

# Distributed Systems

Raffaele Castagna

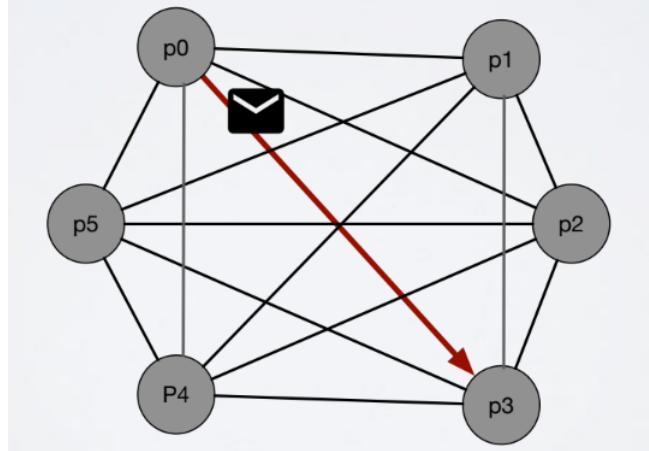
Academic Year 2025-2026

## Indice

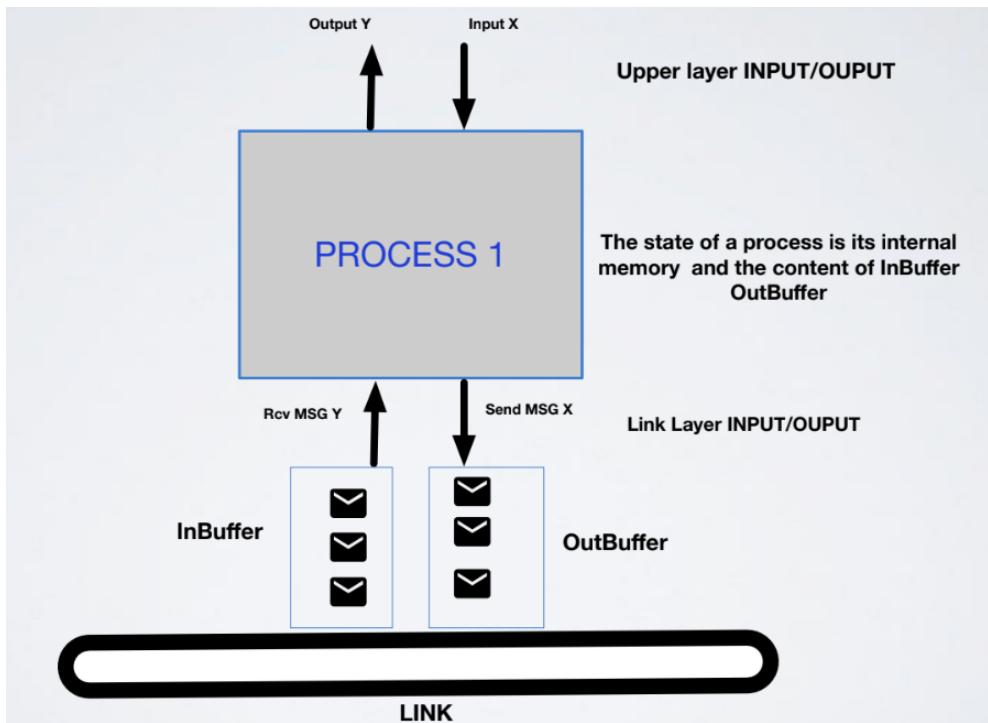
|          |  |          |
|----------|--|----------|
| <b>1</b> | <b>Introduction</b>                            | <b>2</b> |
| 1.1      | Synchronous Vs Asynchronous . . . . .          | 4        |
| 1.2      | Failures . . . . .                             | 4        |
| 1.3      | Different types of links . . . . .             | 5        |
| <b>2</b> | <b>Time in Distributed Systems</b>             | <b>8</b> |
| 2.1      | Clock Synchronization . . . . .                | 8        |
| 2.2      | Algorithms for synchronization . . . . .       | 10       |
| 2.3      | Logical Time . . . . .                         | 12       |
| 2.4      | Lamport's Mutual Exclusion Algorithm . . . . . | 16       |
| 2.5      | Failure Detectors . . . . .                    | 20       |

# 1 Introduction

**Definition 1.** In a system we have  $n$  processes in  $\Pi : p_0 \dots, p_{n-1}$  each with a distinct identity they communicate by utilizing a communication graph  $G : (\Pi, E)$ , the communication is done by exchanging messages.



**Definition 2.** A process is a (possibly infinite) State Machine (I/O Automaton).



Each process has multiple qualities:

- A set of internal states  $Q$
- A set of initial states  $Q_i \subset Q$
- A set of all possible messages  $M$  in the form  $\langle \text{sender}, \text{receiver}, \text{payload} \rangle$
- Multiset of delivered messages  $InBuf_j$
- Multiset of inflight messages  $OutBuf_j$

We can formally describe this as follows: (this isn't part of the exam btw)

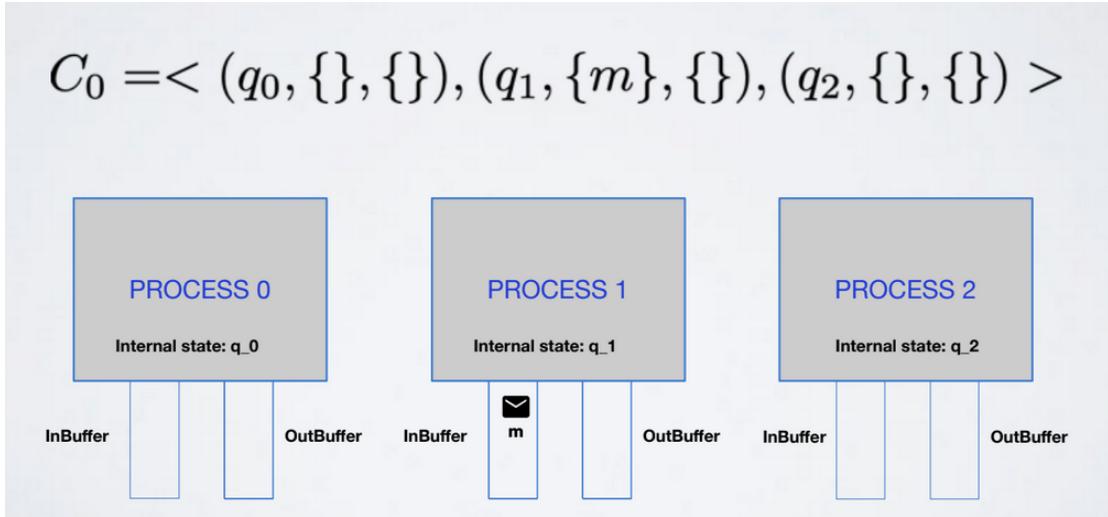
$$P_j(q \in Q \cup Q_{in}, InBuf_j) = (q' \in Q, SendMsg \subset M)$$

$$OutBuf_j = OutBuf_j \cup SendMsg$$

$$InBuf_j = \emptyset$$

To execute a process we have an adversary that schedules a set of events (scheduler), these events may be for example a delivery (e.g.  $Del(m,i,j)$ ) or it can be one step of the step machine of process i ( $Exec(i)$ )

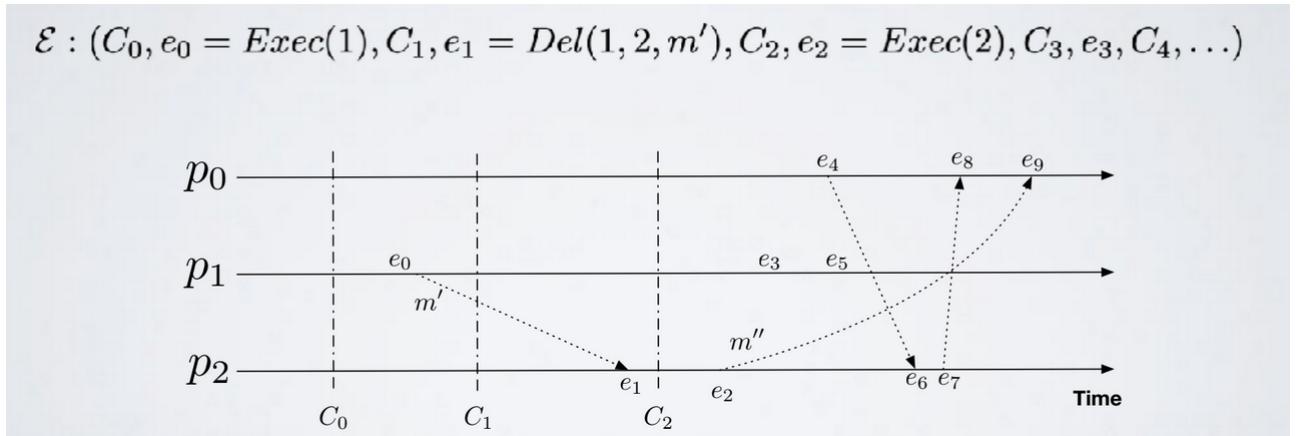
**Definition 3.** A configuration  $C_t$  is a vector of  $n$  components, component  $j$  indicates the state of process  $j$ .



An event is **enabled** in configuration  $c$  if it can happen.

**Definition 4.** An execution is an infinite sequence that alternates configurations and events:  $(C_0, e_0, C_1, e_1, C_2, e_2, \dots)$  such that each event  $e_t$  is enabled in configuration  $C_t$  and  $C_t$  is obtained by applying  $e_{t-1}$  to  $C_{t-1}$

It may be useful to visualize how an execution involving multiple processes works, here we have an example:



**Definition 5.** A *fair execution* is an execution  $E$  where each process  $p_i$  executes an infinite number of local computations ( $\text{Exec}(i)$  events are not finite) and each message  $m$  is eventually delivered (we can't stall messages)

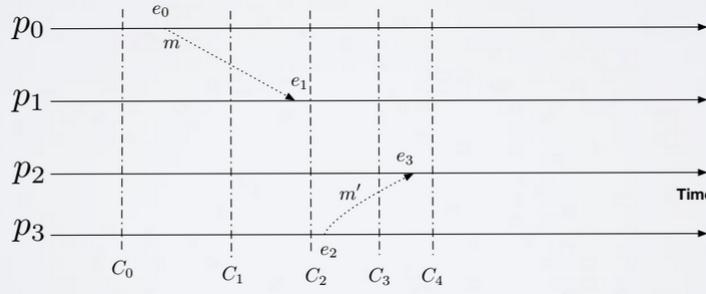
We will always use fair executions unless stated otherwise.

