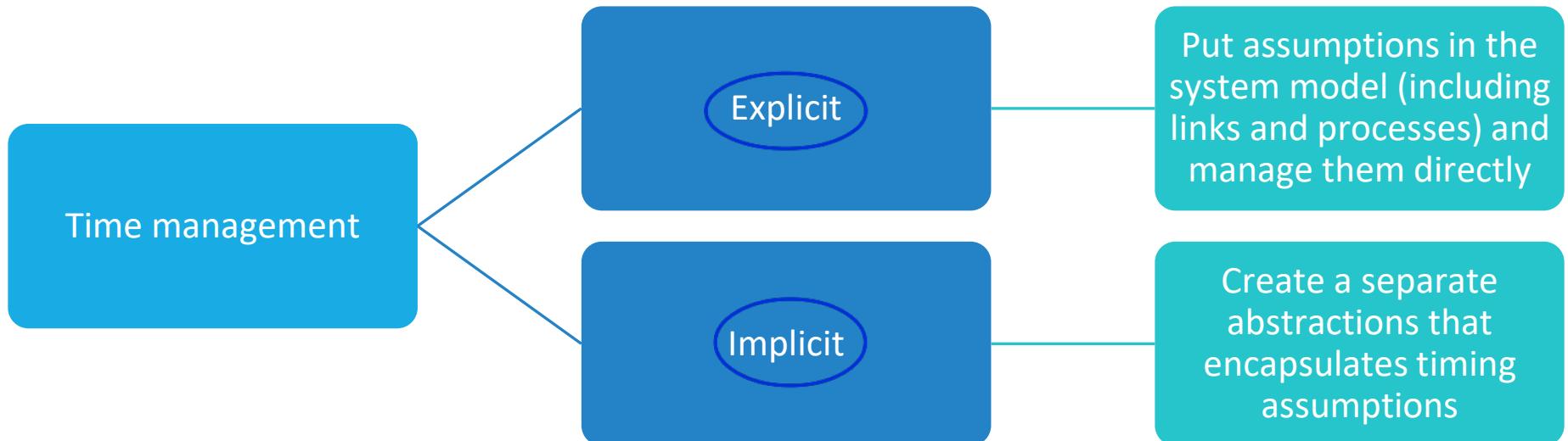
Partial synchrony

- requires abstract timing assumptions (after an unknown time t the system becomes synchronous)
 \downarrow
 asynchronous most of the time
 then synchronous

Recap on Timing Assumptions

NOTE

- *manipulating time inside a protocol/algorithim is complex* and the correctness proof may become very involved and sometimes prone to errors



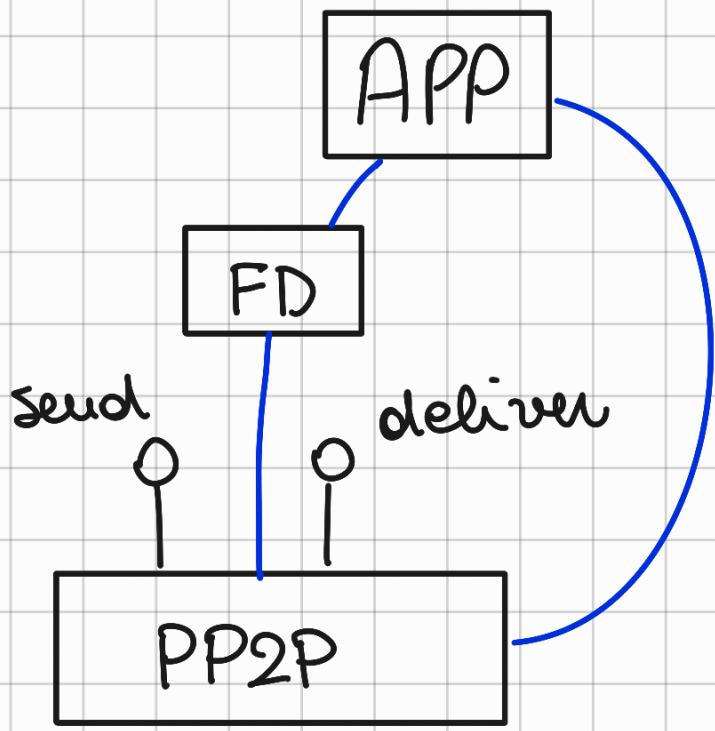
Failure Detector Abstraction

A Failure Detector is an oracle providing information about the failure state of a process in the system

