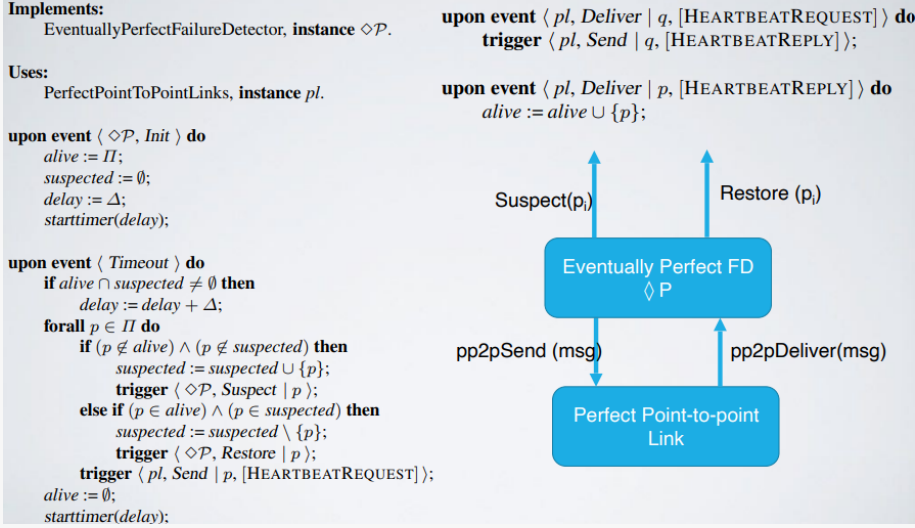


5. TIME

EVENTUALLY PERFECT FAILURE DETECTORS $\diamond P$

Algorithm:



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. After that time a message sent by a correct process p to another one q will be delivered within a bounded time (the time is MAX_{DELAY} unknown to us). If p was wrongly suspected by q , then q will revise its suspicious. Moreover q increases its delay, if the new delay is less than MAX_{DELAY} , eventually q suspects again p and then it corrects again its delay. After a finite number of errors the delay of q will be greater than MAX_{DELAY} .

Is possible that the delay is big enough that can happen a situation where a process p_2 communicates at process p_1 that he's dead but the message is received after the detection of p_2 dead line from p_1 .

Solution: add branches T_x .

Property of Suspected:

Lemma: take two correct processes, p_1 and p_2 and let $suspected_1$ and $suspected_2$ be the respective sets. There is a time t (stabilization time), after which $suspected_1 = suspected_2$.

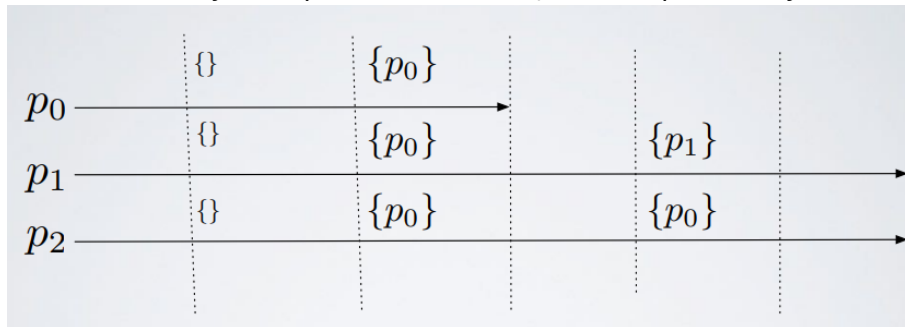
Proof (by contradiction): suppose that exists a p in $suspected_1$ but not $suspected_2$:

- if p is correct this violates the **eventual strong accuracy**.
- if p crashed this violates the **strong completeness**.