**Definition 6.** Given an execution  $E$  and a process  $p_j$ , we define the local view/ local execution of  $E|p_j$  the subset of events in  $E$  that impact  $p_j$

$$\mathcal{E} = (C_0, e_0 = \text{Exec}(0), C_1, e_1 = \text{Del}(0, 1, m), C_2, e_2 = \text{Exec}(3), C_3, e_3 = \text{Del}(3, 2, m'), \dots)$$

$$\mathcal{E}|p_1 = (\text{Del}(0, 1, m), \dots)$$

$$\mathcal{E}|p_2 = (\text{Del}(3, 2, m'), \dots)$$



But these executions do not account for time, so we may have executions that are the same even though the events happened at different times, in case this does happen, we say that two executions are *indistinguishable*.

**Theorem 1.** In the asynch. model there is no distributed algorithm capable of reconstructing the system execution.

## 1.1 Synchronous Vs Asynchronous

We have 3 main types of synchrony:

- Asynchronous Systems
- Eventually Synchronous Systems
- Synchronous systems

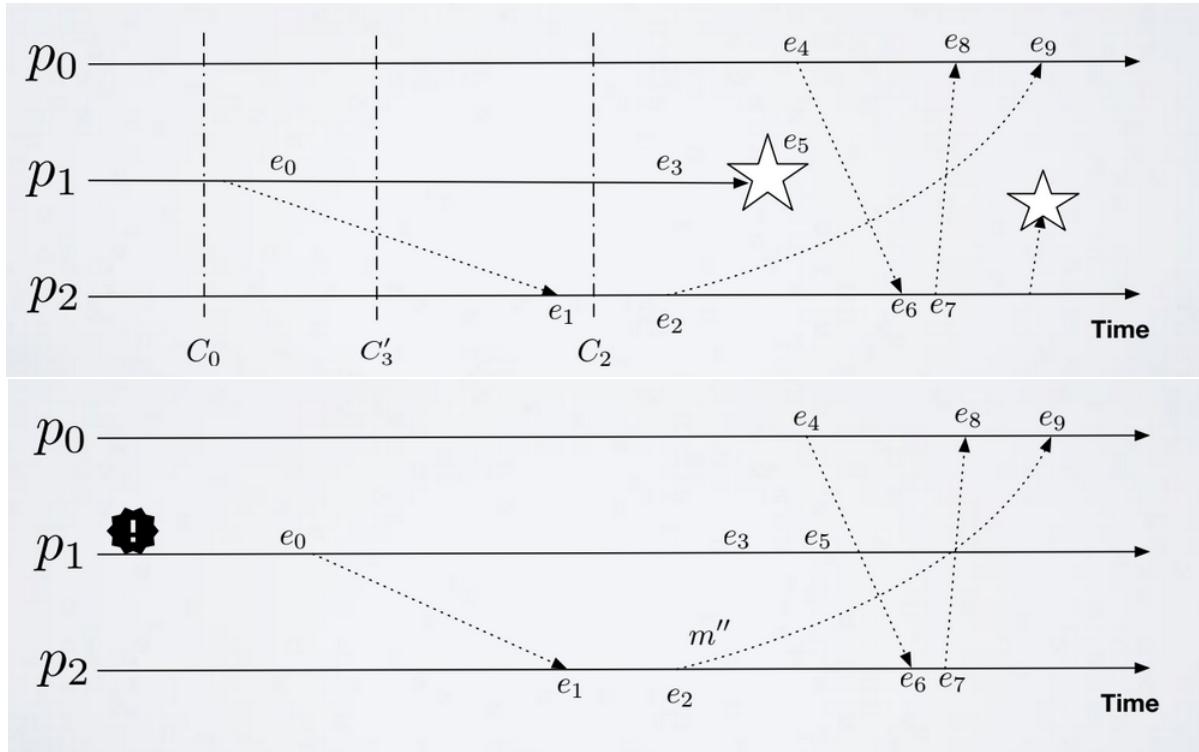
We can say that a system is synchronous if it has a fixed bound on the delay of messages, on the time of actions executed by processes, and a fixed bound between execution of actions.

## 1.2 Failures

We have 2 main models for failures:

- Crash-stop Failures (The program crashes, and doesn't respond)
- Byzantine Failures (The behaviour of the program is random)

We signal crash failures with a star sign and byzantine failures with a !



Byzantine failures are a superset of Crashstop failures, so algorithms that will work on byzantine failures will always work on crash stop failures, but not the contrary.

A process is **correct** if it does not experience a failure, every algorithm has a maximum number  $f$  of failures that it can experience.

### 1.3 Different types of links

**Definition 7.** An **Abstraction** is the formalization of a problem/object, to built one we must define the system model and formalize a problem/object so that there is no ambiguity regarding the properties of our abstraction

For example let us abstract a link: it is something that may lose a message with a certain probability  $pr$ , the messages can be duplicated a finite number of times and they must come from somewhere. Inside a link we have two main events:

- **Requests:**  $\langle \text{Send} — q, m \rangle$  sends message  $m$  to process  $q$
- **Indication:**  $\langle \text{Deliver} — p, m \rangle$  delivers a message  $m$  from process  $p$  (this might just be an identifier, e.g. an ip address or mac address)

Now let us formalize this further via its properties:

- **FL1:** (Fair loss) If a correct process  $p$  sends infinitely often  $m$ , a process  $q$  then delivers  $m$  an infinite number of times. (e.g. suppose we have  $\frac{1}{2}$  probability and we send it over 10 times, then the probability will be  $1 - \frac{1}{2^{10}}$  as events are independentm aka if we try hard enough we get the message)
- **FL2:** (Finite duplication) if a correct process  $p$  send  $m$  a finite number of times to  $q$ , then  $q$  cannot deliver  $m$  an infinite number of times (we'll receive a finite number of duplicated packets)
- **FL3:** (No creation) If a certain process  $q$  sends a message  $m$  with  $\text{send}(p)$ ,then  $m$  was sent by  $p$  to  $q$

Our objective is hiding the probability behind infinity.

A link that respects these properties is called a Fair-lossy link, and it's always behind 2 process  $p$  and  $q$ . We can broadly categorize these properties into two classes:

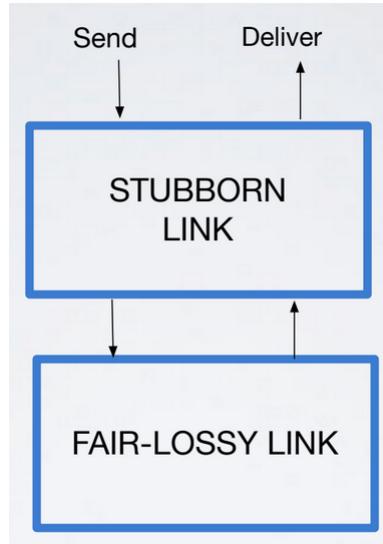
- **Safety:** if a property is violated at time  $t$ , then it cannot be satisfied after that time  $t$ . So if in an execution  $E$  we violated a safety property, then there is a prefix  $E'$  of  $E$  such that any extension  $E'$  also violates the property, an example of safety property is: if we die, we cannot resurrect.
- **Liveness:** these kind of properties cannot be violated in a finite execution, more formally, given an execution  $E$  that does not satisfy a liveness property, there is an extension of  $E$  that satisfies it, informally it just says that something good will eventually happen.

If for example we created a bound on  $FL_2$ , then we have a safety property, as it cannot be extended, as a rule of thumb, if the property is infinite, then it is a liveness property.

Of course there exist badly written properties that try and write both types into a rule, however we should always decompose them (e.g  $q$  will eventually deliver and the delivery is unique, if we decompose it we have "  $q$  eventually delivers,  $m$  is delivered at most once")

We also have what we call **stubborn links**, which inherit  $FL_3$  and add:

- **SL1:** (Stubborn delivery) if a correct process  $p$  sends  $m$  to  $q$ , then  $q$  delivers  $m$  an infinite number of times, hence stubborn.



Our algorithms will always reflect a reactive computing model utilizing handles that consume events or create them, they will always be atomic unless stated otherwise.

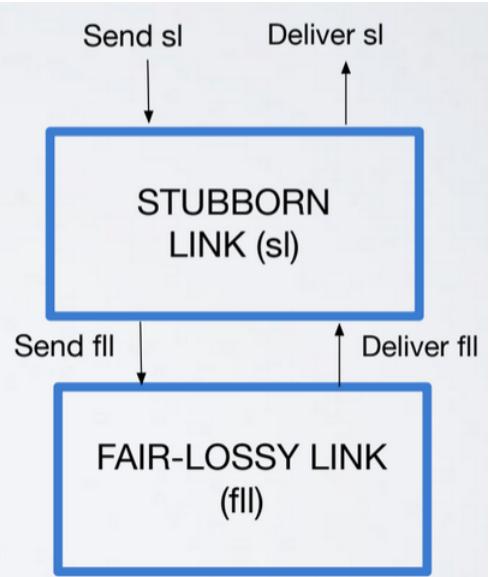
```

upon event <  $sl, Init$  > do
   $sent := \emptyset;$ 
   $starttimer(\Delta);$ 

upon event <  $Timeout$  > do
  forall  $(q, m) \in sent$  do
    trigger <  $fll, Send | q, m$  >;
     $starttimer(\Delta);$ 

upon event <  $sl, Send | q, m$  > do
  trigger <  $fll, Send | q, m$  >;
   $sent := sent \cup \{(q, m)\};$ 

upon event <  $fll, Deliver | p, m$  > do
  trigger <  $sl, Deliver | p, m$  >;
  
```



The initialization event creates a set  $sent$  containing the messages that were sent and then it starts a local timer of delta time (delta is whatever we want), we must remember that this timer is not a global

clock, but local for that process.

In an sl link whenever we get an input we trigger an event to send the message to our fl link, after that we just add it to our set sent.

When the fl wants to send a message to a process it will trigger an sl deliver event.

When a timeout event happens we scan all the messages in the sent set and we send them again, after that we start a timer.

