

Validity

Initially $V = \{v_i\}$

To execute `propose(v_i)`

round $k, 1 \leq k \leq f+1$

- 1: send $\{v \in V : p_i \text{ has not already sent } v\}$ to all
- 2: for all $j, 0 \leq j \leq n-1, j \neq i$ do
- 3: receive S_j from p_j
- 4: $V := V \cup S_j$

`decide(x)` occurs as follows:

- 5: if $k = f+1$ then
- 6: decide $\min(V)$

Validity

Initially $V = \{v_i\}$

To execute `propose(v_i)`

round $k, 1 \leq k \leq f+1$

1: `send` $\{v \in V : p_i \text{ has not already sent } v\}$ to all

2: `for all` $j, 0 \leq j \leq n-1, j \neq i$ `do`

3: `receive` S_j from p_j

4: $V := V \cup S_j$

`decide(x)` occurs as follows:

5: `if` $k = f+1$ `then`

6: `decide` $\min(V)$

- Suppose every process proposes v^*
- Since only crash model, only v^* can be sent
- By Validity of send and receive, v^* sent by correct processes will be received by correct processes
- By Uniform Integrity of send and receive, only v^* can be received
- By protocol, $V = \{v^*\}$
- $\min(V) = v^*$
- `decide(v^*)`

Agreement

Initially $V = \{v_i\}$

To execute `propose(v_i)`

round k , $1 \leq k \leq f+1$

- 1: send $\{v \in V : p_i \text{ has not already sent } v\}$ to all
- 2: for all j , $0 \leq j \leq n-1$, $j \neq i$ do
- 3: receive S_j from p_j
- 4: $V := V \cup S_j$

`decide(x)` occurs as follows:

- 5: if $k = f+1$ then
- 6: decide $\min(V)$

Lemma 1

For any $r \geq 1$, if a process p receives a value v in round r , then there exists a sequence of processes p_0, p_1, \dots, p_r such that $p_r = p$, p_0 is v 's proponent, and in each round p_{k-1} sends v and p_k receives it. Furthermore, all processes in the sequence are distinct.

Proof

By induction on the length of the sequence

Agreement

Initially $V = \{v_i\}$

To execute **propose**(v_i)

round $k, 1 \leq k \leq f+1$

- 1: **send** $\{v \in V : p_i \text{ has not already sent } v\}$ to all
- 2: **for all** $j, 0 \leq j \leq n-1, j \neq i$ **do**
- 3: **receive** S_j from p_j
- 4: $V := V \cup S_j$

decide(x) occurs as follows:

- 5: **if** $k = f+1$ **then**
- 6: **decide** $\min(V)$

Lemma 2:

In every execution, at the end of round $f+1$,
 $V_i = V_j$ for every correct processes p_i and p_j

Agreement follows from Lemma 2, since
 \min is a deterministic function

Agreement

Initially $V = \{v_i\}$

To execute `propose(v_i)`

round k , $1 \leq k \leq f+1$

- 1: send $\{v \in V : p_i \text{ has not already sent } v\}$ to all
- 2: for all j , $0 \leq j \leq n-1$, $j \neq i$ do
- 3: receive S_j from p_j
- 4: $V := V \cup S_j$

`decide(x)` occurs as follows:

- 5: if $k = f+1$ then
- 6: decide $\min(V)$

Lemma 2:

In every execution, at the end of round $f+1$,
 $V_i = V_j$ for every correct processes p_i and p_j

Agreement follows from Lemma 2, since
`min` is a deterministic function

Proof:

- Show that if a correct p has x in its V at the end of round $f+1$, then every correct p has x in its V at the end of round $f+1$

Agreement

Initially $V = \{v_i\}$

To execute $\text{propose}(v_i)$

round $k, 1 \leq k \leq f+1$

- 1: send $\{v \in V : p_i \text{ has not already sent } v\}$ to all
- 2: for all $j, 0 \leq j \leq n-1, j \neq i$ do
- 3: receive S_j from p_j
- 4: $V := V \cup S_j$

$\text{decide}(x)$ occurs as follows:

- 5: if $k = f+1$ then
- 6: decide $\min(V)$

Lemma 2:

In every execution, at the end of round $f+1$, $V_i = V_j$ for every correct processes p_i and p_j

Agreement follows from Lemma 2, since \min is a deterministic function

Proof:

- Show that if a correct p has x in its V at the end of round $f+1$, then every correct p has x in its V at the end of round $f+1$
- Let r be earliest round x is added to the V of a correct p . Let that process be p^*
- If $r \leq f$, then p^* sends x in round $r+1 \leq f+1$; every correct process receives x and adds x to its V in round $r+1$

Agreement

Initially $V = \{v_i\}$

To execute $\text{propose}(v_i)$

round $k, 1 \leq k \leq f+1$

- 1: send $\{v \in V : p_i \text{ has not already sent } v\}$ to all
- 2: for all $j, 0 \leq j \leq n-1, j \neq i$ do
- 3: receive S_j from p_j
- 4: $V := V \cup S_j$

$\text{decide}(x)$ occurs as follows:

- 5: if $k = f+1$ then
- 6: decide $\min(V)$

Lemma 2:

In every execution, at the end of round $f+1$, $V_i = V_j$ for every correct processes p_i and p_j

Agreement follows from Lemma 2, since \min is a deterministic function

Proof:

- Show that if a correct p has x in its V at the end of round $f+1$, then every correct p has x in its V at the end of round $f+1$
- Let r be earliest round x is added to the V of a correct p . Let that process be p^*
- If $r \leq f$, then p^* sends x in round $r+1 \leq f+1$; every correct process receives x and adds x to its V in round $r+1$
- What if $r = f+1$?

