Paxos

Always safe

- Live during periods of synchrony
- Leader (primary) responsible for proposing the consensus value
- Features Dijkstra as a cheese inspector

Paxos

Always safe

- Live during periods of synchrony
- Leader (primary) responsible for proposing the consensus value
- Features Dijkstra as a cheese inspector

💿 Somewhat popular...

- D Part-time Parliament [L98]
- 🗇 Frangipani [TML97]
- Byzantine Paxos [CL99]
- Disk Paxos [GL00]
- Deconstructing Paxos [BDFG01]
- □ Reconstructing Paxos [BDFG01]
- Active Disk Paxos [CM02]
- Separating Agreement & Execution [YMAD 03]
- Byzantine Disk Paxos [ACKM04]
- □ Fast Byzantine Paxos [MA05]
- Fast Paxos [L05]
- Hybrid Quorums [CMLRS06]
- 🗅 Chubby [806]
- Paxos Register [LCAA07]
- 🛛 Zyzzyva [KADCW07]
 -

The Game of Paxos

Processes are competing to write a value in a write-once register

To learn the final value:

- 1. Push the read button and examine the token that falls into the tray
- 2. If the token is green, GAME OVER the final value of the register is stamped on the token!
- If the token is red and stamped with a value, place the token in the slot, set the dial to the same value, and push the write button



The Game of Paxos

Processes are competing to write a value in a write-once register

To learn the final value:

- 1. Push the read button and examine the token that falls into the tray
- 2. If the token is green, GAME OVER the final value of the register is stamped on the token!
- If the token is red and stamped with a value, place the token in the slot, set the dial to the same value, and push the write button
- 4. If the token is red and not stamped, place the token in the slot, set the dial to any value, and push the write button



Quorum Systems

Given a set of servers $\mathcal{U}, |\mathcal{U}| = n$ a quorum system is a set $\mathcal{Q} \subseteq 2^{\mathcal{U}}$ such that $\forall Q_1, Q_2 \in \mathcal{Q} : Q_1 \cap Q_2 \neq \emptyset$

Each Q in Q is a quorum

A R/W Register



store at each server a (v, ts) pair

store at each server

a (v, ts) pair

A R/W Register



Write(x, d)

- \square Ask servers in some Q for their ts
- $\Box \quad \mathsf{Set}\, ts_c > \max(\{ts\} \cup \mathsf{any previous } ts_c)$
- \square Update some Q' with (d, ts_c)

A R/W Register



Write(x, d)

- $\mathsf{Read}(x)$ \Box Ask servers in some Q for their ts \Box Ask servers in some Q for their (v,ts)
- \Box Set $ts_c > \max(\{ts\} \cup any \text{ previous } ts_c h$ Select most recent (v, ts)
- \Box Update some Q' with (d, ts_c)

System Model

O Universe U of servers, |U| = n

Byzantine faulty servers

 \square modeled as a non-empty fail-prone system $\mathcal{B} \subseteq 2^U$ \sqcap no $B \in \mathcal{B}$ is contained in another \Box some $B \in \mathcal{B}$ contains all faulty servers Olients are correct (can be weakened) Ø Point-to-point authenticated and reliable channels

> A correct process q receives a message from another correct process p if and only if p sent it

Masking Quorum System

[Malkhi and Reiter, 1998]

A quorum system Q is a masking quorum system for a fail-prone system \mathcal{B} if the following properties hold:

M-Consistency $\forall Q_1, Q_2 \in \mathcal{Q} \; \forall B_1, B_2 \in \mathcal{B} : (Q_1 \cap Q_2) \setminus B_1 \not\subseteq B_2$

M-Availability $\forall B \in \mathcal{B} \exists Q \in \mathcal{Q} : B \cap Q = \emptyset$

Dissemination Quorum System

A masking quorum system for self-verifying data client can detect modification by faulty server

