# The Atomic Distributed Object Model for Distributed System Verification

## PhD Dissertation Defense

Wolf Honoré
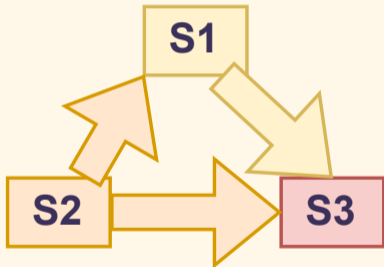
Yale University

August 19, 2022
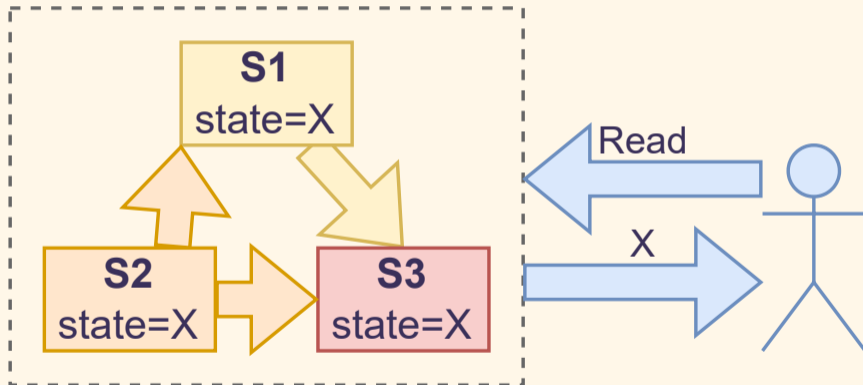
## Roadmap

- **Motivation**
  - What is a distributed system?
  - What is formal verification?
  - Why are they important?
- ADO Overview
- Case Study: Advert
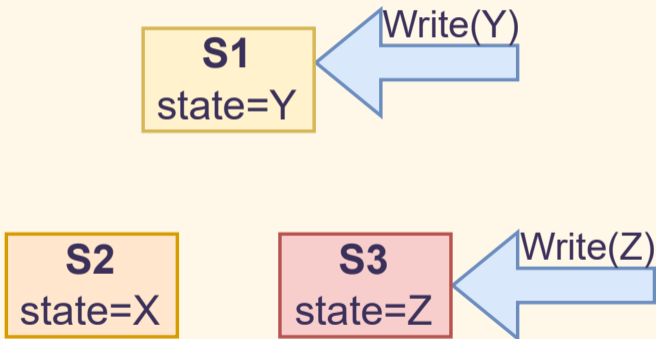- Case Study: Adore
- Case Study: AdoB
- Conclusions

# What is a Distributed System?

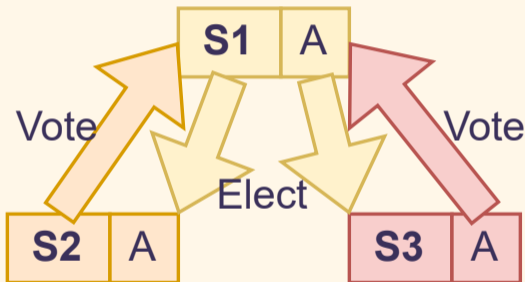# What is a Distributed System?

Motivation
oo●oooooooooooo

ADO Model
ooooooooo

Advert
oooooo

Adore
oooooooooo

AdoB
oooooooooo

Conclusions
ooo

## Replication: Challenges

**S1**
state=Y

Write(Y)

**S2**
state=X

**S3**
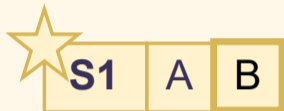state=Z

Write(Z)

