

Distributed Systems

Lecture 4

- o Finish Linearizability
- o CAP
- o Replicated State Machine

Checking linearizability / seq cst / ...

The rules:

I

Given: A history
- Total order

II

Extract a partial order

- Linearizability

$$OP_1 <_H OP_2$$

$$\Leftrightarrow \text{ret } OP_1 < \text{inv } OP_2$$

$$i <_H ii \quad [2 < 3]$$

$$v \quad vi$$

- i 1. p Enq(x) A
- 2. p OK() A
- ii 3. q Enq(y) B
- 4. q OK() B
- iii 5. q Enq(x) A
- 6. q OK() A
- iv 7. p Enq(y) B
- 8. p OK() B
- v 9. p Deq() A
- 10. p OK(y) A
- vi 11. q Deq() B
- 12. q OK(x) B

Reminder

Inv.: object op proc
Ret.: object OK(-) proc

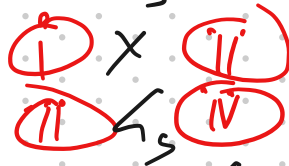
- i 1. p Enq(x) A
- 2. p OK() A

- seq cst

$$op_1 <_s op_2$$

\Leftrightarrow op_1, op_2 from same process & ret $op_1 <_{inv} op_2$

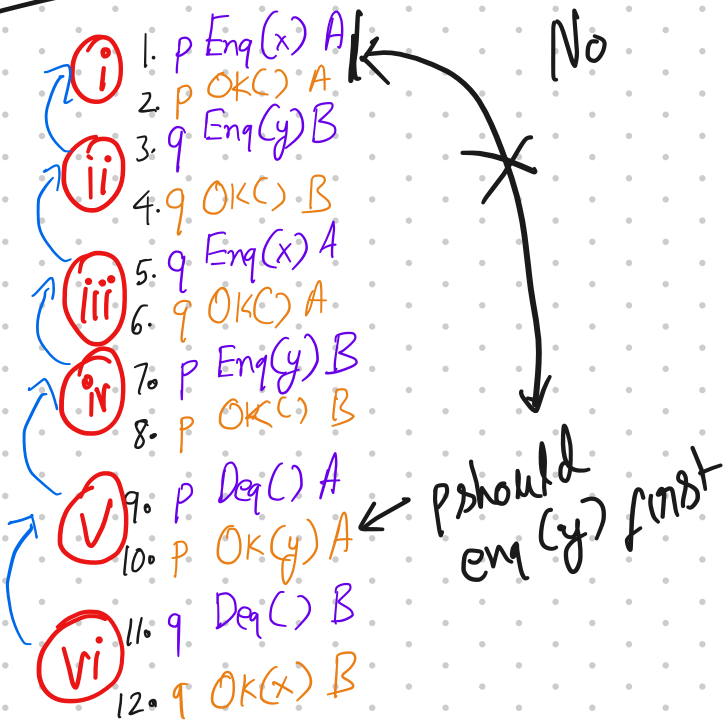
$$i <_s iii \quad [A, 2 < 5]$$



- ii 3. q Enq(y) B
- 4. q OK(x) B
- iii 5. q Enq(x) A
- 6. q OK(x) A
- iv 7. p Enq(y) B
- 8. p OK(x) B
- v 9. p Deq() A
- 10. p OK(y) A
- vi 11. q Deq() B
- 12. q OK(x) B

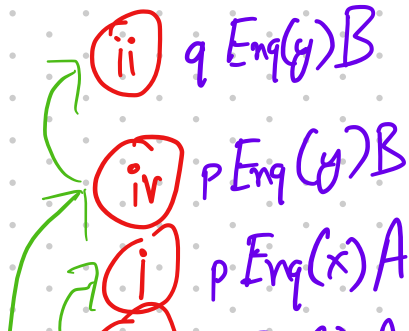
III Find an extension (SUPERSET) to partial order that is a total order on operations & meets ADT's sequential spec