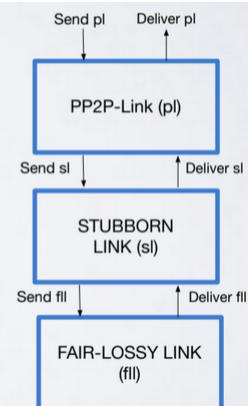
Then we have our **Perfect P2P** links which inherit FL3 but add the following properties:

- **PL1:** Reliable delivery, if a correct process p sends m to q, then q eventually delivers m
- **PL2:** No duplication, a message is delivered at most once.

```
upon event <pl, Init> do
    delivered := ∅;
```

```
upon event <pl, Send | q, m> do
    trigger <sl, Send | q, m>;
```

```
upon event <sl, Deliver | p, m> do
    if m ∉ delivered then
        delivered := delivered ∪ {m};
    trigger <pl, Deliver | p, m>;
```



We initialize it with a delivered empty set, when we want to send a message we create an event and trigger a send event on the sl.

When the sl gets a deliver event we check if m is not delivered, else add it to delivered set, afer that we trigger a pl deliver event.

### Proof of no Duplication P2P

Suppose we send m once, and you receive it twice, the delivery action of a message is guarded by “if  $m \in \text{delivered}$ ”, still suppose we deliver it twice at  $t'$  and  $t$  ( $t < t'$ ), however since the handler is atomic, we have that the set of delivered obtains message m, and therefore when m at  $t'$  is checked at time  $t'$  it is already in delivered, therefore it cannot be delivered again, as this contradicts the fact that trigger  $\langle \text{pl}, \text{deliver} | P, m \rangle$  is executed at (or after) time  $t'$ .

### PL1 Proof

Suppose p sends m and q does not deliver it. There could be two reasons for q to not deliver:

Reason 1: You receive the message from the stubborn link and you drop it.

Reason 2: You don't receive a message m.

If q delivers a message then we execute the handler, the only way to not trigger  $\langle \text{pl}, \text{Deliver} \rangle$  is if it's the if  $m \in \text{delivered}$ , but someone must have delivered it already.

For the second reason we don't get a message at all, but the stubborn link has in turn properties that it can't violate, therefore it is impossible, since the delivery handler would never be triggered.

### Exercise 1

Show that our stubborn algorithm does not work if we change first property to:

*SL(1) If a process p sends m to q, then q delivers m an infinite number of times*

What's missing is the word “correct” therefore the process may crash, so our algorithm with a timer does not work.

Suppose we want to implement this then:

If a process p sends a message m to q, then q delivers m 1 time. Suppose we also have a perfect channel that we need to deliver message m to, but since the process is still not correct, then our link may receive the message and crash, or even before it receives it, it crashes. Therefore it can't be done.

## 2 Time in Distributed Systems

We have different type of "times" in DS:

- **Asynchronous:** No Global-clock, no bound, local-computational steps (Exec) happen at unpredictable time (the adversary e.g. the scheduler), models everything but has a lot of problems
- **Synchronous:** Bound, you can synchronize up to a certain precision, local execution steps happen at certain predetermined interval, and they take a bounded time, models only networks but everything is possible
- **Eventually Synchronous:** we have a threshold, where if its under its in sync, else its async. In an eventually synchronous a safety property does not work if it depends on the delay of the channel.

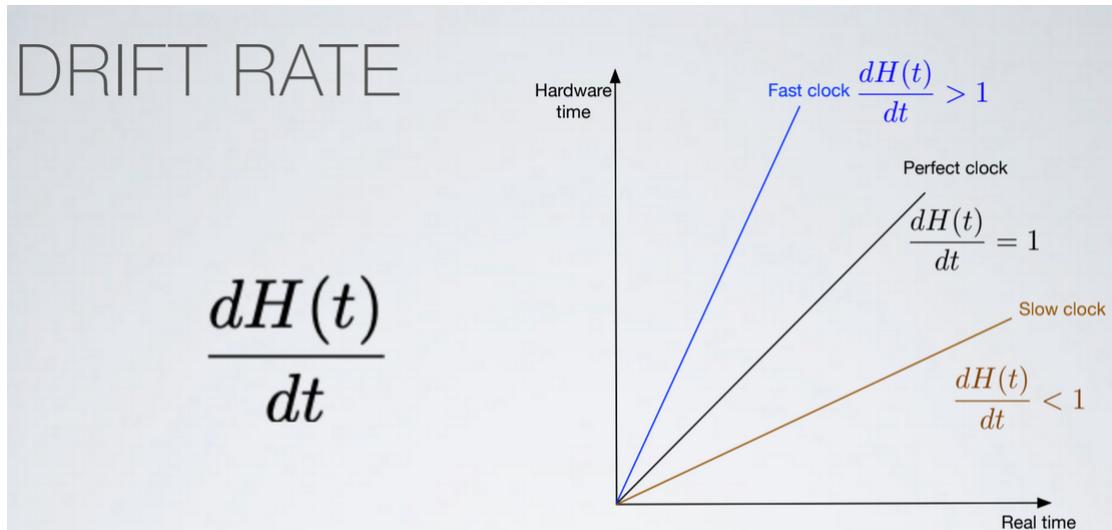
### 2.1 Clock Synchronization

To actually work most applications need an order of actions and synchronization.

To synchronize we need a clock, this clock on most computers works by measuring the oscillation and a counting register that is incremented at every tick of the oscillator, at certain times the OS reads the hardware clock  $H_i(t)$  and produces the local software clock  $C_i(t) = \alpha * H_i(t) + \beta$ . The hardware clock is:

$$H_i(t) = \int_0^t h_i(\tau) d\tau$$

Between two clocks we may have a **skew**, which is the  $|C_i(t) - C_j(t)|$ , we may also have a **drift rate** which is the gradual misalignment of synchronized clocks caused by slight inaccuracies of time-keeping mechanisms, more formally it is the derivate of the hardware clock over the derivative of the time.



To synchronize two clocks we have 2 strategies, either bring one in the future or bring it in the past by modifying beta, we never go into the past as we may have done something before the sync and it would cause confusion. We can however slow the clock that is in the future.

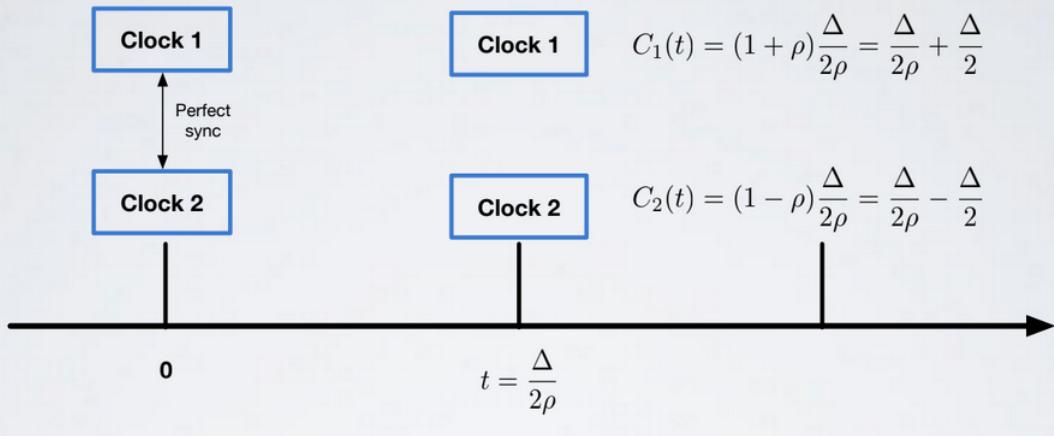
# CORRECT CLOCK (1/2)

An hardware clock  $H$  is correct if its drift rate is within a limited bound of  $\rho > 0$  (e.g.  $10^{-5}$  secs/sec).

$$1 - \rho \leq \frac{dH(t)}{dt} \leq 1 + \rho$$

In presence of a correct hardware clock  $H$  we can measure a time interval  $[t, t']$  (for all  $t' > t$ ) introducing only limited errors.

$$(1 - \rho)(t_1 - t_0) \leq H(t_1) - H(t_0) \leq (1 + \rho)(t_1 - t_0)$$

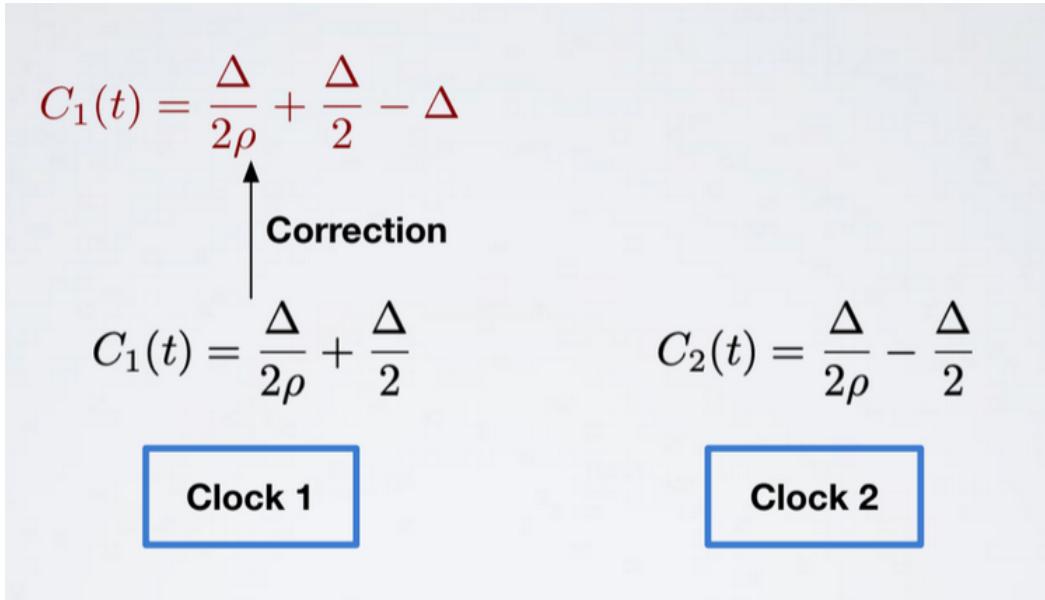


$$C_1(t) = \frac{\Delta}{2\rho} + \frac{\Delta}{2}$$

$$C_2(t) = \frac{\Delta}{2\rho} - \frac{\Delta}{2}$$

**Clock 1**

**Clock 2**



All software clocks have to be monotone:

$$t' > t \implies C(t') > C(t)$$

So if real time  $t'$  is greater than  $t$ , we have that the software clocks reflect that, to actually reflect this, we cannot choose a value  $\beta$  in the negatives, but rather a value  $\alpha$  that is  $< 1$  by *masking oscillations* as well changing the value  $\beta$ .