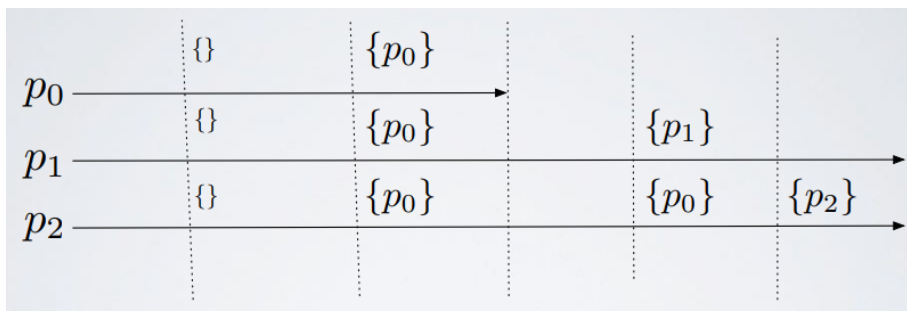
Leader Election

Sometimes, we may be interested in knowing one process that is alive instead of monitoring failures. In this case we can use a different oracle (called **leader election module**) that reports a process that is alive.

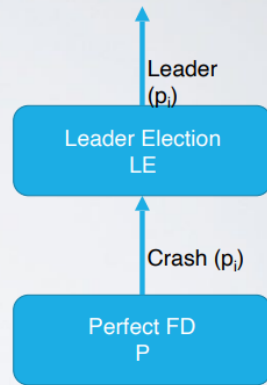
- **Name:** LeaderElection, **instance** le .
- **Indication:** $\langle le, Leader \mid p \rangle$: indicates that process p **is elected as a leader**.
- **Properties:**
 - **LE1** (eventual detection): either there is no correct process, or some correct process is eventually elected as the leader (liveness property).
 - **LE2** (accuracy): if a process is leader, then all previously elected leaders have crashed



When p_0 dies it's okay because before he dies everyone know that he's the leader, but in p_2 after p_0 dead, the leader isn't changed. (broke the eventual detection property).



In this case the accuracy property is broken because p_2 changes the leader even if p_1 is still alive.

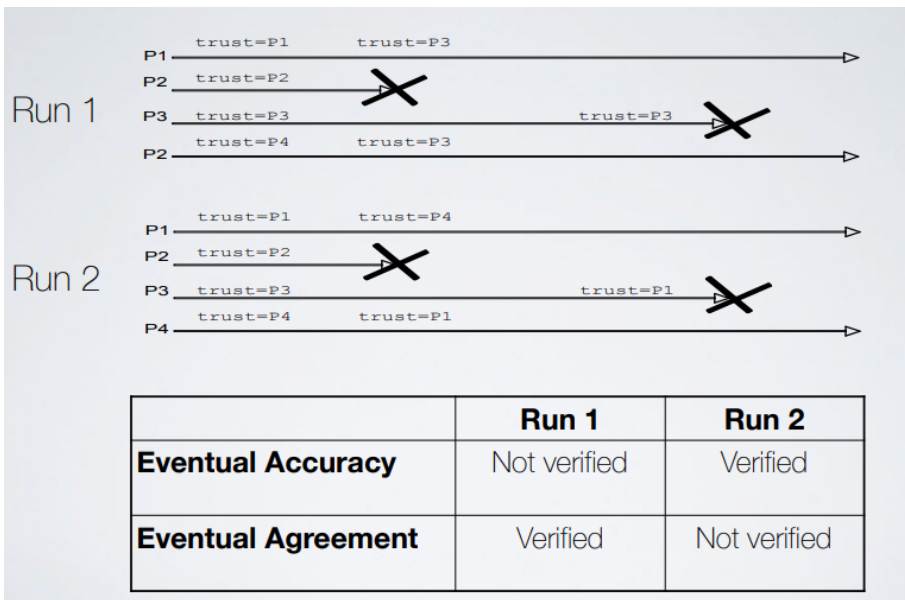
Algorithm:**Implements:**LeaderElection, **instance** le .**Uses:**PerfectFailureDetector, **instance** \mathcal{P} .**upon event** $\langle le, Init \rangle$ **do** $suspected := \emptyset;$ $leader := \perp;$ **upon event** $\langle \mathcal{P}, Crash \mid p \rangle$ **do** $suspected := suspected \cup \{p\};$ **upon** $leader \neq \maxrank(\Pi \setminus suspected)$ **do** $leader := \maxrank(\Pi \setminus suspected);$ **trigger** $\langle le, Leader \mid leader \rangle;$  $(\Pi \setminus suspected)$ is the subset of alive processes.**Correctness:**

- **Eventual detection:** from the strong completeness of P (perfect failure).
- **Accuracy:** from the strong accuracy of P and the total order on the ranks (IDs) of processes.