D-Consistency

 $\forall Q_1, Q_2 \in \mathcal{Q} \; \forall B \in \mathcal{B} : (Q_1 \cap Q_2) \not\subseteq B$

D-Availability

 $\forall B \in \mathcal{B} \exists Q \in \mathcal{Q} : B \cap Q = \emptyset$

f-threshold Masking Quorum Systems $\forall Q_1, Q_2 \in \mathcal{Q} : |Q_1 \cap Q_2| \ge 2f + 1$ $\mathcal{Q} = \left\{ Q \subseteq U : |Q| = \left\lceil \frac{n+2f+1}{2} \right\rceil \right\} \qquad \mathcal{Q} = \left\{ Q \subseteq U : |Q| = \left\lceil \frac{n+f+1}{2} \right\rceil \right\}$ $\frac{n}{n \ge 4f + 1}$ $n \geq 3f + 1$

A safe read/write protocol

Client c executes:

Write(d)

- \rightarrow Ask all servers for their current timestamp t
- ← Wait for answer from |Q| different servers Set ts_c > max({t} ∪ any previous ts_c)
- \rightarrow Send (d,ts_c) to all servers
- ← Wait for |Q| acknowledgments

Read()

- \rightarrow Ask all servers for latest value/timestamp pair
- ← Wait for answer from |Q| different servers
- Select most recent (v,ts) for which at detailed the provers agree (if aby)



A simple observation

Oclient c (with current threshold f) executes:

Write(d)

- \rightarrow Ask all servers for their current timestamp t
- ← Wait for answer from |Q| different servers Set ts_c > max({t} ∪ any previous ts_c)
- \rightarrow Send (d,ts_c) to all servers
- ← wait for Rel acknowledgements

Read()

- → Ask all servers for latest value/timestamp pair
- ← Wait for answer from |Q| different servers Select most recent (v,ts) for which at least f + 1 answers agree (if any)

(Asynchronous) Authenticated Reliable channels

A correct process qreceives a message from another correct process p if and only if p sent it

A-Masking Quorum Systems

A quorum system Q is an a-masking quorum system for a fail-prone system B if the following properties hold for Q_r and Q_w :

AM-Consistency

 $\forall Q_r \in \mathcal{Q}_r \; \forall Q_w \in \mathcal{Q}_w \; \forall B_1, B_2 \in \mathcal{B} \\ (Q_r \cap Q_w) \setminus B_1 \not\subseteq B_2:$

AM-Availability

 $\forall B \in \mathcal{B} \exists Q_r \in \mathcal{Q}_r : B \cap Q_r = \emptyset$

Tradeoffs

	best known n	confirmable	non-confirmable
at the second second	self-verifying	$3f\!+\!1$	2f + 1
	generic	$4f\!+\!1$	$3f\!+\!1$

	radeott	S
best known n	confirmable	non-confirmable
self-verifying and generic	$3f\!+\!1$	2f + 1

Lower bound: never two rows again!

PBFT: A Byzantine Renaissance

@ Practical Byzantine Fault-Tolerance (CL99, CL00)

- 🗖 first to be safe in asynchronous systems
- 🗖 live under weak synchrony assumptions –Byzantine Paxos!
 - ust PBFT uses MACs instead of public key cryptography
- uses proactive recovery to tolerate more failures over system lifetime: now need no more than f failures in a "window"

BASE (RCL 01)

uses abstraction to reduce correlated faults



The General Idea

Primary

- Ø Primary-backup + quorum system
 - 🗋 executions are sequences of views
 - □ clients send signed commands to primary of current view
 - primary assigns sequence number to client's command
 - primary writes sequence number to the register implemented by the quorum system defined by all the servers (primary included)

What could possibly go wrong? 😲

The Primary could be faulty!

- > could ignore commands; assign same sequence number to different requests; skip sequence numbers; etc
- Backups monitor primary's behavior and trigger view changes to replace faulty primary
- Backups could be faulty!
 - > could incorrectly store commands forwarded by a correct primary
 - 🗅 use dissemination Byzantine quorum systems [MR98]
- Taulty replicas could incorrectly respond to the client!