- it is a software module to be used together with process and link abstractions
- It encapsulates timing assumptions of either partially synchronous or fully synchronous system
- *Stronger are the timing assumption, more accurate the information provided by a failure detector will be*

A failure detector is generally described by two properties:

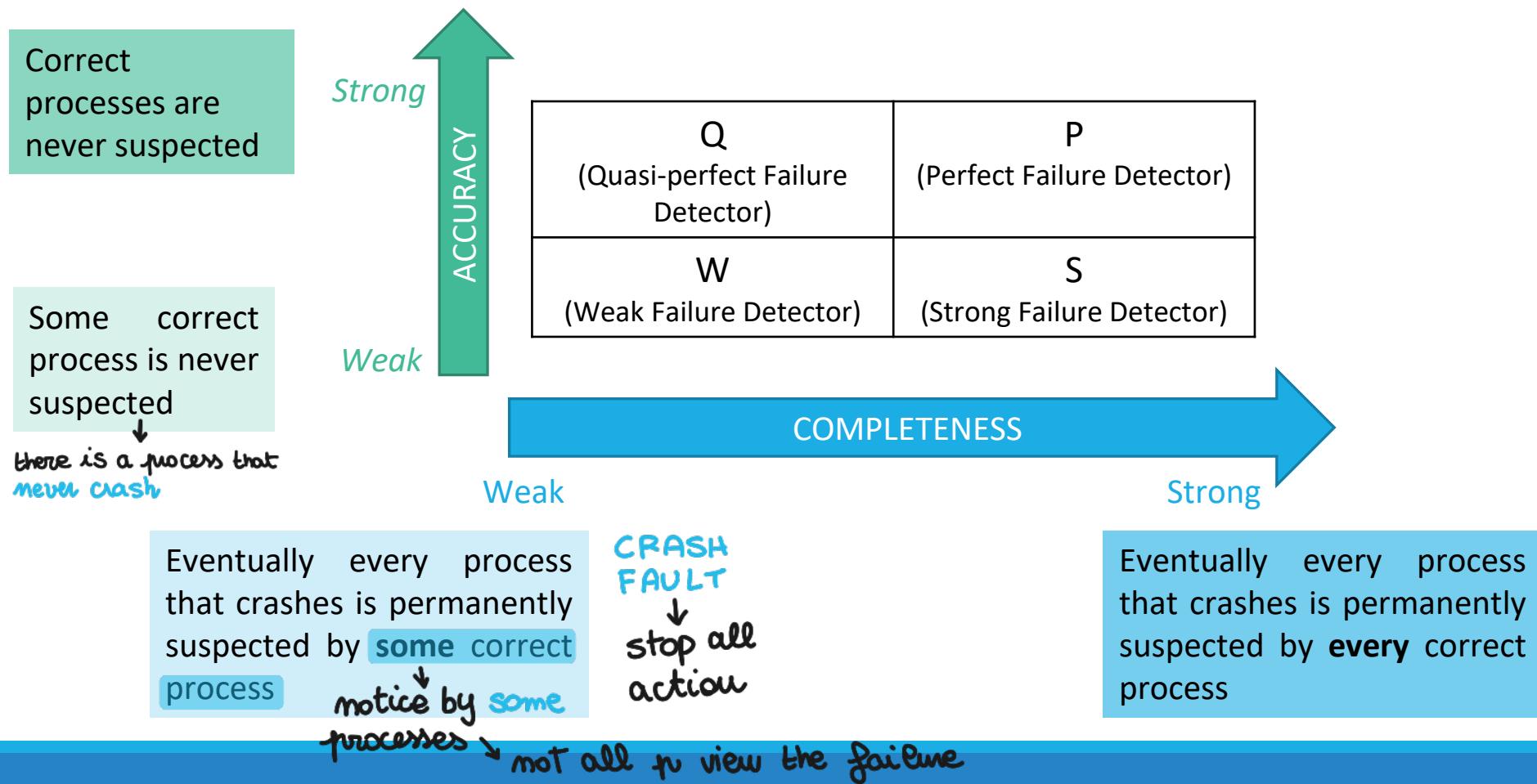
- Accuracy (informally, it is the ability to avoid mistakes in the detection)
- Completeness (informally, it is ability to detect all failures)



ACCURACY : avoid mistakes in FD

COMPLETENESS : detect all failures

Failure Detectors Classification



Failure Detectors Classification

OBSERVATION

- W guarantees that there is at least one correct process that is never suspected.
- Practically speaking, this could be difficult to achieve and thus it is worth to consider accuracy properties to hold *eventually*