Eventual Leader Election Ω

- **Name:** *EventualLeaderDetector*, **instance** Ω .
 - **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.
- Both of properties are **liveness**.

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



We can build an eventual leader election using **crash-stop process abstraction**:

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

Algorithm:

Implements:

EventualLeaderDetector, **instance** Ω .

Uses:

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

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

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

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

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

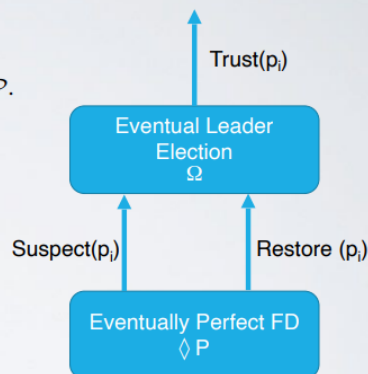
upon event $\langle \diamond P, \text{Restore} \mid p \rangle$ do

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

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

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

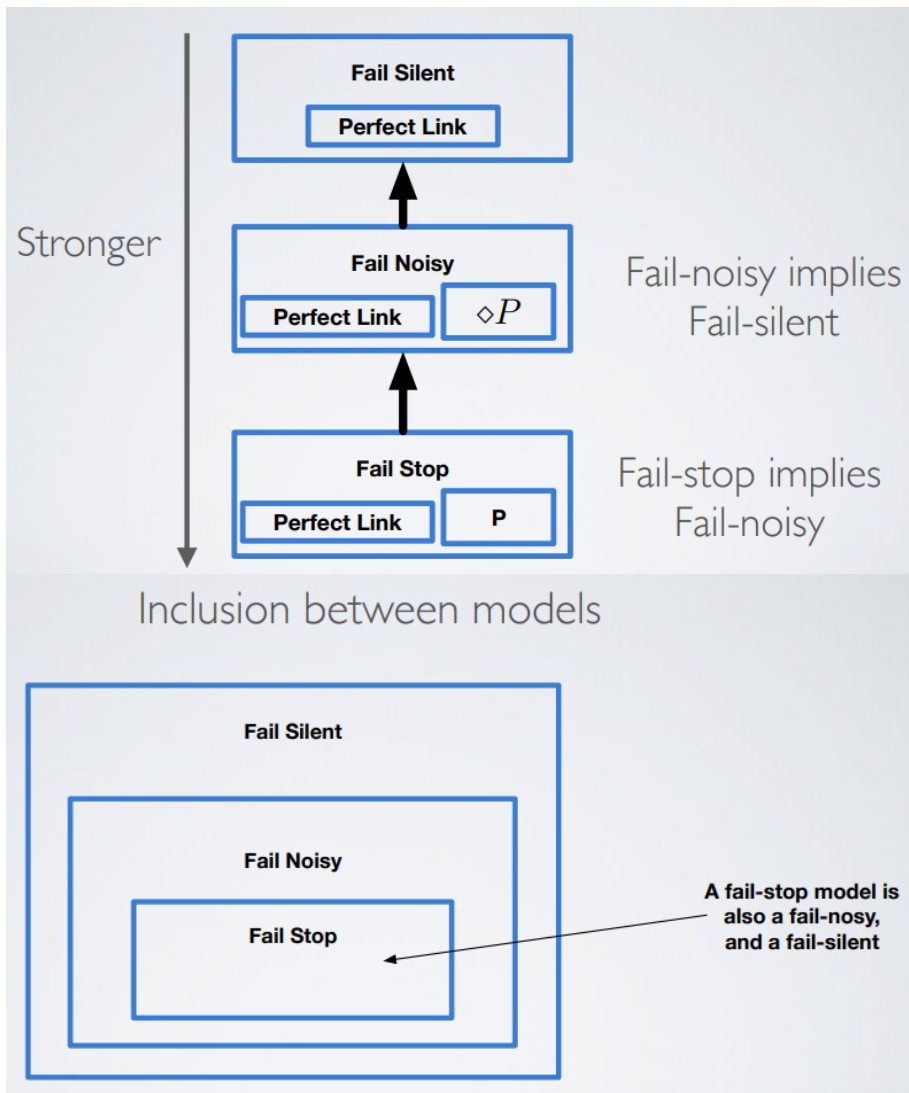
trigger $\langle \Omega, \text{Trust} \mid \text{leader} \rangle;$