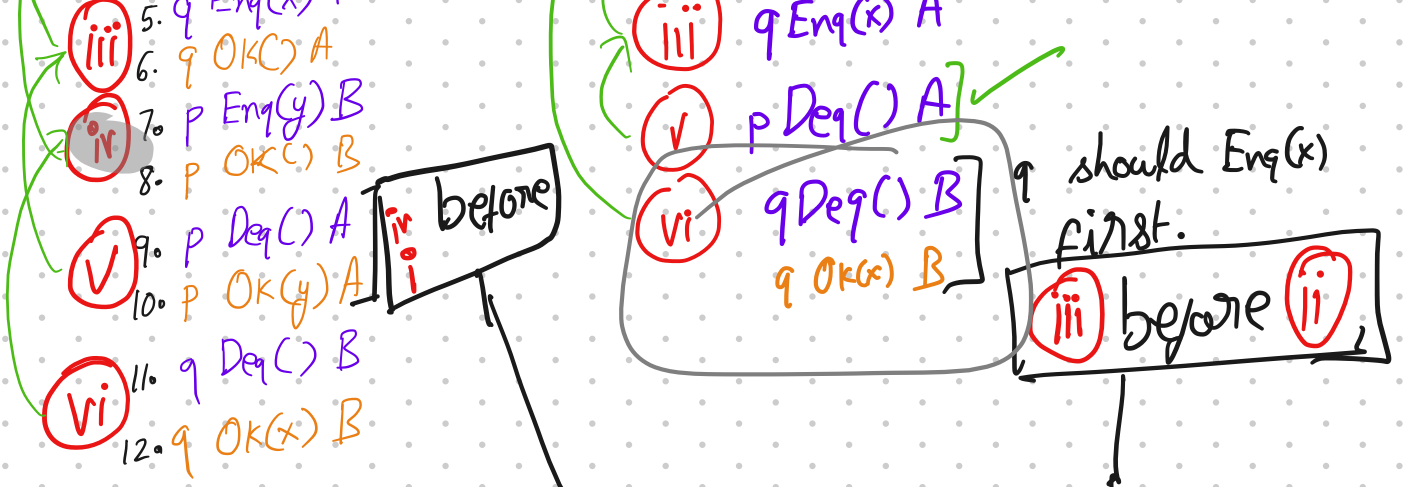
Linearizability



Serializable Seq Cst

- i 1. p Enq(x) A
- 2. p OK(x) A
- ii 3. q Enq(y) B
- 4. q OK(x) B
- iii 5. q Enq(x) A

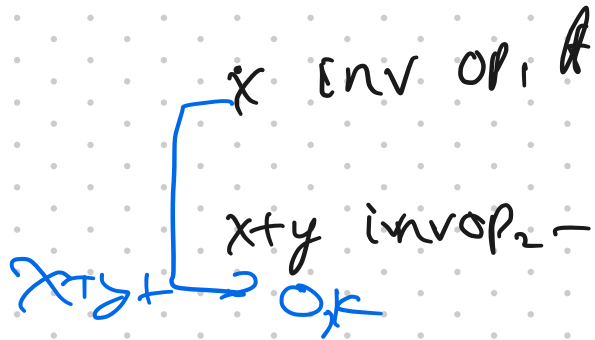




Not possible while preserving \leq_s
 \Rightarrow Not seq cst

Observe: seq cst - Fine for computed total order to not preserve history total order

Linearizable - Real-Time



Locality

```

add(item) {
  ctr.inc()
  q.enq(item)
  ret ok()
}
  
```

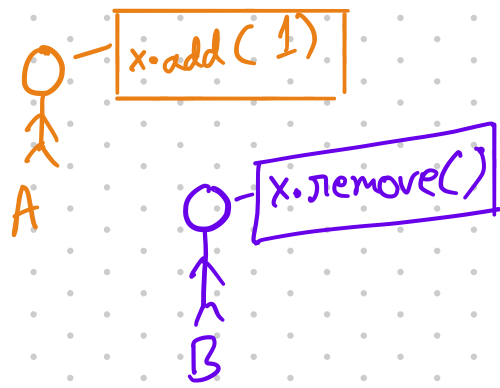
```

remove() {
  if ctr.count() > 0 {
    ctr.dec()
    return q.deq()
  } else {
    return Ok(empty)
  }
}
  
```

\perp on item

}

}



[ctr	inc()	A
[ctr	ok()	A
[ctr	count()	B
[ctr	ok()	B
[ctr	dec()	B
[ctr	ok()	B
[q	deq()	B ←
[q	ok()	B ←
[q	eng(i)	A
[q	ok()	A

System with counter + queue.
- linearizable?