Agreement

Initially $V = \{v_i\}$

To execute `propose(v_i)`

round k , $1 \leq k \leq f+1$

- 1: send $\{v \in V : p_i \text{ has not already sent } v\}$ to all
- 2: for all j , $0 \leq j \leq n-1$, $j \neq i$ do
- 3: receive S_j from p_j
- 4: $V := V \cup S_j$

`decide(x)` occurs as follows:

- 5: if $k = f+1$ then
- 6: decide $\min(V)$

Lemma 2:

In every execution, at the end of round $f+1$, $V_i = V_j$ for every correct processes p_i and p_j

Agreement follows from Lemma 2, since \min is a deterministic function

Proof:

- Show that if a correct p has x in its V at the end of round $f+1$, then every correct p has x in its V at the end of round $f+1$
- Let r be earliest round x is added to the V of a correct p . Let that process be p^*
- If $r \leq f$, then p^* sends x in round $r+1 \leq f+1$; every correct process receives x and adds x to its V in round $r+1$
- **What if $r = f+1$?**
- By Lemma 1, there exists a sequence of distinct processes $p_0, \dots, p_{f+1} = p^*$
- Consider processes p_0, \dots, p_f
- $f+1$ processes; only f faulty
- one of p_0, \dots, p_f is correct, and adds x to its V before p^* does it in round r

CONTRADICTION!

Terminating Reliable Broadcast

- Validity** If the sender is correct and broadcasts a message m , then all correct processes eventually deliver m
- Agreement** If a correct process delivers a message m , then all correct processes eventually deliver m
- Integrity** Every correct process delivers at most one message, and if it delivers $m \neq SF$, then some process must have broadcast m
- Termination** Every correct process eventually delivers some message

TRB for benign failures

Sender in round 1:

1: send m to all

Process p in round $k, 1 \leq k \leq f+1$

1: if delivered m in round $k-1$ then

2: if $p \neq \text{sender}$ then

3: send m to all

4: halt

5: receive round k messages

6: if received m then

7: deliver(m)

8: if $k = f+1$ then halt

9: else if $k = f+1$

10: deliver(SF)

11: halt

TRB for benign failures

Sender in round 1:

1: send m to all

Process p in round k , $1 \leq k \leq f+1$

1: if delivered m in round $k-1$ then

2: if $p \neq \text{sender}$ then

3: send m to all

4: halt

5: receive round k messages

6: if received m then

7: deliver(m)

8: if $k = f+1$ then halt

9: else if $k = f+1$

10: deliver(SF)

11: halt

Terminates in $f+1$ rounds

How can we do better?

Find a protocol whose round complexity is proportional to t – the number of failures that actually occurred – rather than to f – the max number of failures that may occur

Early stopping: the idea

- Suppose processes can detect the set of processes that have failed by the end of round i
- Call that set $faulty(p, i)$
- If $|faulty(p, i)| < i$ there can be no active dangerous chains, and p can safely deliver SF

Early Stopping: The Protocol

$faulty(p, k) \equiv$ set of processes that failed to send a message to p in some round

1: if $p = \text{sender}$ then value := m else value := ?

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then value := ?

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Termination

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then value := m else value := ?

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then value := ?

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Termination

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then $\text{value} := m$ else $\text{value} := ?$

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: $\text{value} := v$

8: deliver value

9: if $p = \text{sender}$ then $\text{value} := ?$

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: $\text{value} := \text{SF}$

12: deliver value

13: if $k = f+1$ then halt

- ◉ If in any round a process receives a value, then it delivers the value in that round
- ◉ If a process has received only "?" for $f+1$ rounds, then it delivers SF in round $f+1$

Validity

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then value := m else value := ?

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then value := ?

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Validity

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then value := m else value := ?

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then value := ?

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

- If the sender is correct then it sends m to all in round 1
- By Validity of the underlying send and receive, every correct process will receive m by the end of round 1
- By the protocol, every correct process will deliver m by the end of round 1

Agreement – 1

Lemma 1:

For any $r \geq 1$, if a process p delivers $m \neq \text{SF}$ in round r , then there exists a sequence of processes p_0, p_1, \dots, p_r such that $p_0 = \text{sender}$, $p_r = p$, and in each round $k, 1 \leq k \leq r$, p_{k-1} sent m and p_k received it. Furthermore, all processes in the sequence are distinct, unless $r = 1$ and $p_0 = p_1 = \text{sender}$

Lemma 2:

For any $r \geq 1$, if a process p sets value to SF in round r , then there exist some $j \leq r$ and a sequence of distinct processes $q_j, q_{j+1}, \dots, q_r = p$ such that q_j only receives “?” in rounds 1 to j , $|faulty(q_j, j)| < j$, and in each round $k, j+1 \leq k \leq r$, q_{k-1} sends SF to q_k and q_k receives SF

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then value := m else value := ?

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then value := ?

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Agreement - 2

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then value := m else value := ?

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then value := ?