{ EVENTUAL STRONG ACCURACY

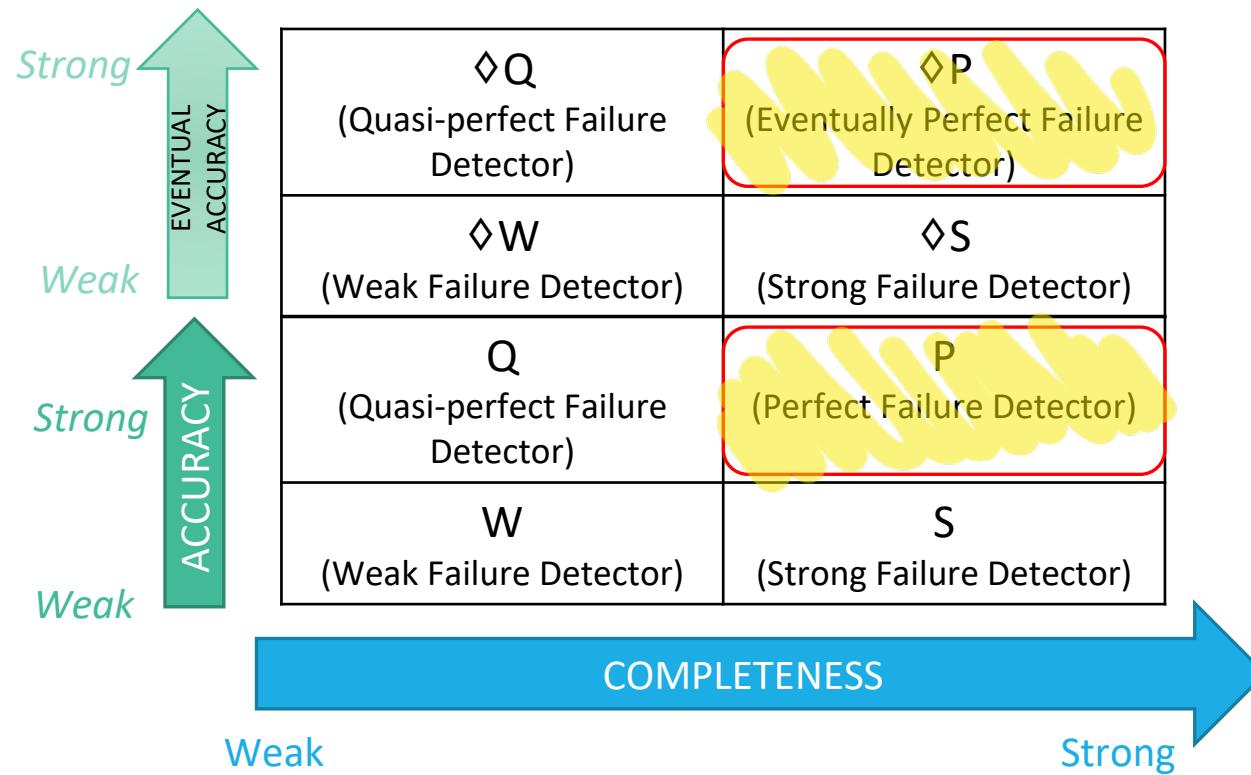
- There is a time after which correct processes are not suspected

EVENTUALLY

{ EVENTUAL WEAK ACCURACY

- There is a time after which some correct process is not suspected

Failure Detectors Classification



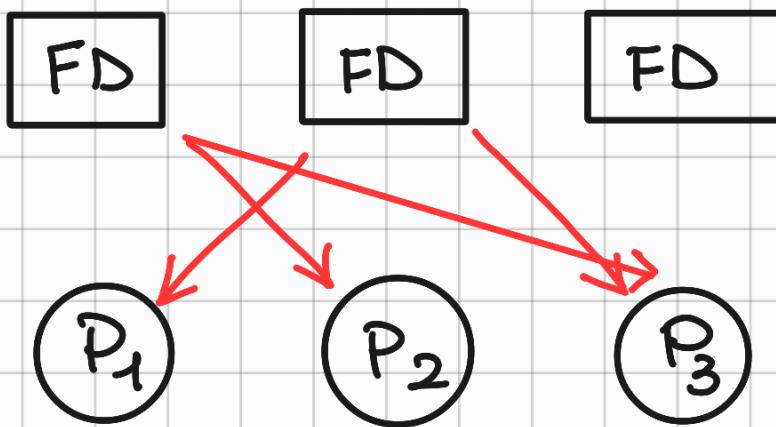
Perfect Failure detectors (P)