System with new adt remove ⇒
if empty return empty
else

Comparing Consistency Models

A stronger than B

For any history H: $H \vdash A \Rightarrow H \vdash B$
but not $H \vdash B \Rightarrow H \vdash A$

A equivalent to B

$H \vdash A \Leftrightarrow H \vdash B$

Linearizability

$$OP_1 <_H OP_2$$

$$\Leftrightarrow \text{ret } OP_1 < \text{inv } OP_2$$

Seq Cst

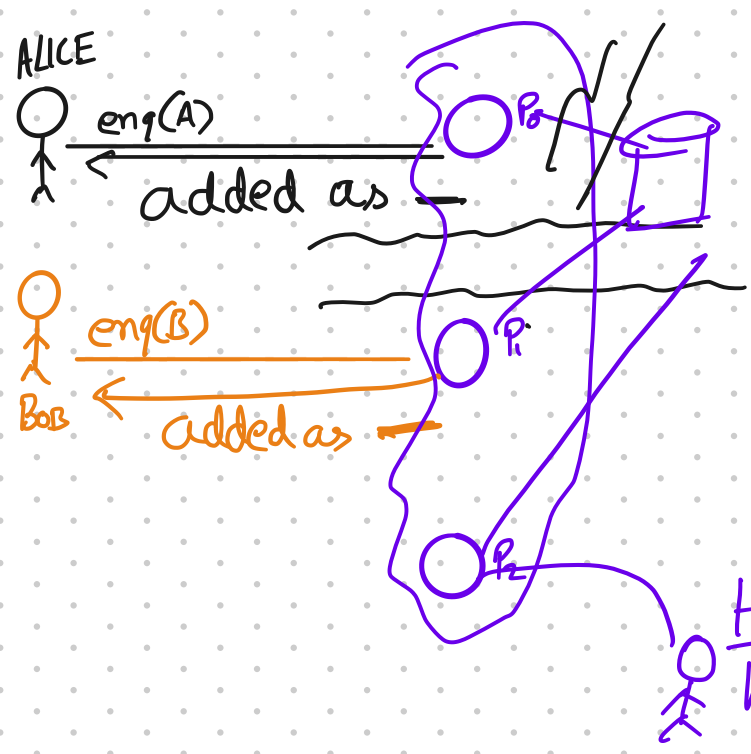
$$OP_1 <_S OP_2$$

$$\Leftrightarrow OP_1, OP_2 \text{ from same process}$$

$$\text{ret } OP_1 < \text{inv } OP_2$$

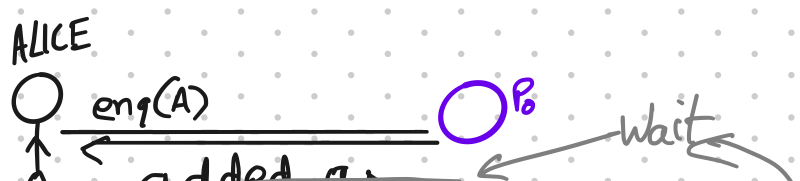
CAP

CORE MESSAGE: Linearizability comes at a cost

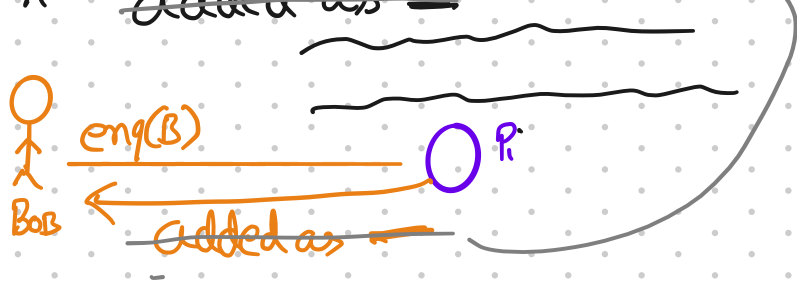


- To be linearizable
Need to order
enq(A), enq(B)
- To order, need to
know about
the other

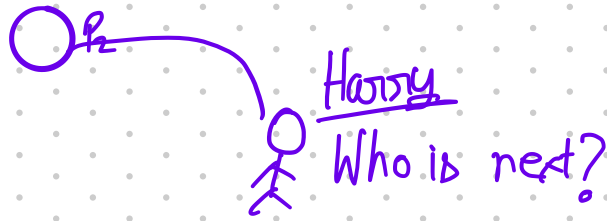
Choices



wait => Not available



[Not sure how long it will be]