What could possibly go wrong? 😲

- The Primary could be faulty!
 - > could ignore commands; assign same sequence number to different requests; skip sequence numbers; etc
 - Backups monitor primary's behavior and trigger view changes to replace faulty primary
- Backups could be faulty!
 - > could incorrectly store commands forwarded by a correct primary
 - 🗋 use dissemination Byzantine quorum systems [MR98]
- Taulty replicas could incorrectly respond to the client!
 - \Box Client waits for f+1 matching replies before accepting response

What could possibly go wrong? 😲

- The Primary could be faulty!
 - > could ignore commands; assign same sequence number to different requests; skip sequence numbers; etc
 - Backups monitor primary's behavior and trigger view changes to replace faulty primary
- Backups could be faulty!
 - > could incorrectly store commands forwarded by a correct primary
 - 🛭 use dissemination Byzantine quorum systems [MR98]
- Taulty replicas could incorrectly respond to the client!
 - \Box Client waits for f+1 matching replies before accepting response
- Carla Bruni could start singing!

Me, or your lying eyes?

- Algorithm steps are justified by certificates
 - Sets (quorums) of signed messages from distinct replicas proving that a property of interest holds
- **\textcircled{O}** With quorums of size at least 2f+1
 - Any two quorums intersect in at least one correct replica
 - Always one quorum contains only non-faulty replicas

PBFT: The site map

Solution

□ How the protocol works in the absence of failures – hopefully, the common case

View changes

 $\hfill\square$ How to depose a faulty primary and elect a new one

Garbage collection

 $\hfill\square$ How to reclaim the storage used to keep certificates

Recovery

 \Box How to make a faulty replica behave correctly again

Normal Operation

Three phases: Pre-prepare assigns sequence number to request ensures fault-tolerant consistent ordering of requests within views Commit ensures fault-tolerant consistent ordering of requests across views

- O Each replica *i* maintains the following state:
 - 🗆 Service state
 - □ A message log with all messages sent or received
 - \square An integer representing *i*'s current view

<section-header> Client issues request verequest, o,t, c>oc Primary Backup 1 Backup 2 Backup 3

Client issues request request of the operation primary Backup 1 Backup 2 Backup 3

Client	issues	request
$(REQUEST, o, t, c) = \sigma_c \\ f \\ timestamp$		
Primary	section of	
Backup 1		
Backup 2		
Backup 3		

<pre> </pre>	${\tt UEST,} o,t,c >_{\sigma_c} f$ client id		
Primary			
Backup 1			
Backup 2			
Backup 3			

Cl	ient	issues	reques	5†
REQUE	$ST_{,o,t,c} arrow_{\sigma_c}$			
Primary				
Backup 1				
Backup 2				
Backup 3				

Pre-prepare		
	Primary multicasts < <pre-prepare,<math>v,n,d > \sigma_p, m > 0</pre-prepare,<math>	
Primary		
Backup 1		
Backup 2		
Backup 3		

	Primary multicasts < <pre-prepare,<math>v,n,d > \sigma_p,m ></pre-prepare,<math>	
Primary		
Backup 1		
Backup 2		
Backup 3		

Pre-	bre	pare
		Pare

Primary multicasts «PRE-PREPARE, v, n, d > op, m>

	Primary multicas	ts < <pre_prepa< th=""><th>client's r $\langle RE, v, n, d > \sigma_{\pi}, m >$</th><th>request</th><th></th></pre_prepa<>	client's r $\langle RE, v, n, d > \sigma_{\pi}, m >$	request	
Primary			, , , , , , , , , , , , , , , , , , ,		
3ackup 1					
Backup 2					

Pre-prepare
Primary multicasts < <pre-prepare,<math>v,n,d \neq \sigma_p,m></pre-prepare,<math>
Primary
Backup 1
Backup 2
Backup 3