System model

- synchronous system
- crash failures

based on behaviour: WE KNOW THE TOTAL EXECUTION
some processes can stop their execution

Using its own clock and the bounds of the synchrony model, a process can infer if another process has crashed



Detect a crash: send messages to all processes and check the response if it arrive or not

Perfect failure detectors (P)

Specification

Module 2.6: Interface and properties of the perfect failure detector

Module:

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

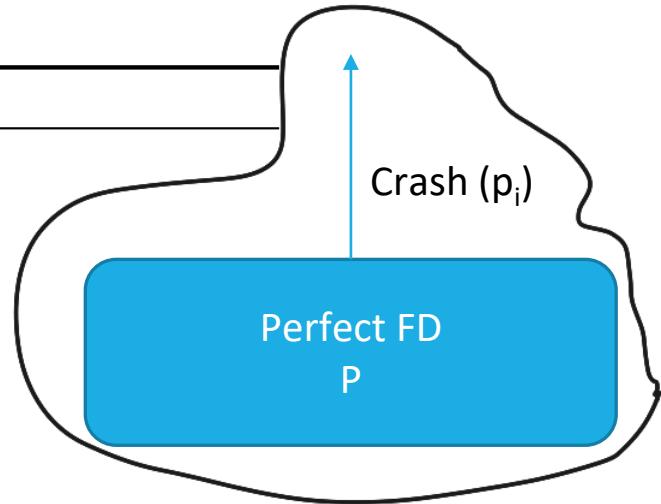
Events:

Indication: $\langle \mathcal{P}, \text{Crash} \mid p \rangle$: Detects that process p has crashed.

Properties:

PFD1: *Strong completeness*: Eventually, every process that crashes is permanently detected by every correct process.
not in a specific

PFD2: *Strong accuracy*: If a process p is detected by any process, then p has crashed.



Perfect failure detectors (P)

Implementation

Algorithm 2.5: Exclude on Timeout

Implements:

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

Uses:

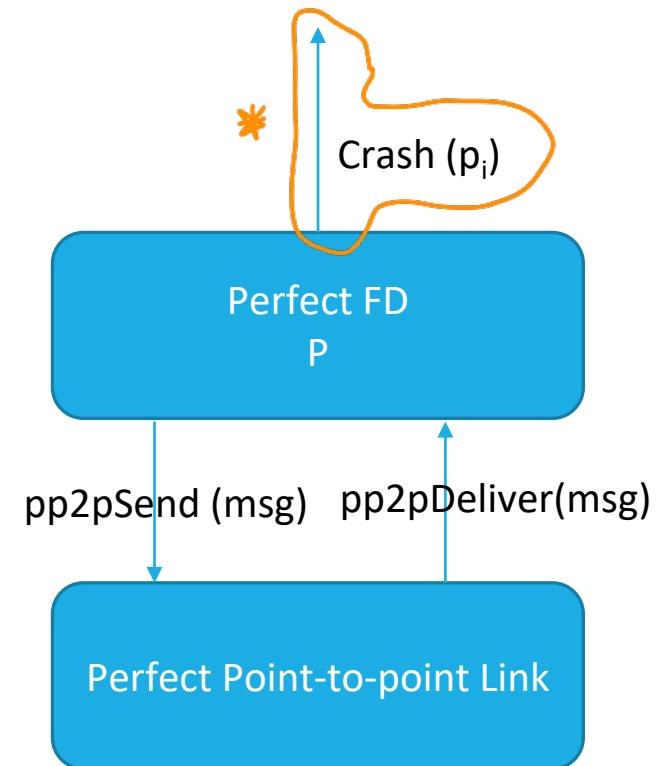
PerfectPointToPointLinks, **instance** pl .

```
upon event <  $\mathcal{P}$ , Init > do
    alive :=  $\Pi$ ; all processes of system
    detected :=  $\emptyset$ ;
    startimer( $\Delta$ );

upon event < Timeout > do
    forall  $p \in \Pi$  do
        if  $(p \notin alive) \wedge (p \notin detected)$  then
            detected := detected  $\cup \{p\}$ ;
            trigger <  $\mathcal{P}$ , Crash |  $p$  >; *
            trigger <  $pl$ , Send |  $p$ , [HEARTBEATREQUEST] >;
            alive :=  $\emptyset$ ;
            startimer( $\Delta$ );

upon event <  $pl$ , Deliver |  $q$ , [HEARTBEATREQUEST] > do
    trigger <  $pl$ , Send |  $q$ , [HEARTBEATREPLY] >

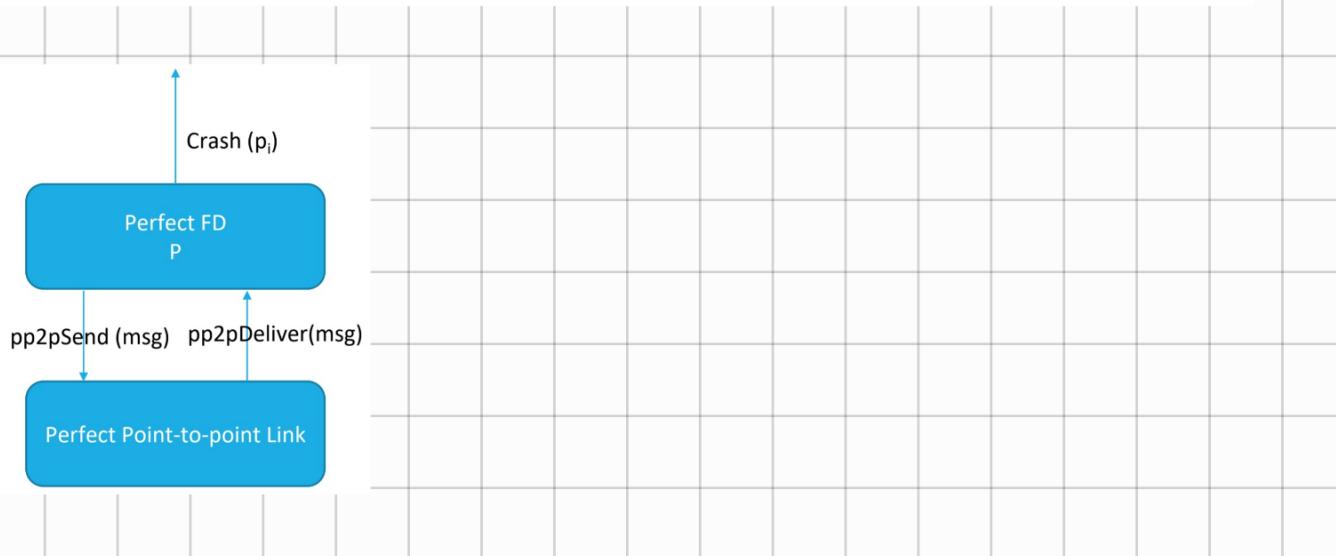
upon event <  $pl$ , Deliver |  $p$ , [HEARTBEATREPLY] > do
    alive := alive  $\cup \{p\}$ ;
```