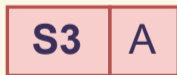## Consensus: Reaching Agreement



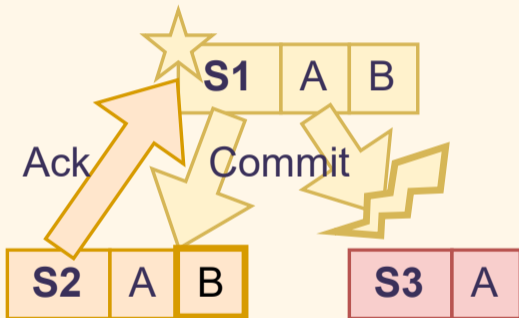election:
**S1** collects
votes

## Consensus: Reaching Agreement



local update:
**S1** applies B

# Consensus: Reaching Agreement



commit:
**S1** replicates B

2 out of 3 is
sufficient

# What Can Go Wrong?



COMPANY ANNOUNCEMENTS

## Today's outage for several Google services

Jan 24, 2014 · 1 min read

Share

**B** Ben Treynor
VP, Engineering

CRYPTO WORLD

## Hackers have stolen $1.4 billion this year using crypto bridges. Here's why it's happening

PUBLISHED WED, AUG 10 2022·4:36 PM EDT | UPDATED WED, AUG 10 2022·5:25 PM EDT

POSTED ON OCTOBER 4, 2021 TO NETWORKING & TRAFFIC

### Update about the October 4th outage

Dropbox.Tech    Topics ⌄    Developers    Jobs ↗

## Outage post-mortem

// By Akhil Gupta • Jan 12, 2014
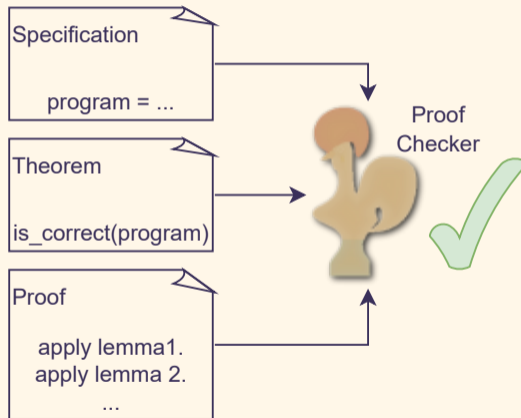
Summary of the Amazon EC2 and Amazon RDS Service
Disruption in the US East Region

April 29, 2011

FACEBOOK

# Formal Verification: Proving Correctness



Specification

program = ...

Theorem

is_correct(program)

Proof

apply lemma1.
apply lemma 2.
...

Proof Checker

## Abstraction Layers

## Abstraction Layers

**Application**

**Theorem**
write(x) followed by
read() returns x

Motivation
○○○○○○○●○○○○○○

ADO Model
○○○○○○○○

Advert
○○○○○○

Adore
○○○○○○○○○○

AdoB
○○○○○○○○○○

Conclusions
○○○

# Abstraction Layers

**Protocol**　　　　　**Application**

←──────────────────────────→

**S1**　commit
ldr=true

**S2**　　**S3**

**Theorem**
Only leaders can commit

# Abstraction Layers



**Network**   **Protocol**   **Application**

**S1**   **S1** commit
         ldr=true

**S2** ← **S3**   **S2**   **S3**

**Theorem**
Messages are
delivered in order

# Network-Based Models



Individual servers with local logs

| S1 | A | B | C |

| S2 | A | B |

| S3 | A | B |

C uncommitted

S2 adds D

| S1 | A | B | C |

| S2 | A | B | D |

| S3 | A | B |

S1 and S2 disagree

S2 replicates its log

One packet dropped

| S1 | A | B | C |

| S2 | A | B | D |

| S3 | A | B | D |

D committed

# State Machine Replication (SMR)

Single log

A | B  →  Log grows atomically  →  A | B | D

Uncommitted
commands hidden

## Abstraction Spectrum

## Prior Consensus Verification Work

| | |
|---|---|
| IronFleet (SOSP '15) | Semi-automates refining network-level specifications with SMT. |
| Verdi (PLDI '15) | Transforms simplified network specifications into more fault-tolerant equivalents. |
| Paxos Made EPR (OOPSLA '17) | Reduces the safety of Paxos to a decidable first-order logic. |
| Velisarios (ESOP '18) | Proves PBFT's safety using happens-before relations on network events. |
| Aneris (ESOP '20) | Supports modular network-based specifications with thread-level concurrency. |

## Contributions

► ADO Model: A novel, protocol-level model for consensus.

## Contributions

- ► ADO Model: A novel, protocol-level model for consensus.
- ► Compositional distributed application reasoning.