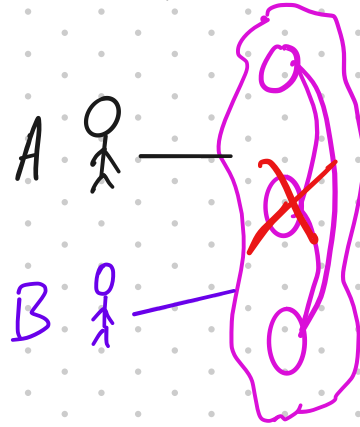
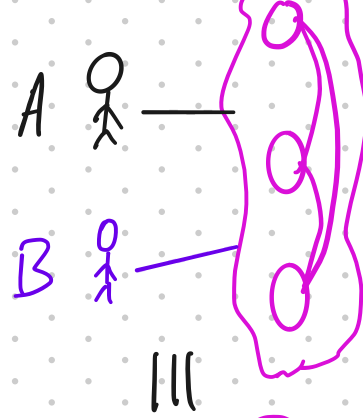
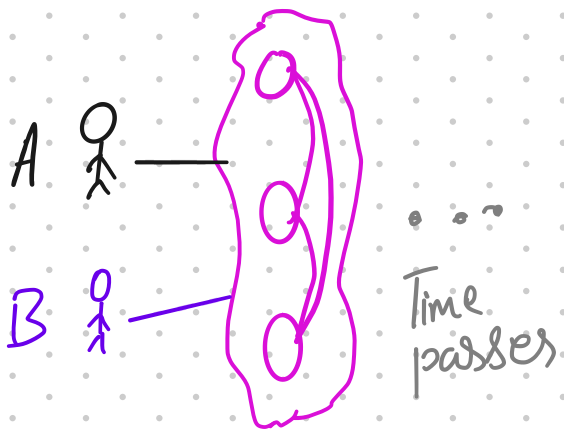
Give up on linearizability

Replicated State Machines

- PROBLEM: How To BUILD A FAULT TOLERANT ???

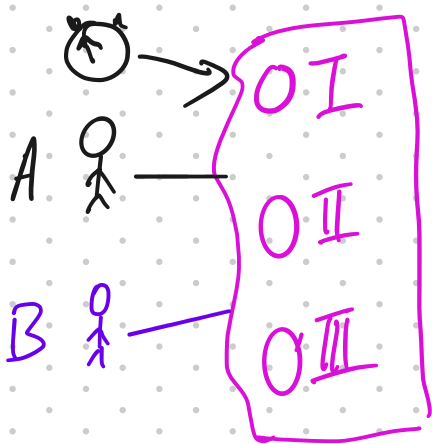
??? ≡ Airline reservation system (Lamport)
 Key value store (Chubby, Etd, ...)
 Database (...)
 Ice cream stand (...)

What we want



Note: Many possible solutions BUT

- Hard to design a correct one.
- So instead: identify a recipe.