When we synchronize with external time, we synchronize with UTC as it is the international standard and we'll synchronize via satellites.

We can define **external synchronization** when we synchro with UTC, so each process is synchronized with an authoritative external source, this means that the difference between any computer and the external source is below a certain bound D.

If we don't synchro with the external world we are doing **Internal synchronization**, a set of processes is internally synchronized if the difference between their clocks is less than D. If we have external we also have internal.

## 2.2 Algorithms for synchronization

For external synchronization we have christian's algorithm, we utilize a Server S that receives a signal from an UTC source, works (probabilistically) in async systems, its based on RTTs, and only if RTTs are small and respect the required accuracy will the response be considered. We can also divide to know the difference between when the server sent the info and when I received it, so just add  $time\_response + \frac{RTT}{2}$  but this assumes that RTT is symmetric. Another problem is the fact that we have a single point of failure (server) in real life we just synchronize with multiple servers.

## Case 1

- Reply time is greater than estimate ( $RTT/2$ )

- Assume it is equal to  $RTT - min$

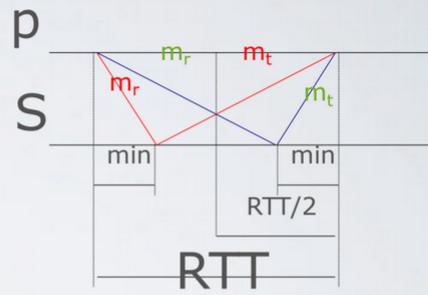
$$\Delta = RTT/2 - (RTT - min) = \frac{-RTT + 2min}{2} = -\frac{RTT}{2} + min = -\left(\frac{RTT}{2} - min\right)$$

## Case 2

- Reply time is smaller than estimate ( $RTT/2$ )
- Assume it is equal to  $min$

$$\Delta = \frac{RTT}{2} - min$$

Accuracy is  $\pm (RTT/2 - min)$

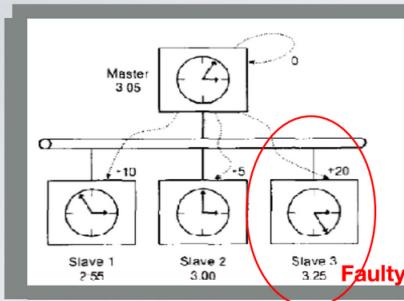


For internal clock synchronization we have the Berkley algorithm, in which we have a master/slave architecture, where the master process  $p_m$  sends a message with timestamp  $t_1$  (local clock) to each process including itself, after that when each process receives a message from the master, it sends back a reply with a timestamp  $t_2$  (local clock value), the master then receives the reply, reads its local clock and computes the difference between the clocks. We will then compute an average value of all non-faulty processes (difference is not more than a certain  $y$ ) and the master will compute Deltas for each computer.

Let's see an example:

We have the master, with a time of 3:05, then we have slave 1 with 2:55, then slave 2 with 3:00 and slave 3 with 3:25. The master asks the time for everyone, receiving the answers. The master will now compute the delta between the master and itself, which is obviously zero, then He will do the same for  $\Delta M_1 = -10$ , then  $\Delta M_2 = -5$  and then  $\Delta M_3 = 20$ . By using those quantities, we'll compute an average difference, but before doing so the masters will remove processes who are likely faulty, so processes who surpass a certain threshold. By setting the threshold, for example, at 10, we will remove  $\Delta M_3$ , because it's  $20 > 10$ . We will now have that  $Avg = 0 - 10 - 5 / 3 = -5$ . We'll tell anyone now that, in order to move to the correct value, the process will have to move to  $\Delta M_i - Avg$ . For example, for the slave 1 he'll have to move to  $-5 - (-10) = +5$ . The slave 2 will have to move of  $-5 - (-5) = 0$ . We compute also for slave 3, by having  $-5 - 20 = -25$ . Also for the master, by having  $-5 - 0 = -5$ . Each process will now be at 3:00.

## EXAMPLE



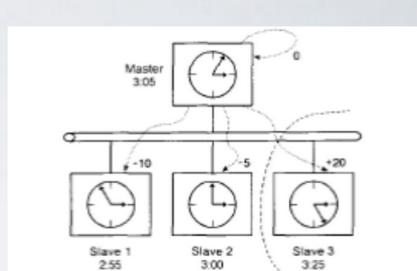
Measuring the differences

$$\Delta p_m = 3:05 - 3:05 = 0$$

$$\Delta p_1 = 2:55 - 3:05 = -10$$

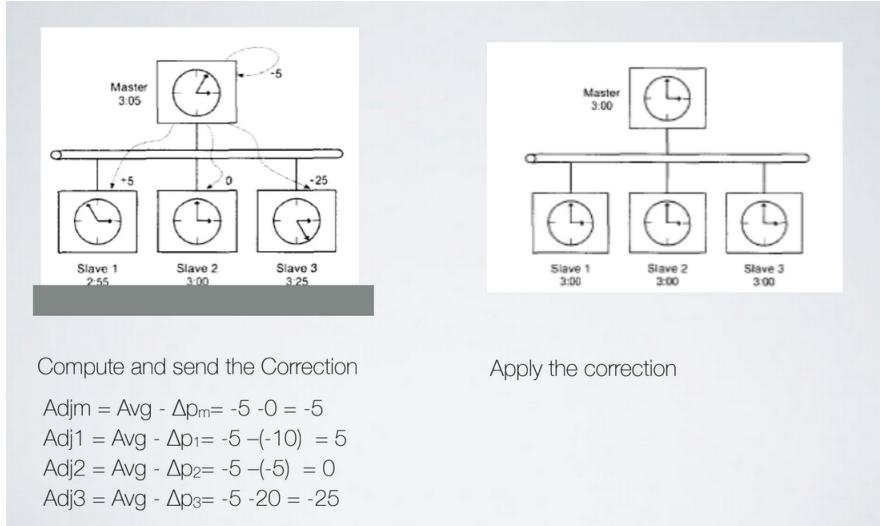
$$\Delta p_2 = 3:00 - 3:05 = -5$$

$$\Delta p_3 = 3:25 - 3:05 = 20$$



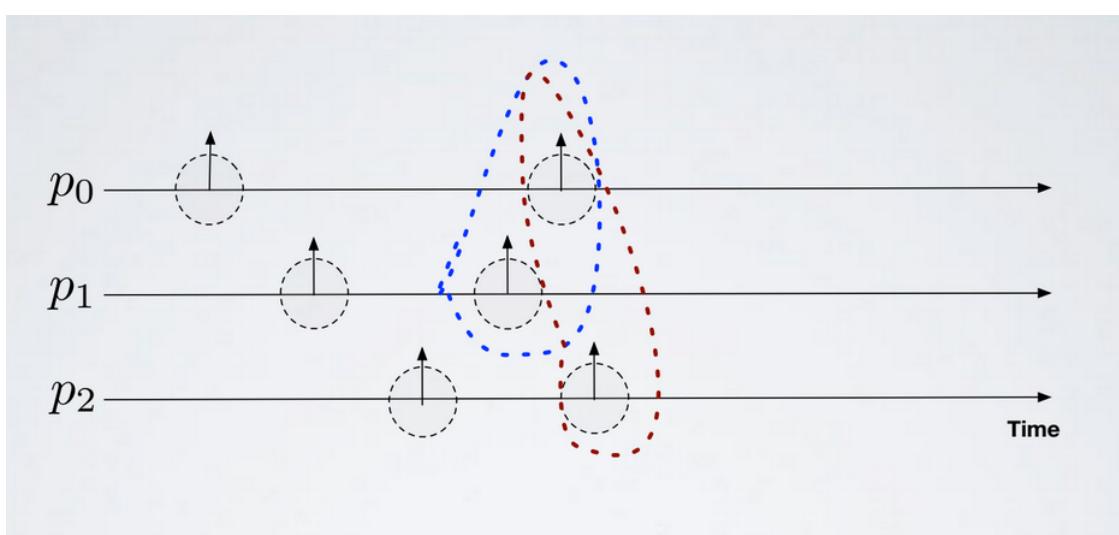
Computing the average

$$Avg = (0 - 10 - 5) / 3 = -5$$



For simplicity here we can go back in time, though in real life we should slow down the clock. The accuracy of Berkley depends on the maximum round trip time, as the master does not consider values associated with RTTs higher than a threshold, this algorithm is also not fault tolerant for the master, if the slaves are faulty the algorithm still works, if the master crashes, eventually a new one will be elected.

In real life we use **NTP** (Network Time Protocol), it synchronizes computers with UTC. The servers are organized in a hierarchy. The server at the top has the clock. The second level (stratum) synchronize directly with server 1, with an algorithm similar to Christian's one. The third stratum only synchronize with stratum 2. We can have more strata, like 4 and so on. Error propagates with the descension of stratum (so stratum 3, for example, will have a cumulative error of  $e_1$  and  $e_2$ ). We do so such that the bandwidth of the server 1 is not completely saturated. Clock synchronization has a limit, we can do it only on synchronous system, we have to have predictable delays in the channel. Even on synchronized systems, we could have a problem with clocks, given by the bounded accuracy. The problems take place when the bound of the error overlaps with the bound of the error of another process, but not guarantee the order.



### 2.3 Logical Time

Logical time is a concept of time we can implement also in asynchronous systems. It takes care of the causal relationship between events. We'll talk about two methods given by Scalar/Lamport's clocks, and then about Vector Clocks, used to measure logical time. We'll also see Lamport's Mutual Exclusion

algorithm.

Let's do an example.

Let's say we want to build a chat application. The easiest thing to do would be using a server-client structure. If we are in a scenario with crash failures, if the server dies the application stops working. Another way is not using a server, i got lost sorry.

Let's say we have a chat application with a specific pattern of messages. We have Alice asking a question, then Mike responds, then Alice asks another question and then another entity answers. This would be the real order of messages.

Because of the fact that the system is asynchronous, a broadcast message might arrive with a certain delay to other entities.