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Lemma 3:

It is impossible for p and q , not necessarily correct or distinct, to set value in the same round r to m and SF, respectively

Agreement - 2

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then $\text{value} := m$ else $\text{value} := ?$

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then $\text{value} := ?$

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Lemma 3:

It is impossible for p and q , not necessarily correct or distinct, to set value in the same round r to m and SF, respectively

Proof

By contradiction

Suppose p sets value = m and q sets value = SF

By Lemmas 1 and 2 there exist

p_0, \dots, p_r

q_j, \dots, q_r

with the appropriate characteristics

Since q_j did not receive m from process p_{k-1} $1 \leq k \leq j$ in round k q_j must conclude that p_0, \dots, p_{j-1} are all faulty processes

But then, $|faulty(q_j, j)| \geq j$

CONTRADICTION

Agreement – 3

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then value := m else value := ?

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then value := ?

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Agreement – 3

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then value := m else value := ?

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then value := ?

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Proof

If no correct process ever receives m , then every correct process delivers SF in round $f+1$

Agreement – 3

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then $\text{value} := m$ else $\text{value} := ?$

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then $\text{value} := ?$

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Proof

If no correct process ever receives m , then every correct process delivers SF in round $f+1$

Otherwise, let r be the earliest round in which a correct process delivers value $\neq \text{SF}$

$r \leq f$

□ By Lemma 3, no (correct) process can set value differently in round r

□ In round $r+1 \leq f+1$, that correct process sends its value to all

□ Every correct process receives and delivers the value in round $r+1 \leq f+1$

$r = f+1$

□ By Lemma 1, there exists a sequence $p_0, \dots, p_{f+1} = p_r$ of distinct processes

□ Consider processes p_0, \dots, p_f

⦿ $f+1$ processes; only f faulty

⦿ one of p_0, \dots, p_f is correct-- let it be p_c

⦿ To send v in round $c+1$, p_c must have set its value to v and delivered v in round $c < r$

CONTRADICTION

Integrity

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then value := m else value := ?

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: value := v

8: deliver value

9: if $p = \text{sender}$ then value := ?

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: value := SF

12: deliver value

13: if $k = f+1$ then halt

Integrity

Let $faulty(p, k)$ be the set of processes that have failed to send a message to p in any round $1, \dots, k$

1: if $p = \text{sender}$ then $\text{value} := m$ else $\text{value} := ?$

Process p in round $k, 1 \leq k \leq f+1$

2: send value to all

3: if delivered in round $k-1$ then halt

4: receive round k values from all

5: $faulty(p, k) := faulty(p, k-1) \cup \{q \mid p \text{ received no value from } q \text{ in round } k\}$

6: if received value $v \neq ?$ then

7: $\text{value} := v$

8: deliver value

9: if $p = \text{sender}$ then $\text{value} := ?$

10: else if $k = f+1$ or $|faulty(p, k)| < k$ then

11: $\text{value} := \text{SF}$

12: deliver value

13: if $k = f+1$ then halt

👁 At most one m

□ Failures are benign, and a process executes at most one deliver event before halting

👁 If $m \neq \text{SF}$, only if m was broadcast

□ From Lemma 1 in the proof of Agreement

FLP

What about the asynchronous model?

Theorem

There is no deterministic protocol that solves Consensus in a message-passing asynchronous system in which at most one process may fail by crashing

(Fisher, Lynch, and Paterson. Impossibility of distributed consensus with one faulty process. JACM, Vol. 32, no. 2, April 1985, pp. 374-382)

Around FLP in 80 Slides

How can one get around FLP?

Weaken the problem

① Weaken termination

- use randomization to terminate with arbitrarily high probability
- guarantee termination only during periods of synchrony

② Weaken agreement

□ ϵ - agreement

- ▶ real-valued inputs and outputs
- ▶ agreement within real-valued small positive tolerance ϵ

□ k-set agreement

- ▶ **Agreement**: In any execution, there is a subset W of the set of input values, $|W| = k$, s.t. all decision values are in W
- ▶ **Validity**: In any execution, any decision value for any process is the input value of some process

How can one get around FLP?

Constrain input values

- ① Characterize the set of input values for which agreement is possible

Strengthen the system model

- ① Introduce **failure detectors** to distinguish between crashed processes and very slow processes

Paxos

It's (Almost) Everywhere!

Production use of Paxos [\[edit \]](#)

- The Petal project from DEC SRC was likely the first system to use Paxos, in this case for widely replicated global information (e.g., which machines are in the system).^[20]
- Google uses the Paxos algorithm in their Chubby [distributed lock service](#) in order to keep replicas consistent in case of failure. Chubby is used by [BigTable](#) which is now in production in Google Analytics and other products.
- The [Infinit](#) peer-to-peer file system relies on Paxos to maintain consistency among replicas while allowing for quorums to evolve in size.
- [Google Spanner](#) and Megastore use the Paxos algorithm internally.
- The [OpenReplica replication service](#) [↗](#) uses Paxos to maintain replicas for an open access system that enables users to create fault-tolerant objects. It provides high performance through concurrent rounds and flexibility through dynamic membership changes.
- IBM supposedly uses the Paxos algorithm in their [IBM SAN Volume Controller](#) product to implement a general purpose fault-tolerant virtual machine used to run the configuration and control components of the [storage virtualization](#) services offered by the cluster.^[citation needed]
- Microsoft uses Paxos in the [Autopilot cluster management service](#) [📄](#) from Bing.
- [WANdisco](#) have implemented Paxos within their DConE active-active replication technology.^[21]
- [XtreemFS](#) uses a Paxos-based [lease](#) negotiation algorithm for fault-tolerant and consistent replication of file data and metadata.^[22]
- [Heroku](#) uses [Doozerd](#) [↗](#) which implements Paxos for its consistent distributed data store.
- [Ceph](#) uses Paxos as part of the monitor processes to agree which OSDs are up and in the cluster.
- The [Clustrix](#) distributed SQL database uses Paxos for [distributed transaction resolution](#) [↗](#).
- [Neo4j](#) HA graph database implements Paxos, replacing [Apache ZooKeeper](#) from v1.9
- VMware NSX Controller uses Paxos-based algorithm within NSX Controller cluster.
- [Amazon Web Services](#) uses the Paxos algorithm extensively to power its platform.^[23]
- [Nutanix](#) implements the Paxos algorithm in Cassandra for [metadata](#) [↗](#).
- [Apache Mesos](#) uses Paxos algorithm for its [replicated log](#) [↗](#) coordination.
- [Windows Fabric](#) used by many of the [Azure](#) services make use of the paxos algorithm for replication between nodes in a cluster
- [Oracle NoSQL Database](#) leverages Paxos-based automated fail-over election process in the event of a master replica node failure to minimize downtime.^[24]