Time →

Make sure I, II and III always have the same state & behave identically

same state

Behave identically : Deterministic



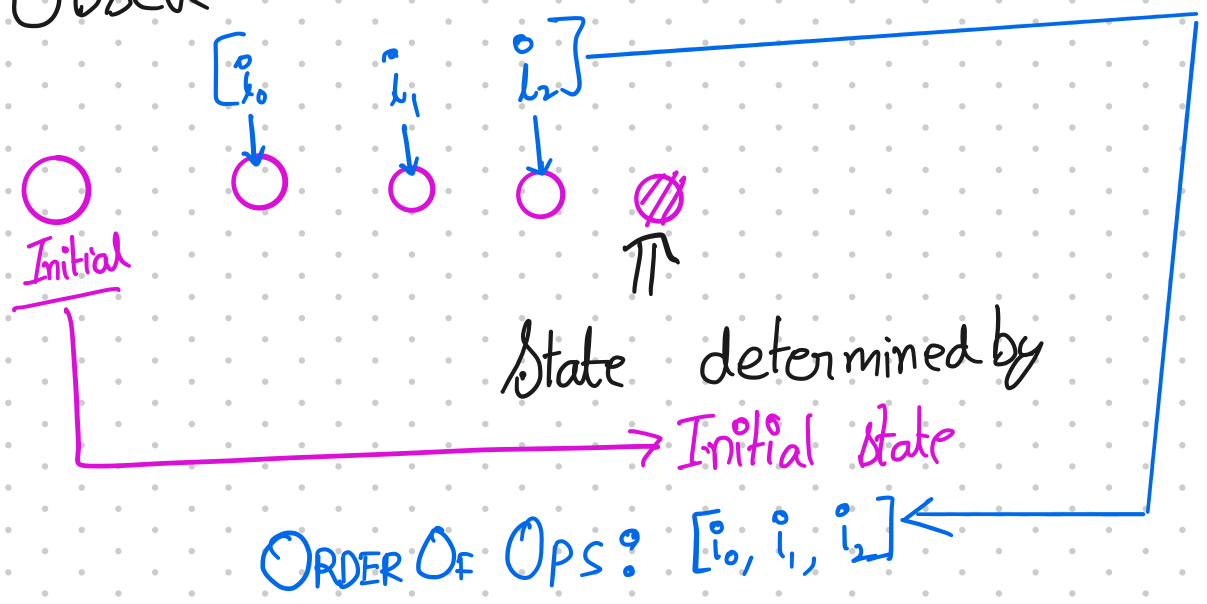
Out can depend on

- Current state
- In

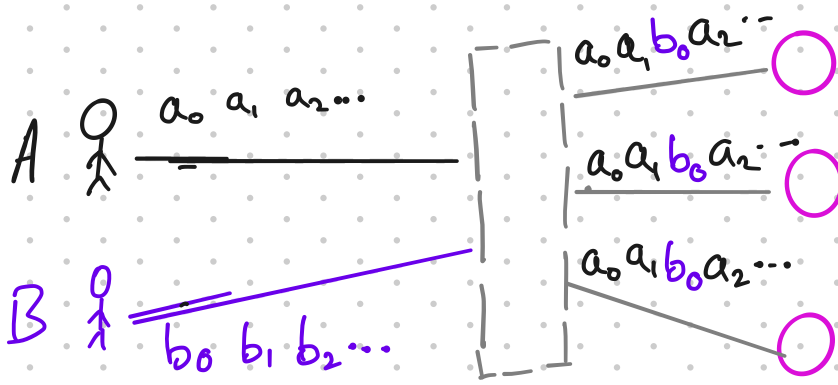
- No dependence on time
- No dependence on observation about the world
- No reading /dev/random
- ...

Assumed by nearly all fault tolerance mechanisms.

Observe



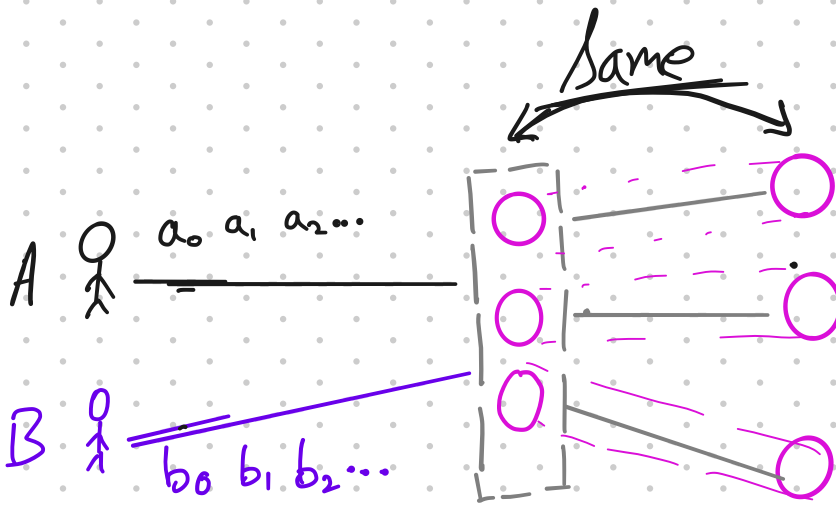
leads to RSMs



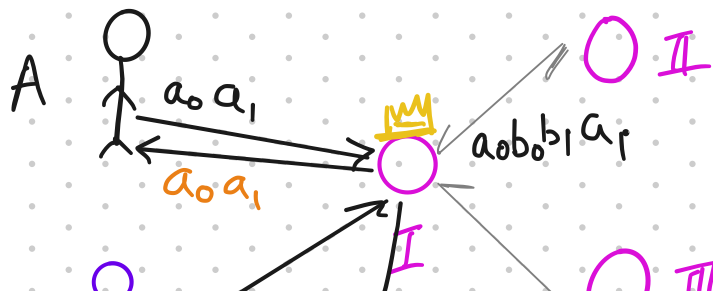
Two requirements
 Agreement: All correct processes get the same set of requests (commands)