Pre-prepare



Each accepted PRE-PREPARE message is stored in the accepting replica's message log (including the Primary's)





Prepare Certificate

O P-certificates ensure total order within views



Prepare Certificate

O P-certificates ensure total order within views

prepare

- Replica produces P-certificate(m,v,n) iff its log holds: \Box The request m
 - \square A PRE-PREPARE for m in view v with sequence number n \square 2f PREPARE from different backups that match the pre-

Prepare Certificate

- P-certificates ensure total order within views
- Replica produces P-certificate(m,v,n) iff its log holds:
 In the request m
 - \square A PRE-PREPARE for m in view v with sequence number n
 - □ 2f **PREPARE** from different backups that match the preprepare
- $\textcircled{\sc original}$ A P-certificate (m,v,n) means that a quorum agrees with assigning sequence number n to m in view v
 - \square NO two non-faulty replicas with P-certificate (m_1,v,n) and P-certificate (m_2,v,n)

P-certificates are not enough

- A P-certificate proves that a majority of correct replicas has agreed on a sequence number for a client's request
- Set that order could be modified by a new leader elected in a view change



Commit Certificate

- C-certificates ensure total order across views
 a can't miss P-certificate during a view change
- O A replica has a C-certificate (m,v,n) if:
 - \square it had a P-certificate (m,v,n)
 - \Box log contains 2f+1 matching COMMIT from different replicas (including itself)
- Replica executes a request after it gets Ccertificate for it, and has cleared all requests with smaller sequence numbers



Aux armes les backups!

- A disgruntled backup mutinies:
 - stops accepting messages (but for VIEW-CHANGE & NEW-VIEW)
 - \square multicasts <VIEW-CHANGE, $v+1, \mathcal{P} >_{\sigma_s}$
 - $\square \mathcal{P}$ contains all P-Certificates known to replica i
- A backup joins mutiny after seeing f+1 distinct VIEW-CHANGE messages
- The Mutiny succeeds if new primary collects a new-view certificate \mathcal{V} , indicating support from 2f+1 distinct replicas (including itself)

On to view v+1: the new primary

- The "primary elect" \hat{p} (replica $v+1 \mod N$) extracts from the new-view certificate \mathcal{V} :
 - \Box the highest sequence number h of any message for which \mathcal{V} contains a P-certificate

On to view v+1: the new primary

- The "primary elect" \hat{p} (replica $v+1 \mod N$) extracts from the new-view certificate \mathcal{V} :
 - \Box the highest sequence number h of any message for which \mathcal{V} contains a P-certificate

On to view v+1: the new primary

- The "primary elect" \hat{p} (replica $v+1 \mod N$) extracts from the new-view certificate \mathcal{V} :
 - \Box the highest sequence number h of any message for which \mathcal{V} contains a P-certificate
 - \square two sets \mathcal{O} and \mathcal{N} :
 - ▶ If there is a P-certificate for n,m in \mathcal{V} , $n \leq h$ $\mathcal{O} = \mathcal{O} \cup \langle \mathsf{PRE}-\mathsf{PREPARE}, v+1, n, m \rangle_{\sigma_{\hat{\sigma}}}$
 - ▷ Otherwise, if $n \le h$ but no P-certificate: $\mathcal{N} = \mathcal{N} \cup \langle \mathsf{PRE}-\mathsf{PREPARE}, v+1, n, null \rangle_{\sigma_s}$