6: PERFECT FAILURE DETECTOR

PFD1: *Strong completeness:* Eventually, every process that crashes is permanently detected by every correct process.

PFD2: *Strong accuracy:* If a process p is detected by any process, then p has crashed.



Algorithm 2.5: Exclude on Timeout

Implements:

PerfectFailureDetector, instance \mathcal{P} .

Uses:

PerfectPointToPointLinks, instance pl .

```
upon event <  $\mathcal{P}$ , Init > do
    alive :=  $\Pi$ ; all processes of system
    detected :=  $\emptyset$ ; is empty
    starttimer( $\Delta$ );

upon event < Timeout > do
    forall  $p \in \Pi$  do
        if (p  $\notin$  alive)  $\wedge$  (p  $\notin$  detected) then
            detected := detected  $\cup$  {p};
            trigger <  $\mathcal{P}$ , Crash | p >;
            trigger < pl, Send | p, [HEARTBEATREQUEST] >;  $\rightarrow$  PING
            trigger < pl, Send | p, [HEARTBEATREQUEST] >; : pi, are you alive?
            alive :=  $\emptyset$ ; every p is stopped
            ...
            starttimer( $\Delta$ );

upon event < pl, Deliver | q, [HEARTBEATREQUEST] > do
    trigger < pl, Send | q, [HEARTBEATREPLY] >; : pi says: I'm alive
    } if pi CRASH
    } this part doesn't
    } exist

upon event < pl, Deliver | p, [HEARTBEATREPLY] > do
    alive := alive  $\cup$  {p}; update "alive"
```

TIMEOUT: is equal to RTT + time processing
(generate request,)
...