The Part-Time Parliament

- Parliament determines laws by passing sequence of numbered decrees
- Direct democracy: Citizens/Legislators leave and enter the chamber at arbitrary times
- No centralized records: each legislator carries a ledger recording the approved decrees



Government 101

- ① No two ledgers contain contradictory information
- ① If a majority of legislators are in the Chamber and no one enters or leaves the Chamber for a sufficiently long time, then
 - any decree proposed by a legislator is eventually passed
 - any passed decree appears on the ledger of every legislator

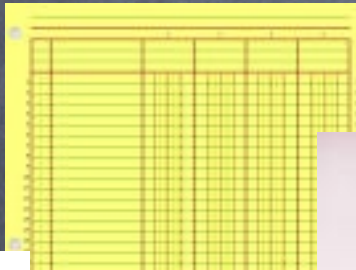
“In a world...”

Political intrigue!

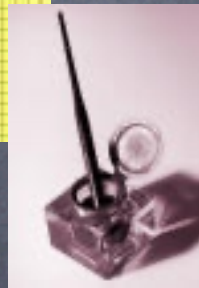
Funky equipment!



messengers!



ledger



pen with indelible ink



hourglass



Λαμπων

Δΐκστρα

Λισκωφ

Mysterious
characters

Σθνειδερ

Σκεεν

Πνυελι

Δωλεφ

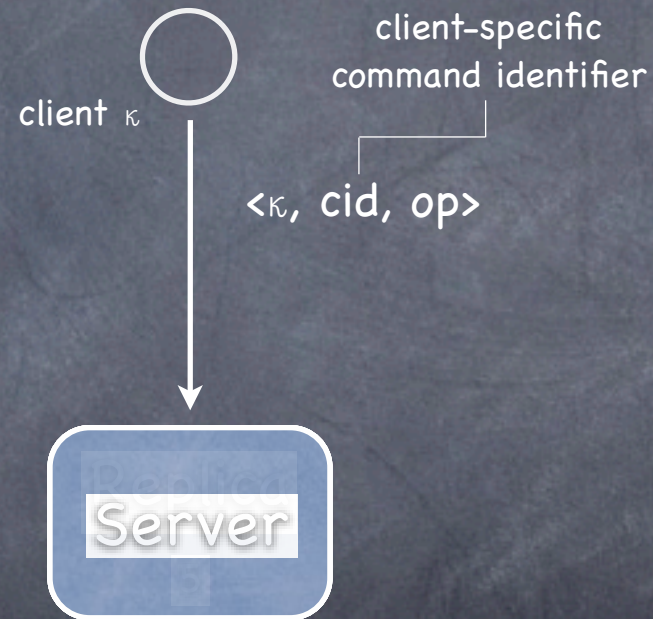
Δφωρκ



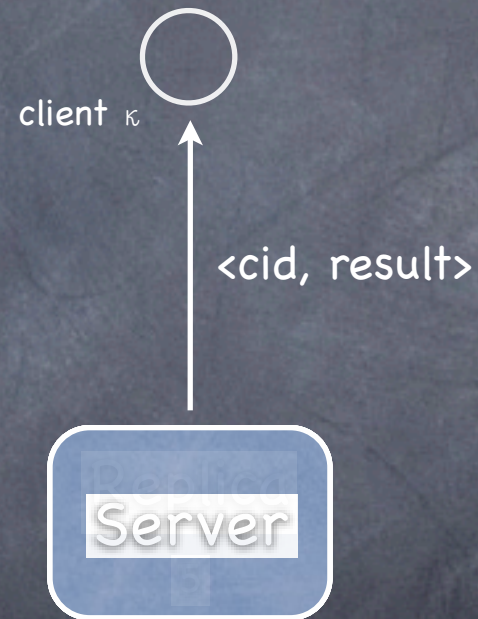
Back to the future

- 👁️ A protocol for state machine replication in an asynchronous environment that admits crash failures (key ideas already present in earlier work on *Viewstamped Replication* by Oki and Liskov)
- 👁️ Messages:
 - ❑ between correct endpoints are eventually received
 - ❑ can be lost and duplicated, but not corrupted