On to view v+1: the new primary

- The "primary elect" \hat{p} (replica $v+1 \mod N$) extracts from the new-view certificate \mathcal{V} :
 - \Box the highest sequence number h of any message for which $\mathcal V$ contains a P-certificate
 - \square two sets \mathcal{O} and \mathcal{N} :
 - ▶ If there is a P-certificate for n,m in \mathcal{V} , $n \le h$ $\mathcal{O} = \mathcal{O} \cup \langle \mathsf{PRE}-\mathsf{PREPARE}, v+1, n, m \rangle_{\sigma_{\hat{n}}}$
 - ▷ Otherwise, if $n \le h$ but no P-certificate: $\mathcal{N} = \mathcal{N} \cup \langle \mathsf{PRE}-\mathsf{PREPARE}, v+1, n, null \rangle_{\sigma_{\pi}}$

 $\mathfrak{O} \ \hat{p}$ multicasts <NEW-VIEW, $v+1, \mathcal{V}, \mathcal{O}, \mathcal{N} \succ_{\sigma_{\hat{n}}}$

On to view v+1: the backup

- **3** Backup accepts NEW-VIEW message for v+1 if
 - \square it is signed properly
 - \square it contains in $\mathcal V$ a valid VIEW-CHANGE messages for v+1
 - □ it can verify locally that *O* is correct (repeating the primary's computation)
- **a** Adds all entries in \mathcal{O} to its log (so did \hat{p} !)
- ${\it I}$ Multicasts a PREPARE for each message in ${\cal O}$
- Adds all PREPARE to log and enters new view

Garbage Collection

- For safety, a correct replica keeps in log messages about request o until it
 - o has been executed by a majority of correct replicas, and
 - 🗅 this fact can proven during a view change
- Truncate log with Certificate
 - □ Each replica *i* periodically (after processing *k* requests) checkpoints state and multicasts <CHECKPOINT,*n*,*d*,*i*>

Garbage Collection

- For safety, a correct replica keeps in log messages about request o until it
 - o has been executed by a majority of correct replicas, and
 - \Box this fact can proven during a view change
- Truncate log with Certificate

last executed request reflected in state

□ Each replica *i* periodically (after processing *k* requests) checkpoints state and multicasts <CHECKPOINT,*n*,*d*,*i*>

Garbage Collection

- For safety, a correct replica keeps in log messages about request o until it
 - □ o has been executed by a majority of correct replicas, and
 - □ this fact can proven during a view change
- Truncate log with Certificate
 - □ Each replica *i* periodically (after processing *k* requests) checkpoints state and multicasts <CHECKPOINT,*n*,*d*,*i*>

state's digest

Garbage Collection

- For safety, a correct replica keeps in log messages about request o until it
 - o has been executed by a majority of correct replicas, and
 - \Box this fact can proven during a view change
- Truncate log with Stable Certificate
 - □ Each replica *i* periodically (after processing *k* requests) checkpoints state and multicasts <CHECKPOINT,*n*,*d*,*i*>
 - □ 2f+1 CHECKPOINT messages are a proof of the checkpoint's correctness

View change, revisited

Ø A disgruntled backup multicasts

 $\langle VIEW-CHANGE, v+1, n, s, C, P, i \rangle_{\sigma_i}$

View change, revisited

A disgruntled backup multicasts

 $\langle VIEW-CHANGE, v+1, n, s, C, P, i \rangle_{\sigma_i}$

sequence number of last stable checkpoint

View change, revisited

A disgruntled backup multicasts

 $\textbf{<\!VIEW-CHANGE}, v \! + \! 1, \! n, \! s, \! \mathcal{C}, \! \mathcal{P}, \! i \! >_{\!\!\sigma_i} \\ {}_{\!\!\text{last stable checkpoint}}$

View change, revisited

A disgruntled backup multicasts

<VIEW-CHANGE, $v + 1, n, s, \mathcal{C}, \mathcal{P}, i >_{\sigma_i}$ stable certificate for s

View change, revisited

Ø A disgruntled backup multicasts

 $\langle VIEW-CHANGE, v+1, n, s, \overline{C, \mathcal{P}, i} \rangle_{\sigma_i}$

P certificates for requests with sequence number > n

View change, revisited

A disgruntled backup multicasts

<VIEW-CHANGE, $v + 1, n, s, C, P, i >_{\sigma_i}$

 ${\it O}$ \hat{p} multicasts

Citius, Altius, Fortius: Towards deployable BFT

Reducing the costs of BFT replication

- Addressing confidentiality
- Reducing complexity

Reducing the costs of BFT replication

- Who cares? Machines are cheap...
 - Replicas should fail independently in software, not just hardware
 - How many independently failing implementations of non-trivial services do actually exist?