## Contributions

- ▶ ADO Model: A novel, protocol-level model for consensus.
- ▶ Compositional distributed application reasoning.
- ▶ Safety and liveness proofs.
  - ▶ First to support hot reconfiguration.
  - ▶ First to generically support benign and byzantine failures.

## Contributions

- ▶ ADO Model: A novel, protocol-level model for consensus.
- ▶ Compositional distributed application reasoning.
- ▶ Safety and liveness proofs.
  - ▶ First to support hot reconfiguration.
  - ▶ First to generically support benign and byzantine failures.
- ▶ Refinement with multiple protocols.
  - ▶ Paxos (single, multi, vertical, CAS)
  - ▶ Chain Replication
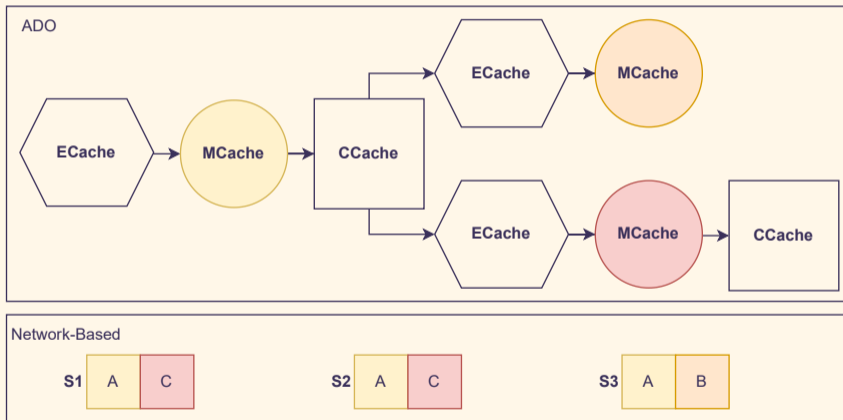  - ▶ Raft
  - ▶ Jolteon

## Acknowledgments

- Jieung Kim: Paxos safety and refinement.
- Ji-Yong Shin: Paxos refinement, OCaml extraction, performance experiments.
- Longfei Qiu: Jolteon refinement.
- Yoonseung Kim: Jolteon refinement.

# Roadmap

- Motivation
- **ADO Overview**
    - Atomic Distributed Objects
    - Global state representation (*cache tree*).
    - Atomic interface (*pull*, *invoke*, *push*).
- Case Study: Advert
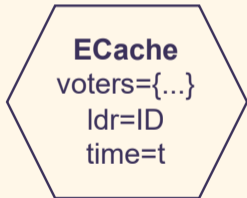- Case Study: Adore
- Case Study: AdoB
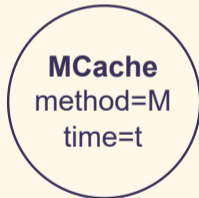- Conclusions

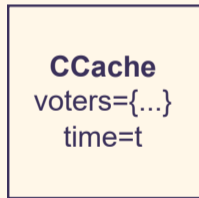# ADO State — Cache Tree

## ADO State — Cache Tree

Created by **pull**
(election)

Created by **invoke**
(local log update)

Created by **push**
(commit)