The big picture



The big picture



The big picture

client κ 

$f+1$

Replica₁

Replica₂

Replica₃

Replica₄

Replica₅



The big picture

client κ 

$f+1$

Replica₁

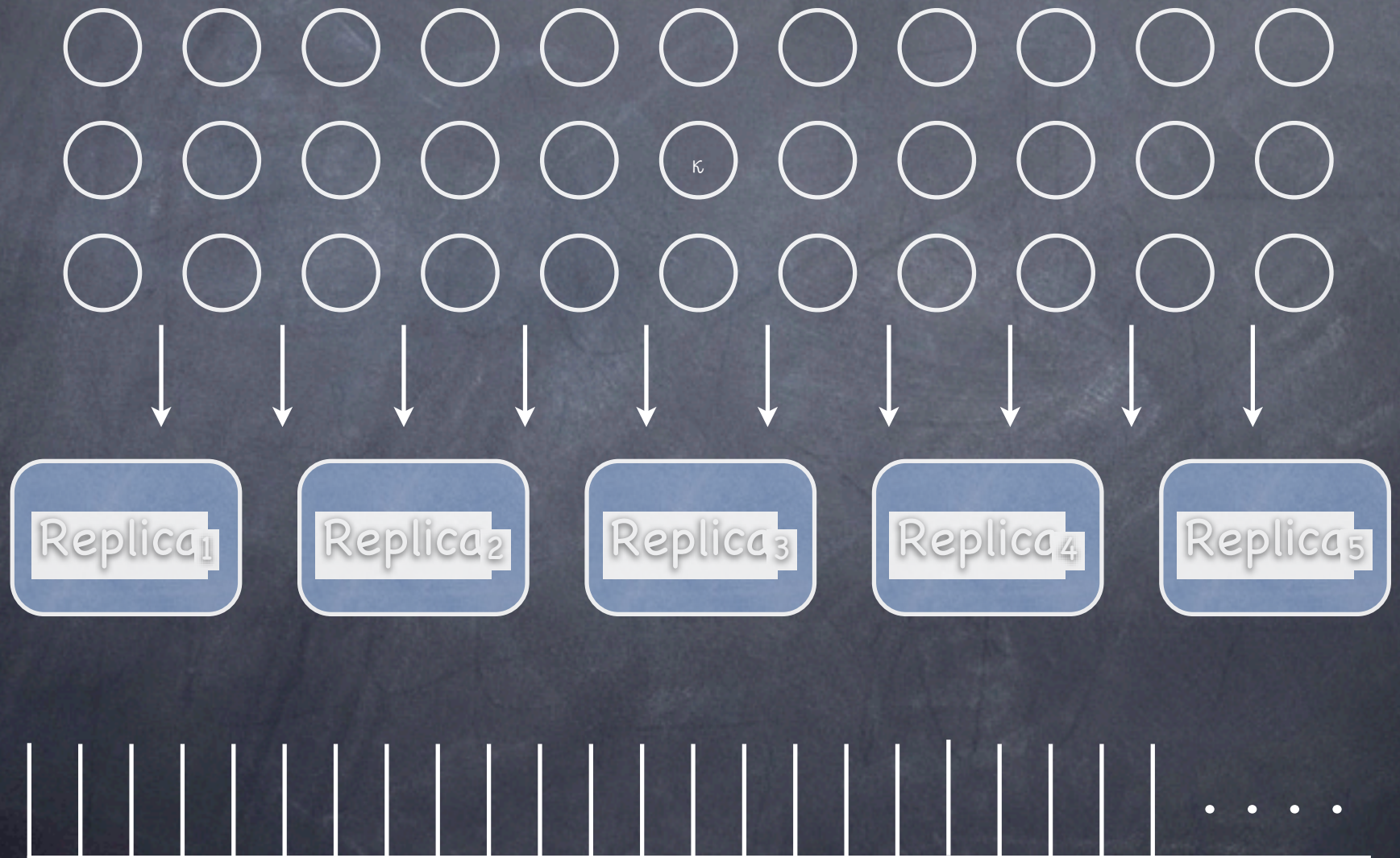
Replica₂

Replica₃

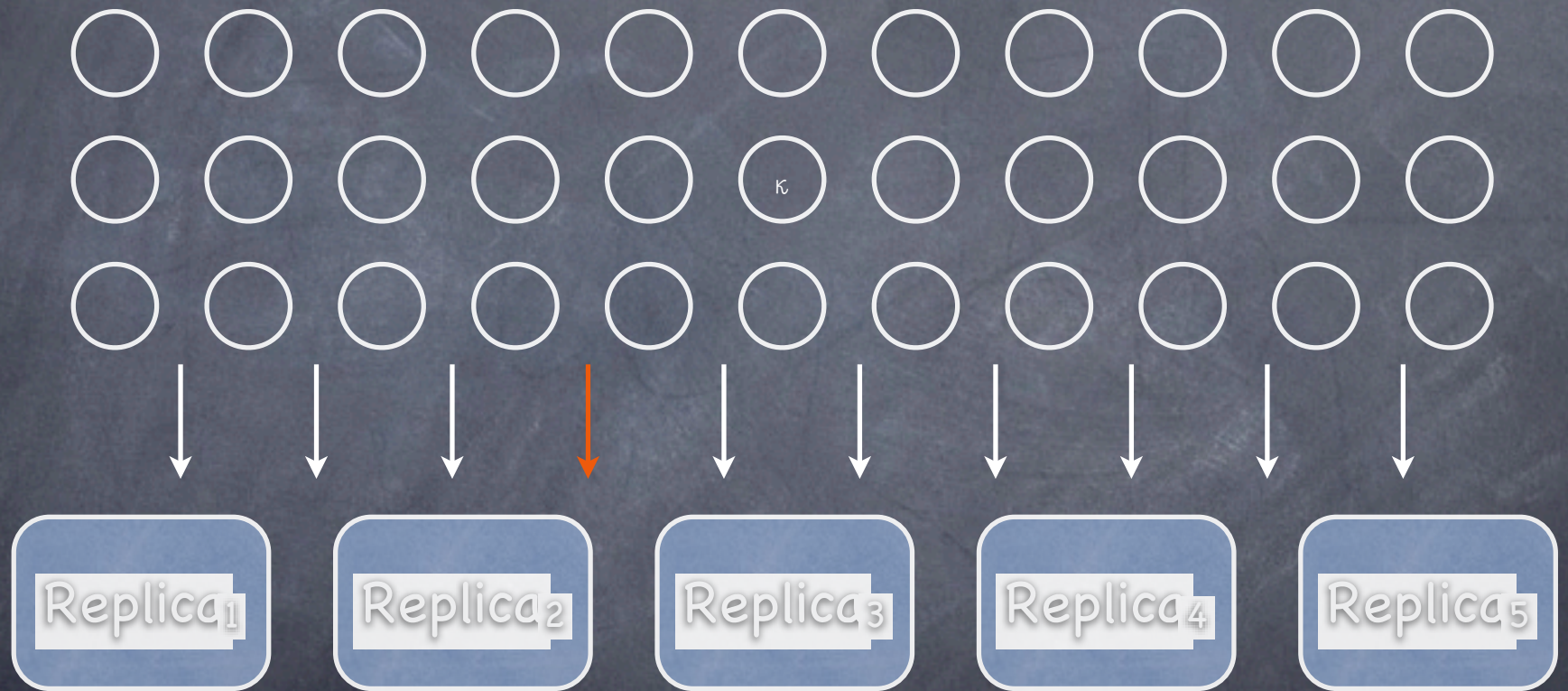
Replica₄

Replica₅

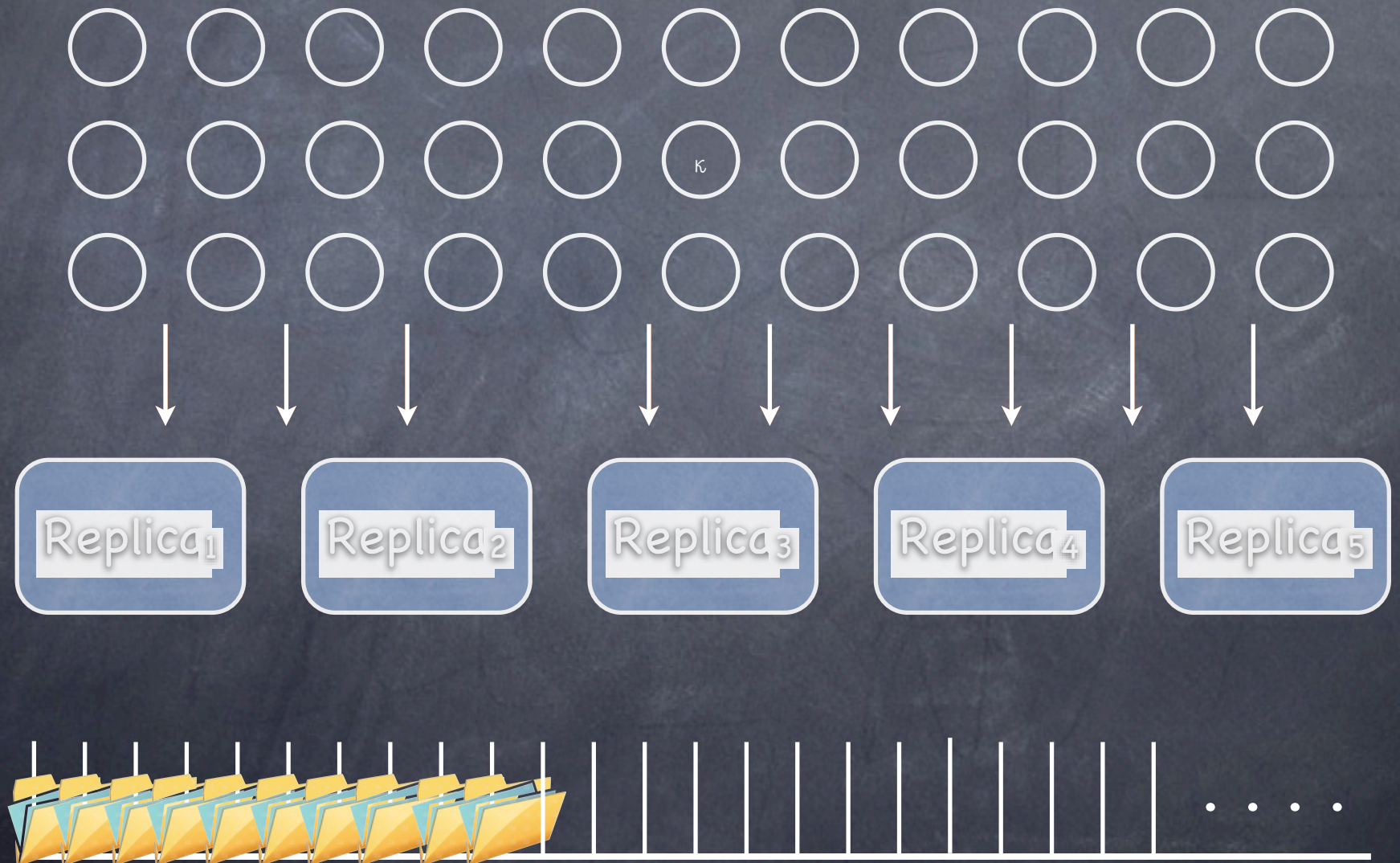
The big picture



The big picture



The big picture



The Cast

👁️ Replicas



👁️ Leaders



👁️ Acceptors



👁️ Commanders



👁️ Scout



Replicas

- ① Receive client requests
 - ① Propose command for lowest unused slot to **leaders**
 - ① Upon decision, execute commands in slot order
 - ① Return result to clients
 - ① **Not necessarily identical at any time!**
- ① Each replica ρ maintains four variables:
 - $\rho.state$: the application state
 - $\rho.slot_num$: next slot for which ρ does not know a decision
 - $\rho.proposals$: set of $\langle slot\ number, command \rangle$ pairs for past proposals
 - $\rho.decisions$: set of $\langle slot\ number, command \rangle$ pairs for decided slots

Replica

```
process Replica(leaders, initial_state)  
  var state := initial_state, slot_num := 1, proposals :=  $\emptyset$ ; decisions :=  $\emptyset$   
  for ever  
    switch receive()  
      case  $\langle$ request, p $\rangle$  :  
        propose(p);  
      case  $\langle$ decision, s, p $\rangle$  :  
        decisions := decisions  $\cup$  {(s, p)}  