Proof:

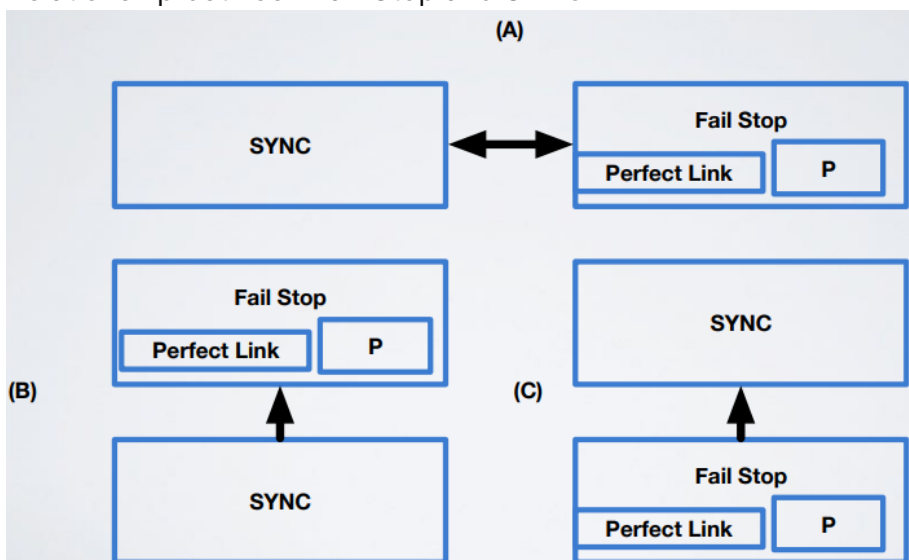
- **ELD1** (Eventual Accuracy): by the strong completeness of the FD we have that eventually suspected set contains all the crashed processes. Thus $\Pi \setminus \text{suspected}$ contains only correct processes (or its empty).
- **ELD2** (Eventual Agreement): for any pair of correct processes, their suspected sets eventually stabilises to the same content (by the property of the FD). If the set are equals $\Pi \setminus \text{suspected}$ returns the same ID on both processes.

THREE MODELS



We can also say that the set of problems solvable in fail-stop includes the problems solvable in other models.

Relationship between Fail-Stop and SYNC:



SYNC is stronger than Fail-Stop (Correct answer is B !!!!possibile domanda esame!!!!).

Problems we can solve with fail-stop is strictly contained in subset of problems we can solve

with SYNC (SYNC can solve problems based on **time** like clock synchronization):



APPLICATION OF FAILURE DETECTOR AND LEADER ELECTOR

Using P to make **Lamport's Mutual Exclusion** fault tolerant.

Events:

- **Request:** from upper layer - requests access to Critical Section (CS).
- **Grant:** to upper layer - grant the access to CS.
- **Release:** from upper layer - release the CS.

Properties:

- **Mutual Exclusions:** at any time t , at most one non-crashed process is inside the CS.
- **Liveness:** if a correct process p requests access, then it eventually enters the CS.
- **Fairness:** if a correct process p requests access before a process q , then q cannot access the CS before p .

ORIGINAL ALGORITHM

Critical points:

- 1) if you crashed you cannot release
- 2) If you crashed you cannot ack

```

1: upon event INIT
2:   Requests = Acks = {}
3:   scalar_clock = 0
4:   my_req = ⊥
5:   Π = {p0, p1, ..., pn-1}
6: ▷ Request access to CS from upper layer
7: upon event REQUEST
8:   scalar_clock = scalar_clock + 1
9:   my_req = (REQ, ts = < i, scalar_clock >)
10:  for all pj ∈ Π do
11:    SEND FIFOPERFECTLINK(pj, req_msg) ▷ Send a REQ containing my ID (i) and ts (scalar_clock) to all
    p ∈ Π
12: ▷ Release CS from upper layer
13: upon event RELEASE
14:   Requests = Requests \ {req_msg}
15:   scalar_clock = scalar_clock + 1
16:   for all pj ∈ Π do
17:     SEND FIFOPERFECTLINK(pj, (RLS, ts = < i, scalar_clock >))
18: ▷ ts(x) < ts(y) when 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 ∈ Requests : ts(req) < ts(my_req) ∧ ∀p ∈ Π : ∃m ∈ Acks | ts(m) > ts(my_req)
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 ∪ {m}
25:     scalar_clock = scalar_clock + 1
26:     SEND FIFOPERFECTLINK(sender(m), (ACK, ts = < i, scalar_clock >))
27:   else if m is a ACK then
28:     Acks = Acks ∪ {m}
29:   else if m is a RLS ∧ ∃req ∈ Request_set : sender(req) = sender(m) then
30:     Requests = Requests \ {req}
  
```

Patch 1:

```

1: upon event INIT
2:   Requests = Acks = {}
3:   scalar_clock = 0
4:   my_req = ⊥
5:   Π = {p0, p1, ..., pn-1}
6:   Crashed = {}
7: ▷ Crash handler, the crash event is generated by a perfect failure detector P
8: upon event CRASH(process pj)
9:   Crashed = Crashed ∪ {pj}
10:  for all req ∈ Requests such that pj is the sender of req do
11:    remove req from Requests
12: ▷ ts(x) < ts(y) when 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
13: upon event #req ∈ Requests : ts(req) < ts(my_req) ∧ ∀p ∈ (Π \ Crashed) : ∃m ∈ Acks | ts(m) > ts(my_req)
14:   trigger event GRANTED
  
```

If the system receive a request from p_j after the crash detection we have a deadlock.

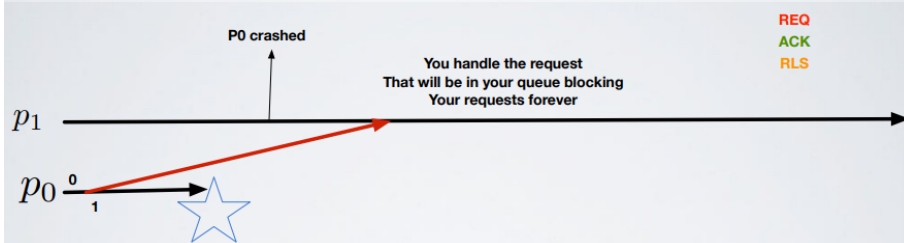
Patch 2:

```

1: upon event INIT
2:    $Requests = Acks = \emptyset$ 
3:    $scalar\_clock = 0$ 
4:    $my\_req = \perp$ 
5:    $\Pi = \{p_0, p_1, \dots, p_{n-1}\}$  ▷ Set of all processes
6:    $Crashed = \{\}$ 

7: ▷ Crash handler, the crash event is generated by a perfect failure detector P
8: upon event CRASH(process  $p_j$ )
9:    $Crashed = Crashed \cup \{p_j\}$ 
10:  for all  $req \in Requests$  such that  $p_j$  is the sender of  $req$  do
11:    remove  $req$  from  $Requests$ 

12: ▷  $ts(x) < ts(y)$  when  $scalar\_clock$  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  $y$ 
13: upon event  $\nexists req \in Requests : ts(req) < ts(my\_req) \wedge \forall p \in (\Pi \setminus Crashed) : \exists m \in Acks | ts(m) > ts(my\_req)$ 
14:  trigger event GRANTED
  