Let's say that Bob, at a certain time  $t$ , only receives the messages from Mike and Bob's supervisor, and not those from Alice. Bob would have a limited local view. To avoid this kind of problem we have to formalize the concept of causality (the messages of Mike, ecc are caused by Alice sending the first message).

The simple FIFO ordering, in this situation, will not help us. Our goal is to find a way to timestamp an event, and to do that we have to formalize the concept of causality, and then we timestamp events.

The definition of causal relationship was introduced by Lamport. We say that an event happens before one another if there is a causal relationship between these two.

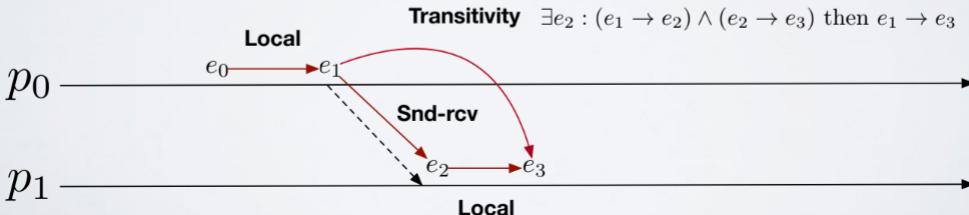
Let's say that I have  $p_0$ , and two events happen on it,  $e_0$  and then  $e_1$ . In this case there is a causal relationship that  $e_0$  happens before  $e_1$  (we write it with an arrow).

With condition B, let's say we have  $e_1$  that's the sending of message  $m$ ,  $e_2$  the receipt. We say that  $e_1$  happens obviously before  $e_2$  (we received a message because someone has sent it).

We can now define a partial ordering between events in the system, we introduce three properties

Two events  $e$  and  $e'$  are related by happened-before relation ( $e \rightarrow e'$ ) if:

- Local ordering:  $\exists p_i | e \rightarrow_i e'$
- Snd-rcv ordering:  $\forall m, send(m) \rightarrow receive(m)$ 
  - $e=send(m)$  is the event of sending a message  $m$
  - $e'=receive(m)$  is the event of receipt of the same message  $m$
- Transitivity:  $\exists e'' : (e \rightarrow e'') \wedge (e'' \rightarrow e') \text{ then } e \rightarrow e'$ 
  - the *happened-before* relation is transitive



The first property is the local ordering, (A) in the slides, and the second receive is (B). For the first, for example, we say that  $e$  and  $e'$  are related if a process  $p_i$  exists by which  $e \rightarrow_i e'$  ( $e$  is related with  $e'$  in  $p_i$ )

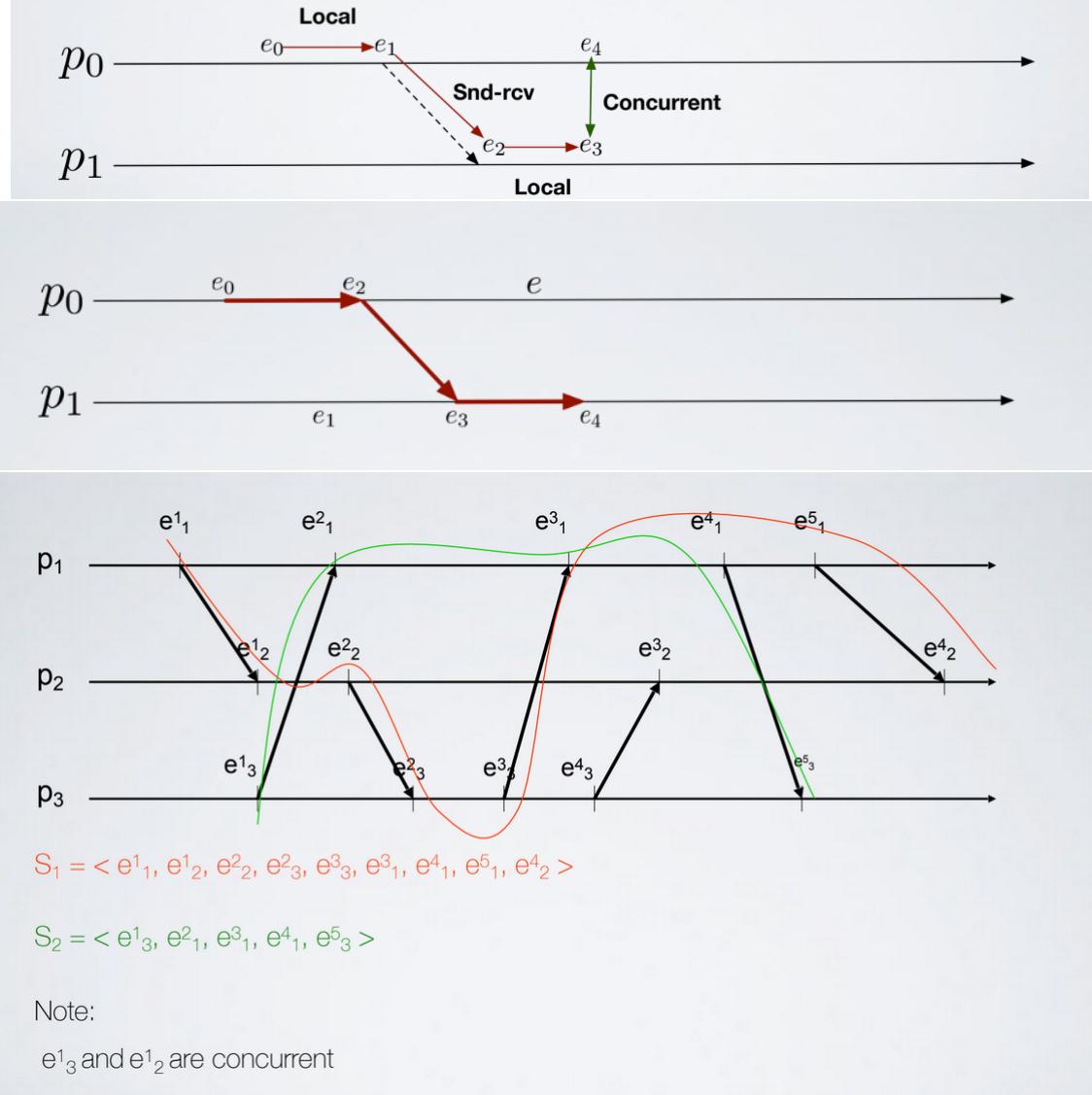
Then we have transitivity. This relationship (in the photo) is telling us that  $e_3$  knows about  $e_0$ .

This relationship is partial, we could have events related by "happened before",

Let's say we have on  $p_0$  another process  $e_4$  happening. We say  $e_3$  and  $e_4$  are concurrent. Two events

are concurrent if they do not know about each other. The only way to have them know is having a message from  $p_0$  to  $p_1$  or viceversa.

- The sequence  $e_1, e_2, \dots, e_n$  may not be unique
- It may exist a couple of events  $\langle e_1, e_2 \rangle$  such that  $e_1$  and  $e_2$  are not in happened-before relation
- If  $e_4$  and  $e_3$  are not in happened-before relation then they are **concurrent** ( $e_4 \parallel e_3$ )
- For any two events  $e_x$  and  $e_y$  in the execution history of a distributed system, either  $e_x \rightarrow e_y$ ,  $e_y \rightarrow e_x$  or  $e_y \parallel e_x$



To see if two events are concurrent or not, we could try to find a path between them.

Let us start with the first tool, the Logical Clock (Lamport's Clock).

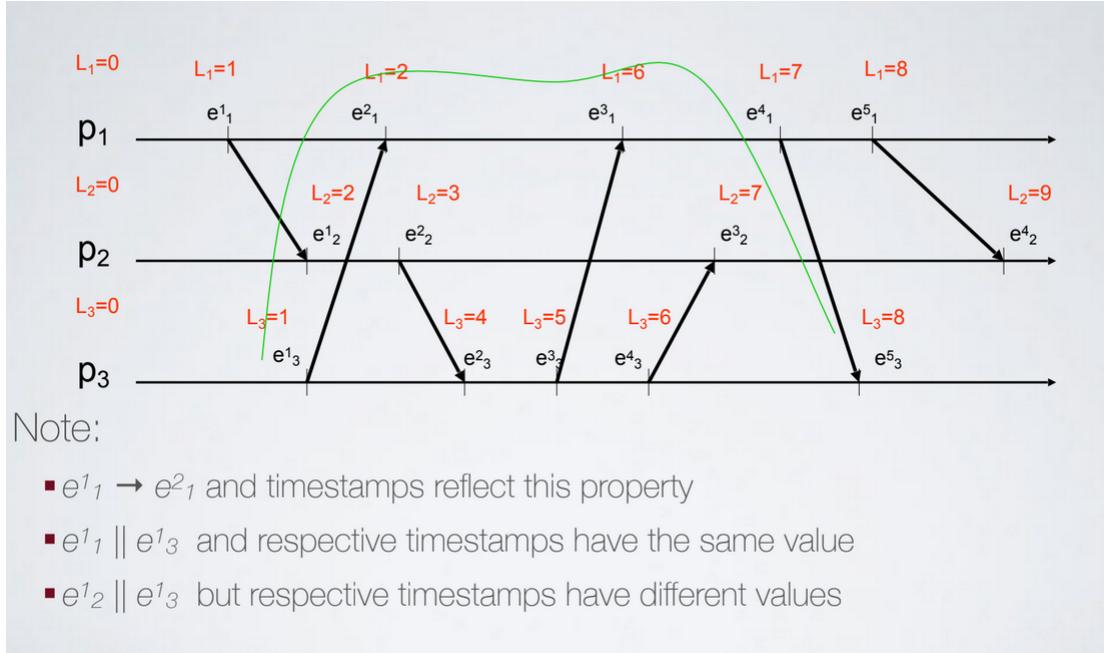
We attach timestamp to each event, and this number will have a special property, indicated as  $L(e)$ . We have that if  $e$  happens before  $e'$ , then  $L(e) < L(e')$ . (Only works with  $\implies$ , not with  $\leftarrow$ ). If the logical clock of  $e$  is smaller than  $e'$ , then this doesn't imply that  $e$  has happened before  $e'$ , they might be concurrent.