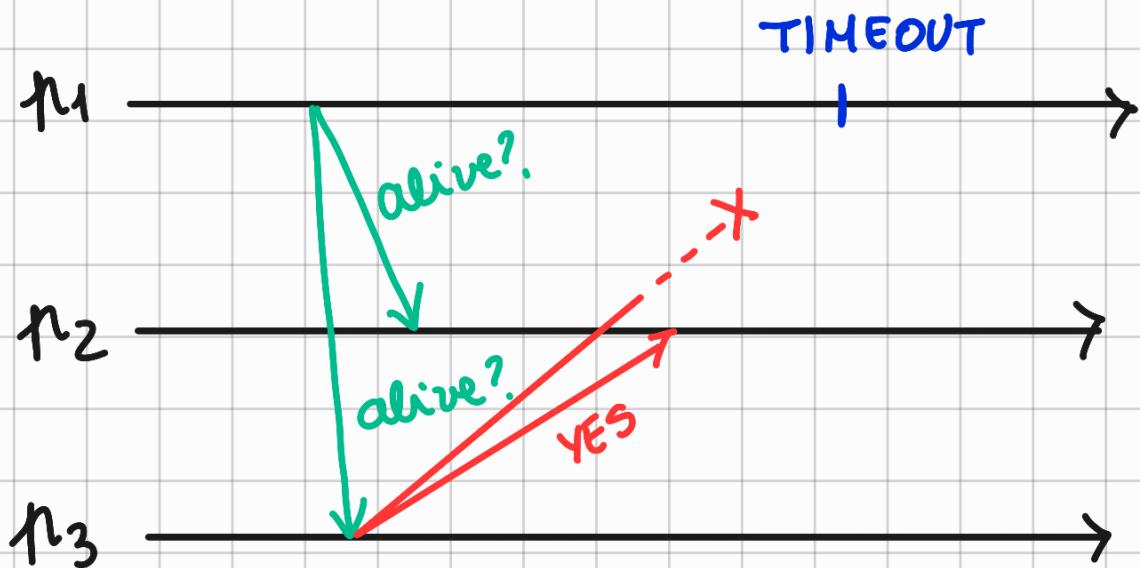
Correctness

- To prove the correctness, we must prove that both *Strong Completeness* and *Strong Accuracy* are satisfied
 - What if we select a timeout too long?
 - What if we select a timeout too short?
 - What if links are fair loss? *m can be lost* *you can invalidate the accuracy* *a process can reply but the m can be lost → CRASHED* *TARGET AS*
- a)** ➤ If a process p_c is in *detected* of a process p_i at time t , is p_c alive at time t ?
- b)** ➤ If a process p_c is in *detected* of a process p_i at time t , is p_c in *detected* of each process p_j at time t ?

a)

DETECTED $\{p_1\}$

|
t



After 2 TIMEOUT all the processes know
the failure of a specific process

Eventually perfect failure detectors ($\diamond P$)

System model

- partial synchrony
- crash failures
- perfect point-to-point links

Crashes can be accurately detected only after a (unknown) time t

- Before time t the system behaves as an asynchronous one
- The failure detector may make mistakes before time t considering correct processes as crashed.
- The notion of detection becomes *suspicious*

Eventually perfect failure detectors ($\diamond P$) Specification

Module 2.8: Interface and properties of the eventually perfect failure detector

Module:

Name: EventuallyPerfectFailureDetector, **instance** $\diamond P$.

Events:

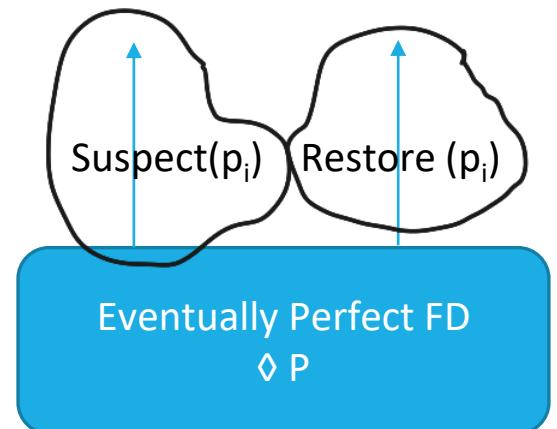
Indication: $\langle \diamond P, \text{Suspect} | p \rangle$: Notifies that process p is suspected to have crashed.

Indication: $\langle \diamond P, \text{Restore} | p \rangle$: Notifies that process p is not suspected anymore.

Properties:

EPFD1: *Strong completeness*: Eventually, every process that crashes is permanently suspected by every correct process.
we don't know when

EPFD2: *Eventual strong accuracy*: Eventually, no correct process is suspected by any correct process.