        while  $\exists p' : \langle$ slot_num, p' $\rangle \in$  decisions do  
          if  $\exists p'' : \langle$ slot_num, p'' $\rangle \in$  proposals  $\wedge$  p''  $\neq$  p' then  
            propose(p'');  
          end if  
          perform(p');  
        end while  
      end switch  
    end for  
  end process
```

```
function perform( $\langle$  $\kappa$ , cid, op $\rangle$ )  
  if  $\exists s : s <$  slot_num  $\wedge$   
     $\langle$ s,  $\langle$  $\kappa$ , cid, op $\rangle \in$  decisions $\rangle \in$  decisions then  
    slot_num := slot_num + 1  
  else  
     $\langle$ next, result $\rangle$  := op(state);  
    atomic  
      state := next;  
      slot_num := slot_num + 1  
    end atomic  
    send( $\kappa$ ,  $\langle$ response, cid, result $\rangle$ );  
  end if  
end function
```

```
function propose(p)  
  if  $\nexists s : \langle$ s, p $\rangle \in$  decisions then  
    s' :=  $\min\{s \mid s \in \mathbb{N}^+ \wedge \nexists p' : \langle$ s, p' $\rangle \in$  proposal  $\cup$  decisions $\}$ ;  
    proposals := proposals  $\cup$  {(s', p)};  
     $\forall \lambda \in$  leaders : send( $\lambda$ ,  $\langle$ propose, s', p $\rangle$ );  
  end if  
end function
```

R4. For each ρ , the variable $\rho.slot_num$ never decreases

Replica

```
process Replica(leaders, initial_state)  
  var state := initial_state, slot_num := 1, proposals :=  $\emptyset$ ; decisions :=  $\emptyset$   
  for ever  