```



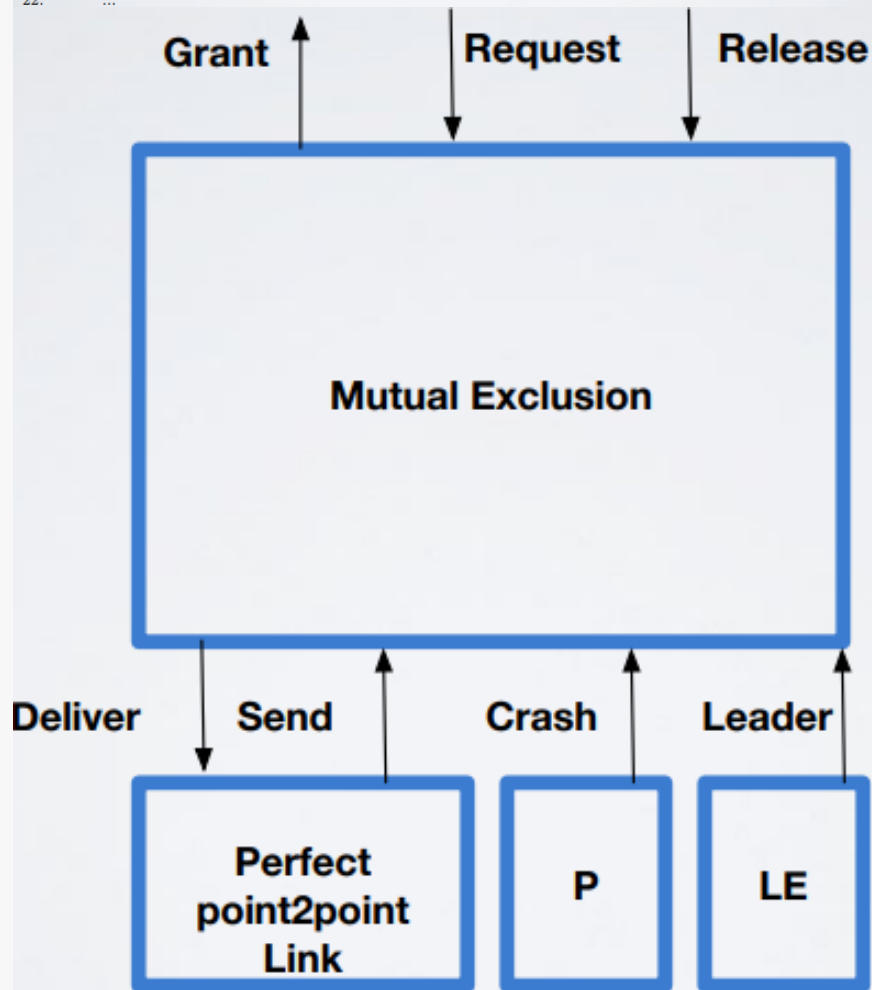
If a process crashes I put a ban on all the events related to him.

Full Algorithm with Patches:

```

1: upon event INIT
2:   Requests = Acks = {}
3:   scalar_clock = 0
4:   my_req = ⊥
5:    $\Pi = \{p_0, p_1, \dots, p_{n-1}\}$ 
6:   Crashed = {}
7:   ▷ Crash handler, the crash event is generated by a perfect failure detector P
8:   upon event CRASH(process  $p_j$ )
9:     Crashed = Crashed  $\cup \{p_j\}$ 
10:    for all req  $\in$  Requests such that  $p_j$  is the sender of req do
11:      remove req from Requests
12:   ▷  $ts(x) < ts(y)$  when 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  $y$ 
13:   upon event  $\nexists req \in Requests : ts(req) < ts(my\_req) \wedge \forall p \in (\Pi \setminus Crashed) : \exists m \in Acks | ts(m) > ts(my\_req)$ 
14:     trigger event GRANTED
15:   upon event DELIVER MESSAGE( $m$ )
16:     if  $m$  is a REQ and the sender of  $m$  is not in Crashed then
17:       scalar_clock = max(clock( $m$ ), scalar_clock) + 1
18:       Request_set = Request_set  $\cup \{m\}$ 
19:       scalar_clock = scalar_clock + 1
20:       SEND FIFOPERFECTLINK(sender( $m$ ), (ACK,  $ts = \langle i, scalar\_clock \rangle$ ))
21:     else if  $m$  is a ACK then
22:       ...

```



Exercise: write an algorithm with those properties:

- Use LE to elect a leader.
- Ask the leader for CS with a request message.
- The leader allows access to CS using FIFO order on requests.
- When done release CS using a release message.
- If the leader detects a crash p :
 - If p is not in CS, it removes the pending request of p (if any).
 - If p is in CS, it acts as p released the CS.

- *Problem*: what to do when a new leader is elected? The old leader was the only one to know who was in CS?