Basic constructions rules of an eventually perfect FD

Use timeouts to suspect processes that did not sent expected messages

A suspect may be wrong

- A process p_i may suspect another one p_j as *the current timeout is too short*

$\Diamond P$ is ready to change its judgment as soon as it receives a message from p_j

- In this case, *the timeout value is updated*

If p_j has actually crashed, p_i does not change its judgment anymore

Eventually perfect failure detectors ($\diamond P$) Implementation

Algorithm 2.7: Increasing Timeout

Implements:

EventuallyPerfectFailureDetector, **instance** $\diamond P$.

Uses:

PerfectPointToPointLinks, **instance** pl .

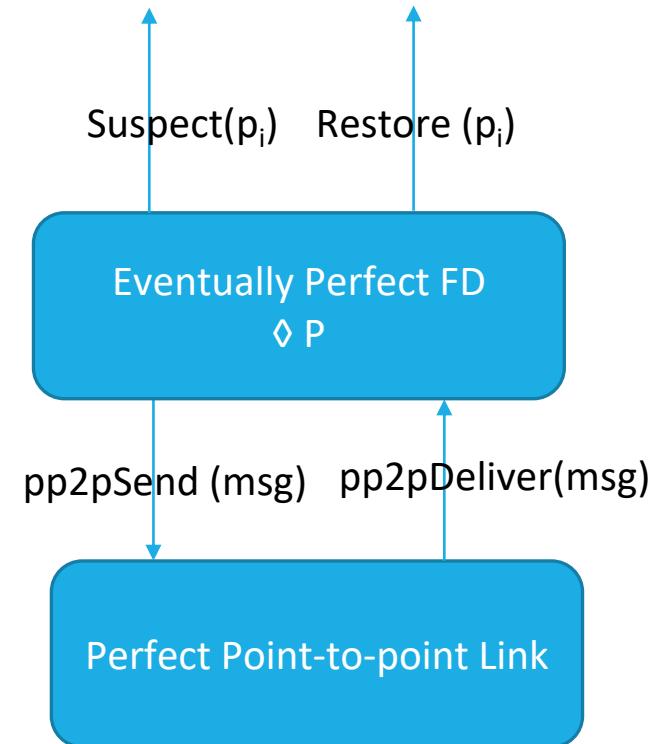
```
upon event <  $\diamond P$ , Init > do
    alive :=  $\Pi$ ;
    suspected :=  $\emptyset$ ;
    delay :=  $\Delta$ ;
    starttimer(delay);

upon event < Timeout > do
    if alive  $\cap$  suspected  $\neq \emptyset$  then
        delay := delay +  $\Delta$ ;
    forall p  $\in$   $\Pi$  do
        if (p  $\notin$  alive)  $\wedge$  (p  $\notin$  suspected) then
            suspected := suspected  $\cup$  {p};
            trigger <  $\diamond P$ , Suspect | p >;
        else if (p  $\in$  alive)  $\wedge$  (p  $\in$  suspected) then
            suspected := suspected  $\setminus$  {p};
            trigger <  $\diamond P$ , Restore | p >;
        trigger < pl, Send | p, [HEARTBEATREQUEST] >;
    alive :=  $\emptyset$ ;
    starttimer(delay);
```



```
upon event < pl, Deliver | q, [HEARTBEATREQUEST] > do
    trigger < pl, Send | q, [HEARTBEATREPLY] >

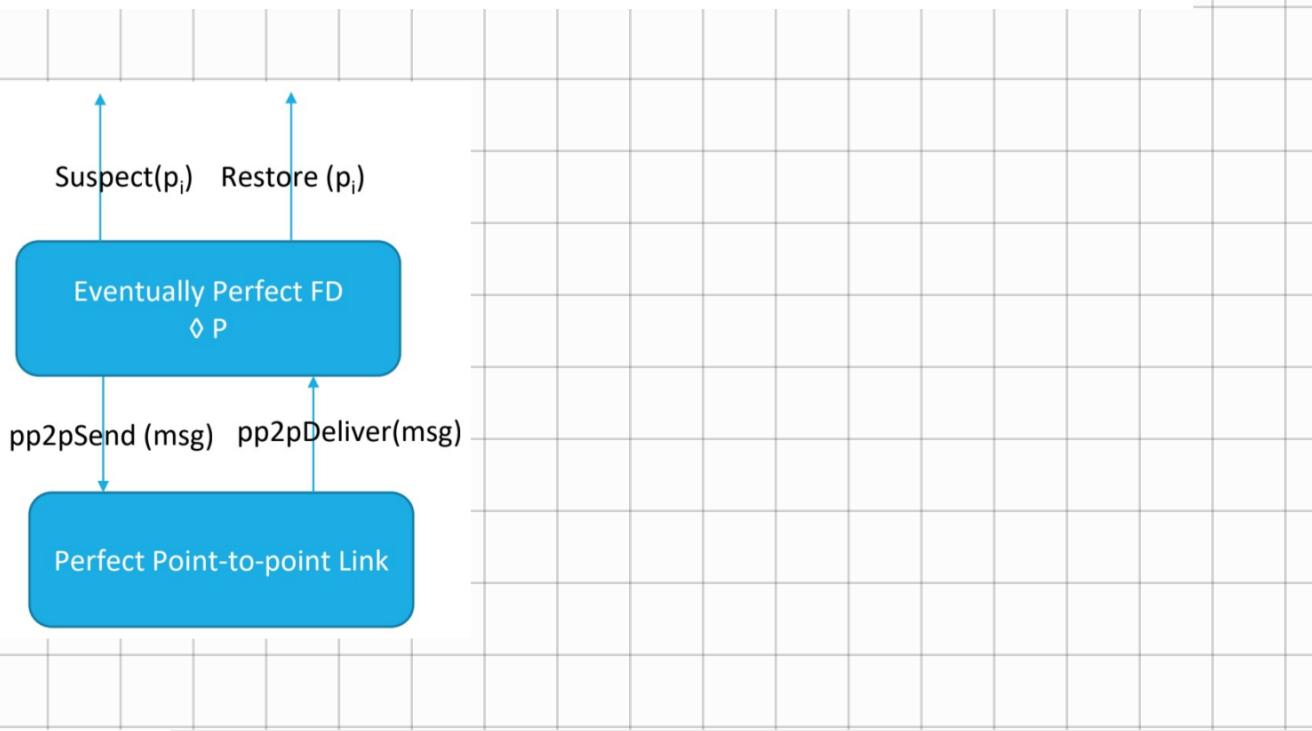
upon event < pl, Deliver | p, [HEARTBEATREPLY] > do
    alive := alive  $\cup$  {p};
```