    switch receive()  
      case  $\langle$ request, p $\rangle$  :  
        propose(p);  
      case  $\langle$ decision, s, p $\rangle$  :  
        decisions := decisions  $\cup$   $\{(s, p)\}$   
        while  $\exists p' : \langle$ slot_num, p' $\rangle \in$  decisions do  
          if  $\exists p'' : \langle$ slot_num, p'' $\rangle \in$  proposals  $\wedge$  p''  $\neq$  p' then  
            propose(p'');  
          end if  
          perform(p');  
        end while  
      end switch  
    end for  
  end process
```

```
function perform( $\langle$  $\kappa$ , cid, op $\rangle$ )  
  if  $\exists s : s <$  slot_num  $\wedge$   
     $\langle$ s,  $\langle$  $\kappa$ , cid, op $\rangle$  $\rangle \in$  decisions then  
    slot_num := slot_num + 1  
  else  
     $\langle$ next, result $\rangle$  := op(state);  
    atomic  
      state := next;  
      slot_num := slot_num + 1  
    end atomic  
    send( $\kappa$ ,  $\langle$ response, cid, result $\rangle$ );  
  end if  
end function
```

```
function propose(p)  
  if  $\nexists s : \langle$ s, p $\rangle \in$  decisions then  
    s' :=  $\min\{s \mid s \in \mathbb{N}^+ \wedge \nexists p' : \langle$ s, p' $\rangle \in$  proposal  $\cup$  decisions $\}$ ;  
    proposals := proposals  $\cup$   $\{(s', p)\}$ ;  
     $\forall \lambda \in$  leaders : send( $\lambda$ ,  $\langle$ propose, s', p $\rangle$ );  
  end if  
end function
```