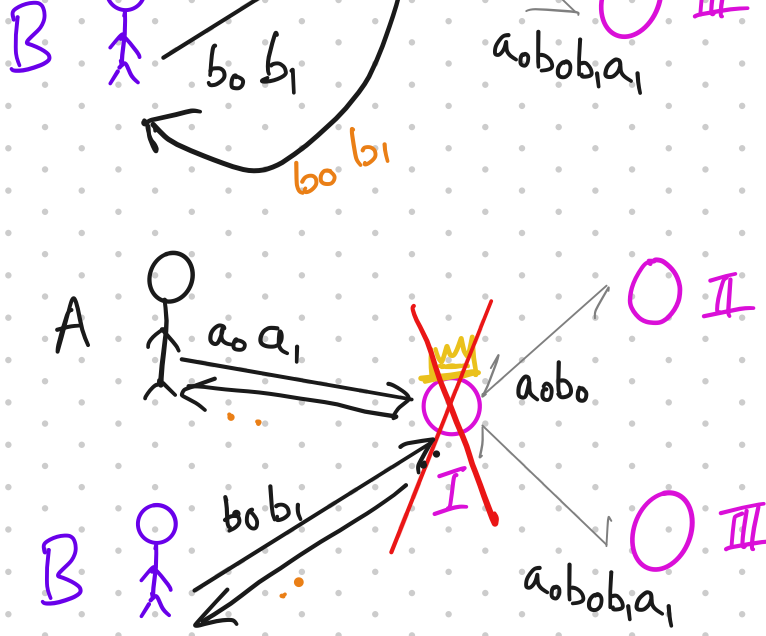
Ordering: All correct processes execute request in the same order.

The asynchronous model makes this hard!
 Formally show this later in class



"Leader" based

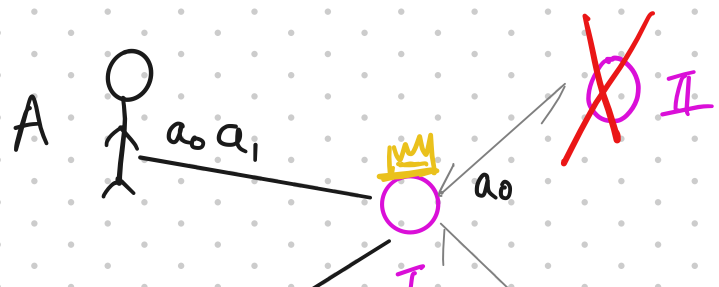


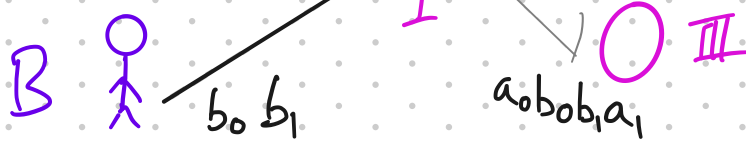


- How to detect leader has failed?
 - Timeout

- Who becomes the next leader?
 - Req: At most one active leader.

Agree on leader





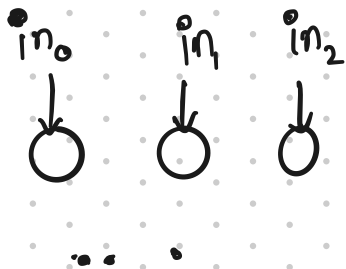
- When is it safe to execute request & return result?

↳ Commit Point

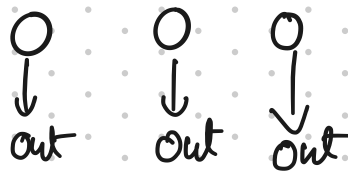
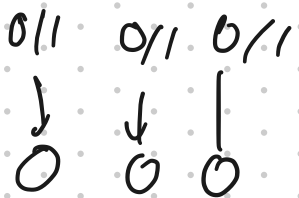
To recover from failures

When? Enough processes agree that they know command.

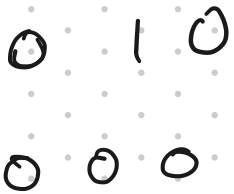
Agreement: Solved by CONSENSUS PROTOCOLS



Agreement: If process p & q output v_p & v_q then $v_p = v_q$



Validity: If process p outputs v_p then $v_p \in \{in_0, in_1, \dots\}$

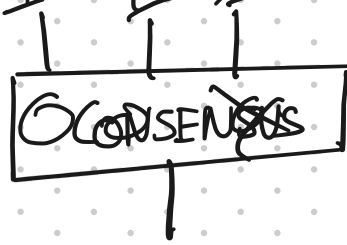


recv m do
return / one
end

Termination: Eventually ~~at~~ correct processes output a decision

- Leader Election





- Sufficiently Replicated

(idx, nref)