6: EVENTUALLY PFD

EPFD1: *Strong completeness:* Eventually, every process that crashes is permanently suspected by every correct process.

EPFD2: *Eventual strong accuracy:* Eventually, no correct process is suspected by any correct process.



Algorithm 2.7: Increasing Timeout

Implements:

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

Uses:

PerfectPointToPointLinks, **instance** pl .

```

upon event <  $\diamond \mathcal{P}$ , Init > do
  alive :=  $\Pi$ ; all processes
  suspected :=  $\emptyset$ ;
  delay :=  $\Delta$ ;
  starttimer(delay);

upon event < Timeout > do
  if alive  $\cap$  suspected  $\neq \emptyset$  then
    delay := delay +  $\Delta$ ;
    forall p  $\in$   $\Pi$  do
      when TIMEOUT expires
        if (p  $\notin$  alive)  $\wedge$  (p  $\notin$  suspected) then
          suspected := suspected  $\cup$  {p};
          trigger <  $\diamond \mathcal{P}$ , Suspect | p >;
          delay or crash?
        else if (p  $\in$  alive)  $\wedge$  (p  $\in$  suspected) then
          suspected := suspected  $\setminus$  {p}; remove
          trigger <  $\diamond \mathcal{P}$ , Restore | p >;
        end
      end
    end
    trigger < pl, Send | p, [HEARTBEATREQUEST] >;
    alive :=  $\emptyset$ ;
    starttimer(delay);
  end
end

upon event < pl, Deliver | q, [HEARTBEATREQUEST] > do
  trigger < pl, Send | q, [HEARTBEATREPLY] >

upon event < pl, Deliver | p, [HEARTBEATREPLY] > do
  alive := alive  $\cup$  {p};
  I'M ALIVE
end
  
```

Annotations:

- To avoid the multiple SUSPECT and RESTORE I make the TIMEOUT longer*
- ARE U ALIVE?*
- WHEN a PROCESS is in SUSPECTED but then REPLAY so it is ALIVE*

Correctness

Strong completeness

- If a process crashes, it will stop to send messages. Therefore the process will be suspected by any correct process and no process will revise the judgement.

Eventual strong accuracy

- After time T the system becomes synchronous. i.e., after that time a message sent by a correct process p to another one q will be delivered within a bounded time. If p was wrongly suspected by q, then q will revise its suspicious.

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

T. Chandra, S. Toueg *Unreliable Failure Detectors for Reliable Distributed Systems*

- <https://dl.acm.org/doi/pdf/10.1145/226643.226647>