R3. For all replicas ρ , $\rho.state$ is the result of applying the operations in $\rho.decisions$ to the initial state in increasing order of slot number s :

$$1 \leq s < slot_num$$

Replica

```
process Replica(leaders, initial_state)  
  var state := initial_state, slot_num := 1, proposals :=  $\emptyset$ ; decisions :=  $\emptyset$   
  for ever  
    switch receive()  
      case  $\langle$ request, p $\rangle$  :  
        propose(p);  
      case  $\langle$ decision, s, p $\rangle$  :  
        decisions := decisions  $\cup$   $\{(s, p)\}$   
        while  $\exists p' : \langle$ slot_num, p' $\rangle \in$  decisions do  
          if  $\exists p'' : \langle$ slot_num, p'' $\rangle \in$  proposals  $\wedge$  p''  $\neq$  p' then  
            propose(p'');  
          end if  
          perform(p');  
        end while  
      end switch  
    end for  
  end process
```

```
function perform( $\langle$  $\kappa$ , cid, op $\rangle$ )  
  if  $\exists s : s <$  slot_num  $\wedge$   
     $\langle$ s,  $\langle$  $\kappa$ , cid, op $\rangle$  $\rangle \in$  decisions then  
    slot_num := slot_num + 1  
  else  
     $\langle$ next, result $\rangle$  := op(state);  
    atomic  
      state := next;  
      slot_num := slot_num + 1  
    end atomic  
    send( $\kappa$ ,  $\langle$ response, cid, result $\rangle$ );  
  end if  
end function
```

```
function propose(p)  
  if  $\nexists s : \langle$ s, p $\rangle \in$  decisions then  
    s' :=  $\min\{s \mid s \in \mathbb{N}^+ \wedge \nexists p' : \langle$ s, p' $\rangle \in$  proposal  $\cup$  decisions $\}$ ;  
    proposals := proposals  $\cup$   $\{(s', p)\}$ ;  
     $\forall \lambda \in$  leaders : send( $\lambda$ ,  $\langle$ propose, s', p $\rangle$ );  
  end if  
end function
```

R2. All commands up to *slot_num* are in the set of decisions

Replica

process *Replica*(*leaders*, *initial_state*)

var *state* := *initial_state*, *slot_num* := 1, *proposals* := \emptyset ; *decisions* := \emptyset

for ever

switch *receive*()

case $\langle \text{request}, p \rangle$:

propose(*p*);

case $\langle \text{decision}, s, p \rangle$:

decisions := *decisions* \cup $\{(s, p)\}$

while $\exists p' : \langle \text{slot_num}, p' \rangle \in \text{decisions}$ **do**

if $\exists p'' : \langle \text{slot_num}, p'' \rangle \in \text{proposals} \wedge p'' \neq p'$ **then**

propose(*p''*);

end if

perform(*p'*);

end while

end switch

end for

end process

function *perform*($\langle \kappa, \text{cid}, \text{op} \rangle$)

if $\exists s : s < \text{slot_num} \wedge$

$\langle s, \langle \kappa, \text{cid}, \text{op} \rangle \rangle \in \text{decisions}$ **then**

slot_num := *slot_num* + 1

else

$\langle \text{next}, \text{result} \rangle := \text{op}(\text{state});$

atomic

state := *next*;

slot_num := *slot_num* + 1

end atomic

send($\kappa, \langle \text{response}, \text{cid}, \text{result} \rangle$);

end if

end function

function *propose*(*p*)

if $\nexists s : \langle s, p \rangle \in \text{decisions}$ **then**

$s' := \min\{s \mid s \in \mathbb{N}^+ \wedge \nexists p' : \langle s, p' \rangle \in \text{proposal} \cup \text{decisions}\};$

proposals := *proposals* \cup $\{(s', p)\};$

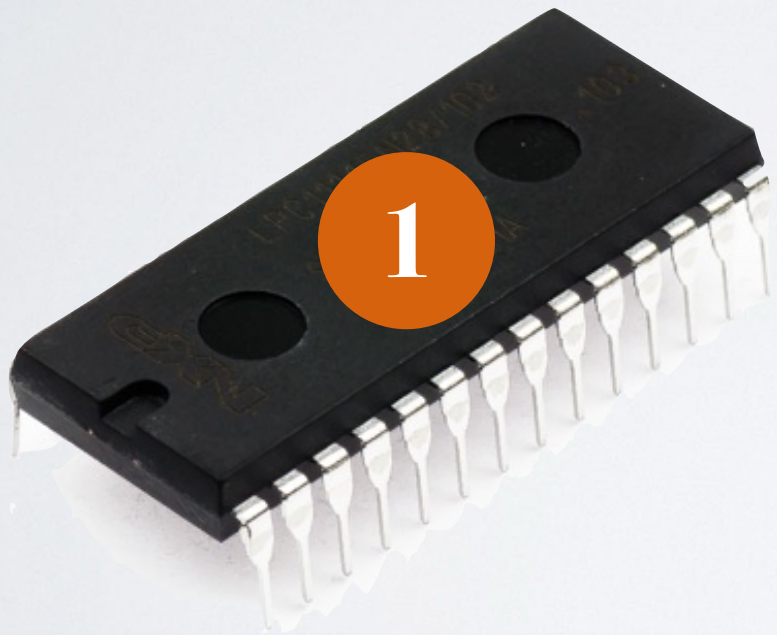
$\forall \lambda \in \text{leaders} : \text{send}(\lambda, \langle \text{propose}, s', p \rangle);$

end if

end function

R1. For any given slot, replicas decide the same command

THE BASIC IDEA: THE WRITE-ONCE REGISTER



Leaders compete to create a permanent mapping between slot numbers and proposals

Participants forever

read register

if register contains a mapping then

decide that mapping and halt

else

attempt to write own mapping