With the logical clock everyone starts with a counter  $L_i$ , initially set to 0. We then increase the counter, by following two rules:

- when a local event happens, we increase the counter by 1
- when we send a message, we also increase the counter by 1. Along side sending the message, we also send the value of the counter ((m, 2) for example)

When I receive a message as process, we'll do a  $1 + \text{Max}(L, L(m))$  (maximum between our logical clock and the value sent with the message), setting our counter to that value.

We have to show that  $e \rightarrow e'$  ( $e$  happens before  $e'$ ), then  $L(e) < L(e')$ . If  $e$  happens before  $e'$ , we have a path between  $e$  and  $e'$ . Whether the possibility, we increase by at least one the value of the logical clock. Then if  $e$  happens before  $e'$ , we'll have that  $L(e')$  is greater at least by 1 than  $L(e)$ .

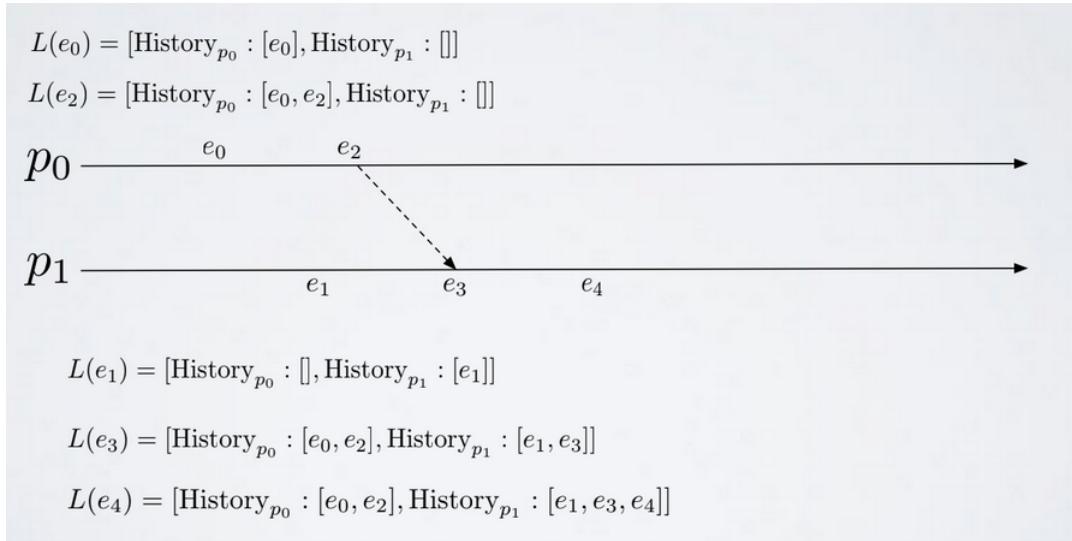


Note:

- $e^{1_1} \rightarrow e^{2_1}$  and timestamps reflect this property
- $e^{1_1} \parallel e^{1_3}$  and respective timestamps have the same value
- $e^{1_2} \parallel e^{1_3}$  but respective timestamps have different values

We would like to have now the  $\iff$  for the relationship of local clock. There exists Vector Clock, which allows us to design a system which has the other direction of the property. We do not have a simple number, but a vector. By comparing vectors, we'll be able to tell if an event happen before another.

As step 1 we construct this vector as a list of events, as step 2 we transform this list of events in a vector. We attach then a data structure attached to an event, which summarizes the history of the events.



We'll have one entry (in the data structure) for each process.

At first the data structure will be  $[], []$  both for  $p_0$  and  $p_1$ . When process zero creates an event  $e_0$ , its history will change. At this point we still do not know nothing about  $e_1$ . Same goes for  $e_1$  for  $p_1$ .

The local event is trivial, the "problem" is when we send a message. The message will include the history of  $p_0$  and  $p_1$  (which at  $e_2$  it's still empty). When  $p_1$  receives the message, he will now know

the history of  $p_0$  and its as well. We'll fill the history of  $p_0$  with the longest one (received with the message), and the history of  $p_1$  will be made by  $e_1$  and  $e_3$ .

What we are doing is that we are constructing the view of the process. Now, how do we compare data structures? We say that  $L(e) < L(e')$  if, for each location of the history, the list is smaller than the other. As example of history of  $e_0$  and  $e_3$ . Some data structures are not comparable, for example  $e_0$  and  $e_1$ , as we have  $[[e_0], []]$  for  $e_0$  and  $[], [e_1]$ . For the first location  $p_0$  is winning, but then for the second one  $p_1$  is winning, they are not comparable. This happens because the events are concurrent. Let's see the proof. We have to show that if  $e \rightarrow e'$  then  $L(e) < L(e')$  (data structure of  $e$  is contained in  $e'$ 's one). If  $e$  happens before  $e'$ , then there is a path between  $e$  and  $e'$ , then the data structure at  $e'$  will contain the one of  $e$ . How to show the other direction? If the data structure of  $e$  is smaller than  $e'$ , then  $e \rightarrow e'$ . If the data structure of  $e$  is less than  $e'$ , then  $L(e)$  is contained in  $L(e')$ , then all events happened at  $e$  then they are contained in  $L(e')$ . Because of the fact that  $e$  is contained also in  $L(e')$ , then there must be a path that somehow has connected  $e$  to  $e'$ , then  $e \rightarrow e'$ .

This data structure captures causality.

What are the disadvantages of using this data structure? The use of a large amount of memory for a lot of events. How do we get something more compact? Let's say that on my system I find two events,  $e$  and  $e'$ . I examine the data structure of  $e$  and  $e'$  and I find that, for a certain location  $x$  in the history, I find something similar.

To have something different would be impossible (like in the slide's example)

$$\exists e : L(e) = [\dots, History_{p_x} = [e_0, e_1, e_2], \dots] \quad \text{Possible?}$$

$$\exists e' : L(e') = [\dots, History_{p_x} = [e_0, e_1, e_3, e_4], \dots]$$

The only one who can append something at the history of  $p_x$  is  $p_x$  itself, so the last event can't be different.

What matters is not the content of the list, but the size of the list. We can then compress the data structure by just using the size of the list.

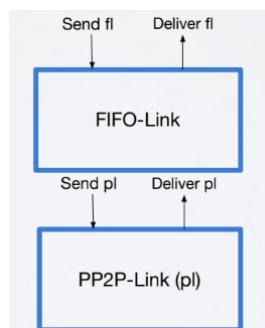
So at first we start with a vector for all zeros, having empty lists for each process.

Let's say we are  $p_1$  and we create a new event, then I increase the value of  $+1$  for the size of  $p_1$ . When another process receives a message, it compares the elements between the lists and take the maximum. We can say that one vector is less than the other ( $V(e) < V(e')$ ) if all components are  $\leq$  and there is one component who is strictly less.

## 2.4 Lamport's Mutual Exclusion Algorithm

In this algorithm, we do not tolerate any crashes, but before talking about it, we need to introduce FIFO links. A FIFO link has the same properties as a P2P link, but also the following property:

**FIFO:** if a process  $p$  delivers  $m$  before  $m'$  then  $m$  was sent before  $m'$



---

**Algorithm 5** FIFO Links - client  $p_i$ 


---

```

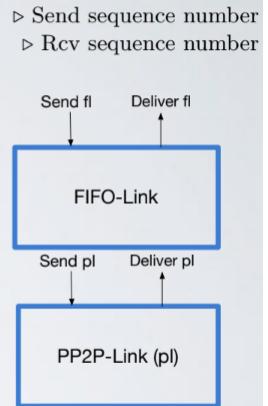
1: upon event INIT
2:   Pending =  $\emptyset$ 
3:   seq = 0
4:   del_seq = 0

5:  $\triangleright$ 
6: upon event SEND FL( $\langle p_d, m \rangle$ )
7:   SEND PERFECT( $p_d, \langle m, seq \rangle$ )
8:   seq = seq + 1

9: upon event DELIVERY FROM PERFECT( $\langle p_s, \langle m, seq \rangle \rangle$ )
10:  Pending = Pending  $\cup \{\langle p_s, \langle m, seq \rangle \rangle\}$ 

11: upon event  $\exists \langle p_s, \langle m, send\_seq \rangle \rangle \in Pending$  SUCH THAT  $send\_seq = del\_seq$ 
12:   DELIVERY FL( $p_s, m$ )
13:   Pending = Pending  $\setminus \{\langle p_s, \langle m, send\_seq \rangle \rangle\}$ 
14:   del_seq = del_seq + 1 typo: here you have to increment del\_seq

```



- **(FIFO)** if a process p delivers m before m', then m was sent before m'.

In the init handler we have a pending set (self-explanatory), a seq number equal to 0 and a del\_seq equal to 0.

In the send event handler we first send the message via the P2P link in the form  $\langle msg, seq\_number \rangle$ , and only after that do we add 1 to the seq number.

When we get a delivery from the perfect link, we just add that message to our pending set.

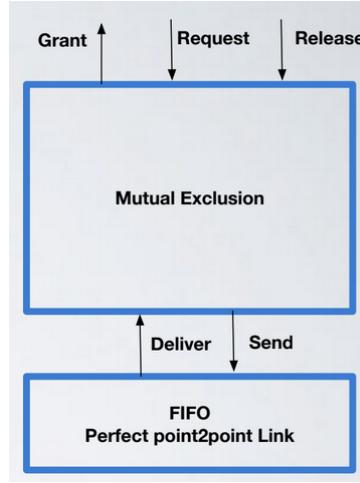
If there exists a message in pending such that the seq = del\_seq, we first deliver the message and after that remove it from our pending list, lastly we just add 1 to our del\_seq

Let us prove the FIFO property, as always let us start by contradiction:

Suppose p delivers in reverse order m' and then m, even if m was sent before m', this means that the timestamp of m is lower than that of m', let t be the instant at which p delivers m', this implies that  $del\_seq = timestamp(m')$ , this is of course, a contradiction, since for  $del\_seq$  to be equal to  $timestamp(m')$ , we should have had  $del\_seq = timestamp(m)$  (remember that  $del\_seq$  grows by one unit each delivery). Thus m has been delivered before t.

Finally let's talk about Lamport's Mutual Exclusion Algorithm, which works only with FIFO links and does not tolerate failures.

The Mutual Exclusion block will have a request in input (a process  $p_1$  asks to enter the critical zone, and it will give it if possible), we also have a release in input (for when we have a go-ahead to actually use a resource) and a Grant in output (give access to a resource to a process), it will communicate with a FIFO link via send and deliver.