Back the old conundrum





A: voter and client share fate!









Rethinking State Machine Replication

Not Agreement + Order

but rather Agreement on Order + Execution

Rethinking State Machine Replication

Not Agreement + Order but rather Agreement on Order + Execution Benefits: 25+1 @3541 state machine replicas

Rethinking State Machine Replication

Not Agreement + Order but rather Agreement on Order + Execution Benefits: 25+1 35441 state machine replicas Melps Replication Marks confidentiality



Separation enables confidentiality

Three design principles:



 \bigcirc

Separation enables confidentiality

Three design principles:

- 1. Use redundant filters for fault tolerance
- 2. Restrict communication
- 3. Eliminate nondeterminism



The Privacy Firewall

- (h+1)²-filter grid tolerates h
 Byzantine failures
- A filter only communicates with filters immediately above or below
- Each filter checks both reply and request certificates
- 🔊 Safe
- h+1 rows \rightarrow one is correct

🛛 Live

- h+1 columns \rightarrow one is correct
- Restricts nondeterminism
- threshold cryptography for replies
- cluster A locks rsn
- controlled message retransmission

Inside the PF



- (h+1)²-filter grid tolerates h
 Byzantine failures
- A filter only communicates with filters immediately above or below
- Each filter checks both reply and request certificates
- 👩 Safe
- h+1 rows \rightarrow one is correct
- 🛛 Live
 - h+1 columns → one is correct
- Restricts nondeterminism
- threshold cryptography for replies
- cluster A locks rsn
- controlled message retransmission

Inside the PF



- (h+1)²-filter grid tolerates h
 Byzantine failures
- A filter only communicates with filters immediately above or below
- Seach filter checks both reply and request certificates
- 💿 Safe
 - h+1 rows \rightarrow one is correct
- Live
- h+1 columns \rightarrow one is correct
- Ø Restricts nondeterminism
- threshold cryptography for replies
- cluster A locks rsn
- controlled message retransmission

Inside the PF



- (h+1)²-filter grid tolerates h
 Byzantine failures
- A filter only communicates with filters immediately above or below
- Seach filter checks both reply and request certificates
- 👩 Safe
- h+1 rows \rightarrow one is correct

🕤 Live

- h+1 columns → one is correct
- Restricts nondeterminism
- threshold cryptography for replies
- cluster A locks rsn
- controlled message retransmission



Inside the PF

- (h+1)²-filter grid tolerates h
 Byzantine failures
- A filter only communicates with filters immediately above or below
- Each filter checks both reply and request certificates
- 👩 Safe
- h+1 rows \rightarrow one is correct
- 🛛 Live
- h+1 columns \rightarrow one is correct
- Restricts nondeterminism
- threshold cryptography for replies
- cluster A locks rsn
- controlled message retransmission

Privacy Firewall guarantees



Output-set confidentiality

Output sequence through correct cut is a legal sequence of outputs produced by a correct node accessed trough an asynchronous, unreliable link

An exciting decade

State machine replication

- Practical Byzantine Fault Tolerance
- □ Reuse of existing (non-deterministic) implementations [SOSP 01]
- □ Reduced replication cost [SOSP 03]
- □ Low-overhead confidentiality [SOSP 03]
- □ High throughput [DSN 04]
- □ Applications: Farsite[OSDI 02], Oceanstore [FAST 03]

Quorums

- □ Fault Scalability (Q/U) [SOSP 05]
- □ Improved performance under contention (HQ) [OSDI 06]