We have 3 properties to respect:

- **Mutual Exclusion:** at any time t, only one process is inside the Critical Section
- **Liveness:** If a process p requests access, then it will eventually enter the CS
- **Fairness:** if the request of process p happens before the request of process q, then q cannot access the CS before p.

When a process wants to enter the CS it sends a request to all other processes, a history is maintained via scalar clocks.

Each transmission and reception event is relevant to the computation, the clock is incremented after each action

---

**Algorithm 1** Lamport's ME Algorithm on process  $p_i$  - MSGS are REQ, ACK, RLS

---

```

1: upon event INIT
2:   Requests = Acks =  $\emptyset$ 
3:   scalar_clock = 0
4:   my_req =  $\perp$ 
5:    $\Pi = \{p_0, p_1, \dots, p_{n-1}\}$                                  $\triangleright$  Set of all processes

6:  $\triangleright$  Request access to CS from upper layer
7: upon event REQUEST
8:   scalar_clock = scalar_clock + 1
9:   my_req = (REQ, ts =  $< i, \text{scalar\_clock} >$ )
10:  for all  $p_j \in \Pi$  do
11:    SEND FIFOUPERFECTLINK( $p_j, \text{req\_msg}$ )  $\triangleright$  Send a REQ containing my ID (i) and ts (scalar_clock) to all
      $p \in \Pi$ 

12:  $\triangleright$  Release CS from upper layer
13: upon event RELEASE
14:   my_req =  $\perp$ 
15:   scalar_clock = scalar_clock + 1
16:   for all  $p_j \in \Pi$  do
17:     SEND FIFOUPERFECTLINK( $p_j, (\text{RLS}, ts = < i, \text{scalar\_clock} >)$ )

18:  $\triangleright$   $ts(x) < ts(y)$  when  $\text{scalar\_clock}$  of x is less than the one of y, or they are equal and the id that sent x is less
     than the id that sent y
19: upon event  $\#req \in \text{Requests} : ts(\text{req}) < ts(\text{my\_req}) \wedge \forall p \in \Pi : \exists m \in \text{Acks} | ts(m) > ts(\text{my\_req}) \wedge \text{sender}(m) = p$ 
20:   trigger event GRANTED

21: upon event DELIVER MESSAGE(m)
22:   scalar_clock = max(clock(m), scalar_clock) + 1
23:   if m is a REQ then
24:     Request_set = Request_set  $\cup \{m\}$ 
25:     scalar_clock = scalar_clock + 1
26:     SEND FIFOUPERFECTLINK(sender(m), (ACK, ts =  $< i, \text{scalar\_clock} >$ ))
27:   else if m is a ACK then
28:     Acks = Acks  $\cup \{m\}$ 
29:   else if m is a RLS  $\wedge \exists req \in \text{Request\_set} : \text{sender}(req) = \text{sender}(m)$  then
30:     Requests = Requests  $\setminus \{req\}$ 

```

---

This is algorithm for process  $p_i$  (all processes have their own private variables)

In init we have a request set and an ACK set, the request set is used to store requests for CS access by other processes, and the ACK sets are the replies we sent to the other processes. After that we initialize our scalar\_clock to 0 (+1 per event), we also have a variable my\_req set to Null and finally a set containing all the other processes  $\prod$

When we want to send a request to access the CS:

We add 1 to our scalar clock, add it to my\_req and send a message via FIFO to all processes in the  $\prod$  set, our requests contain our ID and the scalar clock value, in case the scalar clock value is the same, we use the one with the lower ID, in case they are different the lowest scalar clock wins. We also have a special field with flag REQ.

When we want to release the CS:

Set my\_req to null, add 1 to our scalar\_clock, and for all processes send a message where you say RLS  
 $< p_id, scalar\_clock >$

There's a handler that grants us access to the CS:

The handler tells us that when in the req set there doesn't exist a req with a lower timestamp and for all processes in the systems we have an ACK such that the timestamp of the ACK is higher than the timestamp of my message, only then can we enter the CS.

And finally a handler to process everything (messages received etc.):

When we deliver a message, we update our scalar clock as according to Lamport, after that we check the kind of message we received, if it's a REQ we put it into the Requests set and we increase our scalar clock, finally we send an ACK with our ID and scalar\_clock

If it's an ack we add the message to our acks set.

If it's a RLS, we check if there is a request in the req set where the sender of the currently processed req is equal to the sender of the message we received, then we remove that request from the requests set.

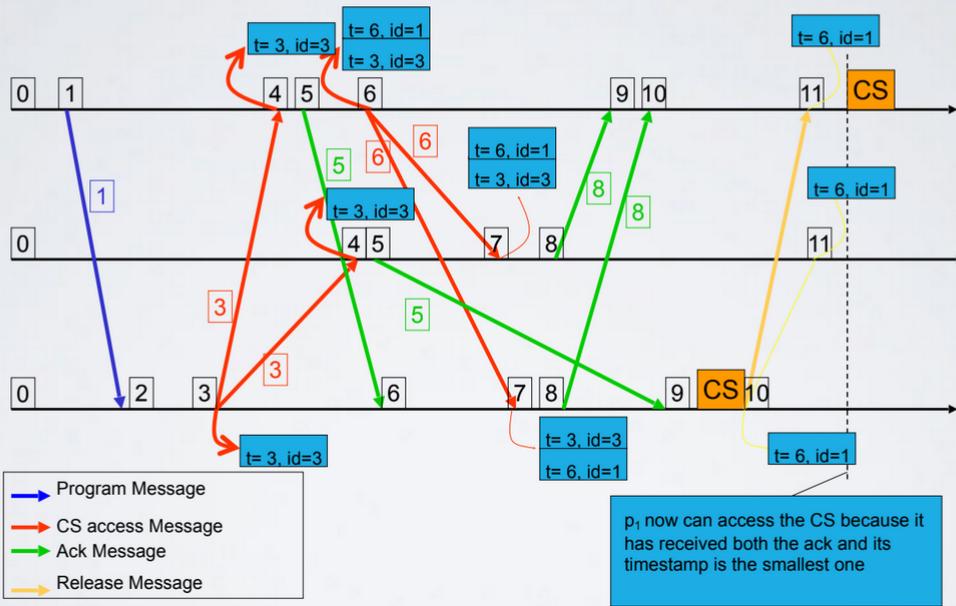
## Proof of Mutual Exclusion

Let us assume we have 2 processes that enter the critical section: p1 and p2, if p2 entered that must mean that it received acks from all other processes and only after it sent the req message, same for p1. But if p1 entered before p2 that must mean that the req message of p1 has a lower timestamp than the req of p2.

## Proof of fairness

Suppose  $p_i$  enters before  $p_j$  even if req of  $p_j$  happened after req  $p_i$ . Since  $p_j$  enters only after the ack of  $p_i$ , by FIFO it sees req  $p_j$  before receiving the ack that allows him to enter. Since req  $p_j$  happens before req  $p_i$  we have  $ts(req, p_j) < ts(req, p_i)$  thus req( $p_j$ ) is not the request with minimal timestamp in the set of  $p_i$ .

# LAMPORT'S ALGORITHM: EXAMPLE



As a measure of complexity of the algo we either have a message complexity (number of messages we need to have), or time unit complexity (we suppose that each action takes 1 unit of time). In the worst case we need to wait  $(N-1) + 2$  (wait for release by everyone else) + req

## 2.5 Failure Detectors

Failure detectors, as the name imply, tell us when a process has crashed, they only work for crash failures and are described by the following two properties:

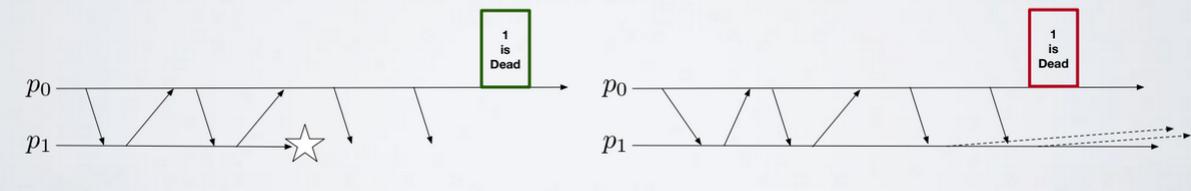
- **Completeness**, the ability to detect all failures
- **Accuracy**, the ability to avoid false suspicions

To actually implement a failure detector we utilize pings.

For now we'll talk about **perfect failure detectors** (P), whose system model is synchronous, it utilizes crash failures and utilizes perfect synchronous P2P links. The basic way to tell if a process has crashed is by utilizing a timeout, in our case let us utilize the clock of our FD, and if the process doesn't respond then we can assume its crashed.

For the purposes of this course, there is no algorithm for FD within an asynchronous system.

In asynch?



# PERFECT FAILURE DETECTORS (P) SPECIFICATION

---

## Module 2.6: Interface and properties of the perfect failure detector

---

**Module:**

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

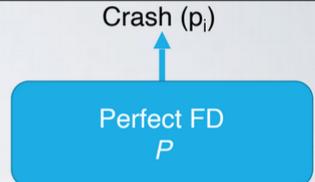
**Events:**

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

**Properties:**

**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.



- **PFD1:** Strong Completeness: Eventually, every process that crashes is permanently detected by every correct process. (If a process crashes after some time, all processes that are correct will signal this crash)
- **PFD2:** Strong Accuracy: If a process  $p$  is detected by any process, then  $p$  has crashed. (If the algo signals that a process has crashed, then it must have crashed)

# PERFECT FAILURE DETECTORS (P) IMPLEMENTATION

---

## Algorithm 2.5: Exclude on Timeout

---

**Implements:**

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

**Uses:**

PerfectPointToPointLinks, **instance**  $pl$ .

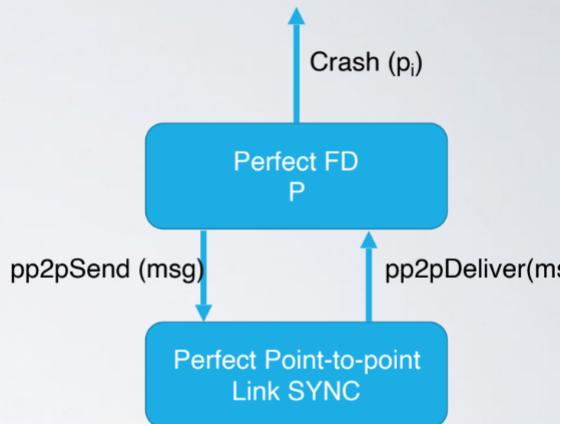
```

upon event <math>\langle \mathcal{P}, \text{Init} \rangle \text{ do} 
  alive := \Pi;
  detected := \emptyset;
  starttimer(\Delta);

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

upon event <math>\langle pl, \text{Deliver} | q, [\text{HEARTBEATREQUEST}] \rangle \text{ do}
  trigger <math>\langle pl, \text{Send} | q, [\text{HEARTBEATREPLY}] \rangle</math>;

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



When we initialize we have an alive set with all the processes, a detected set that is empty (where we store the crashed processes), after that we start a timer with delta.

When we have a deliver event (we get pinged) we reply, at the same time with another handler we add the process that pinged us to the alive set.

When we get a timeout, for every process we first check if it isn't part of the alive set and if it isn't part of the crash set, if both of these conditions are true, we add p to our detected set and we trigger Crash(p).

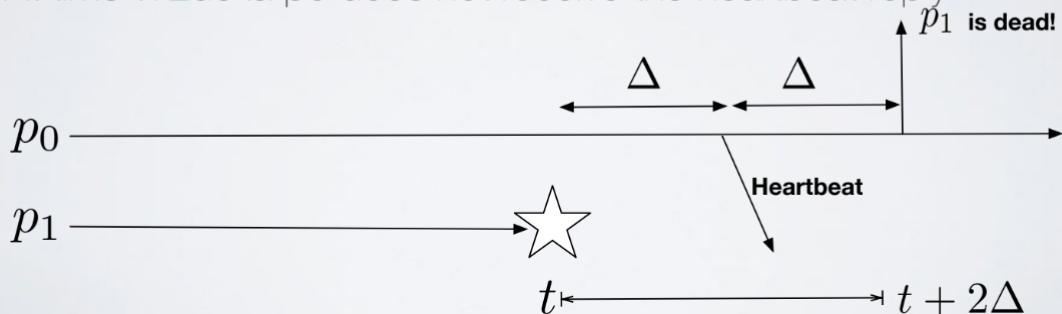
If one of these conditions isn't true then we send a ping to that process.

### Proof of Correctness

Strong completeness:

Assume  $p_1$  crashes at time  $t$ , then  $p_0$  detects it at most by time  $t+2\Delta$ .

- 1) At time  $t+\Delta$  (at most)  $p_0$  sends an heartbeat request
- 2) At time  $t+2\Delta$   $p_0$  does not receive the heartbeat reply



Strong accuracy:

Assuming  $2 \times \text{max\_delay} \leq \Delta$

Proof. By contradiction.

If  $p$  is detected faulty but it was alive, then either:

- My heartbeatREQ did not reach it by the timeout.
- Its heartbeatREPL did not reach me by timeout.

Both contradicts the assumption.

# GIVEN P IS IT POSSIBLE?

|       | $t$ | $t'$ | $t''$ | $t'''$ | $t+x$     |
|-------|-----|------|-------|--------|-----------|
| $p_0$ | {}  | {}   | {}    | {}     | { $p_3$ } |
| $p_1$ | {}  | {}   | {}    | {}     | {}        |
| $p_2$ | {}  | {}   | {}    | {}     | { $p_3$ } |
| $p_3$ | {}  | {}   | ★     |        |           |
| $p_5$ | {}  | {}   | {}    | {}     | {}        |
| $p_6$ | {}  | {}   | {}    | {}     | {}        |

Yes this is possible as strong completeness only says that eventually all crashed processes are detected by all correct processes, it does not say when.

But then considering the algo we have used up until now (exclude on timeout), if at time  $t$  I do not have process  $p$  in my detected set can I state that  $p$  is alive at time  $t$ ? No, as FD do not give us information about current processes, you simply don't know if its dead.

Let us try and improve p, normally each process needs atleast  $2N$  messages ( $N$  pings and  $N$  pong) to detect a crash, so the total number of messages becomes  $2N^2 * N$ .

We could cut the number of messages by half since due to the fact that every process pings all other processes, then we could count the pings as statements of being alive.

A way to actually do this is **Round Based P**, which introduces the concept of rounds.

---

**Algorithm 3** Round-Based, Failure Detector  $\mathcal{P}$  on process  $p_i$

---

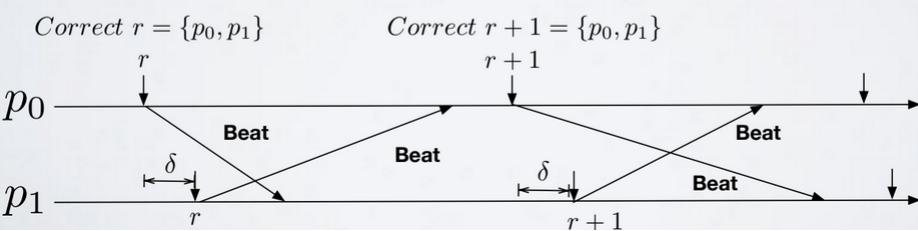
```

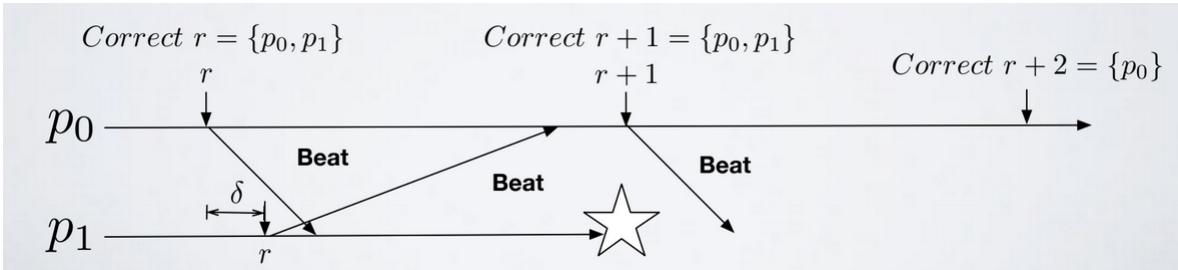
1: upon event INIT
2:    $C_v$                                      ▷ Synchronizer's clock
3:    $T$ 
4:    $Corrects = \{p_0, p_1, p_2, \dots, p_n\}$       ▷ Round length:  $T > 2 \cdot (max\_delay + \delta)$ 
5:    $RoundAlive = Corrects$ 
6:    $r = -1$ 

7: ▷ Actions to execute at each new round
8: upon event  $C_v = kT$  FOR SOME  $k \in \mathbb{N}^+$ 
9:    $r = r + 1$ 
10:  for all  $p_j \in Corrects \setminus RoundAlive$  do
11:     $Corrects = Corrects \setminus \{p_j\}$ 
12:    TRIGGER  $\langle Crash[p_j] \rangle$ 
13:   $RoundAlive = \{p_i\}$                                 ▷ Myself
14:  for all  $p_j \in Corrects$  do
15:    SEND( $< p_i, BEAT >$ )
16: upon event DELIVERY( $\langle p_j, BEAT \rangle$ )
17:    $RoundAlive = RoundAlive \cup \{p_j\}$ 

```

---





FACT 1: If I send a BEAT I am alive.

FACT 2: If I don't send a BEAT, I am dead.

**We have to show that:** If I send a BEAT to p when my round number shows r, then the beat reaches p when its round number shows r.

By the algorithm, I send my BEAT as soon as my round is r.

By the bounded skew of clock synch. If time t is the first time at which my round number is r, then p will be at round r by at most  $t + \delta$

By the round length, p stays in round r at least until

$$t - \delta + T = t - \delta + 2\delta + 2\text{max\_delay} > t + \delta + \text{max\_delay}$$

By the delay of the link, my BEAT reaches p in the interval  $(t + \delta, t + \text{max\_delay}]$

FACT 1: If I send a BEAT I am alive.

FACT 2: If I don't send a BEAT, I am dead.

**We have that:** If I send a BEAT to p when my round number shows r, then the beat reaches p when its round number shows r.

**Strong completeness:** If p dies, I will not receive the BEAT and I will signal it.

**Strong accuracy:** If I do not receive the expected BEAT from p then, The only possible reason is that p did not send it. Therefore p is dead.

So why is this algorithm important? Because with this we can see the equivalence between a synchronous system and round-based systems.

Another question we need to ask ourselves is if decreasing is worth it, certainly we may detect crashes faster, but we may also be wrong and make mistakes, therefore what if we correct those mistakes? This leads us to the concept of **Eventually Perfect Failure Detectors** (denoted with  $\diamond P$ ), which are FD that may make mistakes, but eventually they will stop making them.

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

### Module:

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

### Events:



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

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

### Properties:

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

**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  $p$  may suspect another one  $q$  because the chosen timeout was too short
- $\diamond P$  is ready to change its judgment as soon as it receives a message from  $q$  (updating also the timeout value)
- If  $q$  has actually crashed,  $p$  will not change its judgment anymore.

# EVENTUALLY PERFECT FAILURE DETECTORS ( $\diamond P$ )

**Algorithm 2.7:** Increasing Timeout

**Implements:**

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

**Uses:**

PerfectPointToPointLinks, **instance**  $pl$ .

**upon event**  $\langle \diamond P, Init \rangle$  **do**

```
alive :=  $\Pi$ ;
suspected :=  $\emptyset$ ;
delay :=  $\Delta$ ;
starttimer(delay);
```

**upon event**  $\langle \text{Timeout} \rangle$  **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  $\langle \diamond P, Suspect \mid p \rangle$ ;
    else if  $(p \in alive) \wedge (p \in suspected)$  then
        suspected := suspected  $\setminus \{p\}$ ;
        trigger  $\langle \diamond P, Restore \mid p \rangle$ ;
    trigger  $\langle pl, Send \mid p, [\text{HEARTBEATREQUEST}] \rangle$ ;
alive :=  $\emptyset$ ;
starttimer(delay);
```

**upon event**  $\langle pl, Deliver \mid q, [\text{HEARTBEATREQUEST}] \rangle$  **do**  
**trigger**  $\langle pl, Send \mid q, [\text{HEARTBEATREPLY}] \rangle$ ;

**upon event**  $\langle pl, Deliver \mid p, [\text{HEARTBEATREPLY}] \rangle$  **do**  
 $alive := alive \cup \{p\}$ ;