**ECache**
voters={...}
ldr=ID
time=t

**MCache**
method=M
time=t

**CCache**
voters={...}
time=t

# ADO API — Pull

ADO



Raft

3. Vote

2. Compare with local log

**S1**
time=0 | Empty

**S2**
time=0 | Empty

**S3**
time=0 | Empty

1. Ask for votes

1. Ask for votes (packet dropped)

## ADO API — Pull

# ADO API — Invoke

ADO

ECache
voters={S1,S2}
ldr=S1
time=1

Raft

1. Append to log

**S1**
time=1    A

**S2**
time=1    Empty

**S3**
time=0    Empty

Motivation
0000000000000

ADO Model
0000000

Advert
000000

Adore
0000000000

AdoB
0000000000

Conclusions
000

# ADO API — Invoke

ADO

1. Add MCache

**ECache**
voters={**S1**,**S2**}
ldr=**S1**
time=1

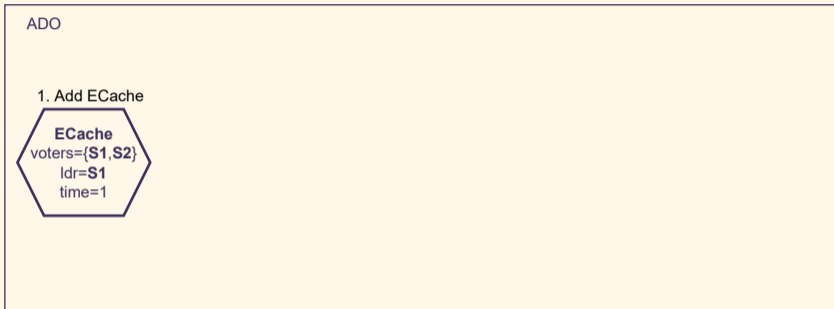**MCache**
method=A
time=1

Raft

**S1**
time=1    A

**S2**
time=1    Empty

**S3**
time=0    Empty

# ADO API — Push

ADO

# ADO API — Push

Motivation
000000000000

ADO Model
00000●00

Advert
000000

Adore
0000000000

AdoB
0000000000

Conclusions
000

# ADO API — Steady State

# ADO API — Steady State

# ADO API — Branching

# ADO API — Branching

# ADO API — Branching

# Safety

# Safety



ECache
voters={S1,S2}
ldr=S1
time=1

MCache
method=A
time=1

CCache
voters={S1,S3}
time=1

ECache
voters={S1,S3}
ldr=S3
time=2

MCache
method=B
time=2

ECache
voters={S1,S2}
ldr=S1
time=3

MCache
method=C
time=3

CCache
voters={S1,S2}
time=3

Election must choose
committed branch

ECache
voters={S1,S2}
ldr=S3
time=4

# Roadmap

- Motivation
- ADO Overview
- **Case Study: Advert**
  - <u>A</u>tomic <u>D</u>istributed Object <u>Ve</u>rification <u>T</u>oolchain
  - Expose partial failures for distributed application optimization.
  - Support ADO composition.
- Case Study: Adore
- Case Study: AdoB
- Conclusions

## Distributed Applications with Partial Failures

> *Partial failure is a central reality of distributed computing. [. . . ] Being robust in the face of partial failure requires some expression at the interface level.*
> *(Jim Waldo. A Note on Distributed Computing. 1994)*

▶ Unavoidable feature unique to distributed systems.

▶ Interact with all aspects of distributed protocols (e.g., leader election and reconfiguration).

▶ Can be used for performance optimizations.
   ▶ TAPIR (SOSP '15): Transactions with out-of-order commits.
   ▶ Speculator (SOSP '05): Speculative distributed file system.

# Distributed Applications with Partial Failures

# Distributed Applications with Partial Failures

# Distributed Applications with Partial Failures

# Distributed Applications with Partial Failures

## Distributed Applications

```
1 ADO KV {
2   shared kv : [string * int] := [];
3   method set(k, v) { this.kv[hash(k)] := (v, len(v)); }
4   method get(k) { return this.kv[hash(k)][0]; }
5   method getmeta(k) { return this.kv[hash(k)][1]; }
6 }
```

## Distributed Applications

```
1 ADO DVec[T] {
2   shared data : [T] := [];
3   method insert(idx, x) { this.data[idx] := x; }
4   method get(idx) { return this.data[idx]; }
5 }
6 ADO DLock {
7   shared owner : option N := None;
8   method tryAcquire() { ... }
9   method release() { ... }
10 }
11 DApp KVLock(lk: DLock, data: DVec[string], meta: DVec[int]) {
12   proc set(k, v) {
13     ... /* acquire, set data, set meta, release */
14   }
15   ... /* get, getmeta */
16 }
```

## Distributed Applications

```
1 DApp KVLock(lk: DLock, data: DVec[string], meta: DVec[int]) {
2   proc set(k, v) {
3     lk.pull();
4
5
6
7
8   }
9 }
```

## Distributed Applications

```
1 DApp KVLock(lk: DLock, data: DVec[string], meta: DVec[int]) {
2   proc set(k, v) {
3     while (lk.pull() == FAIL) {}
4
5
6
7
8   }
9 }
```

## Distributed Applications

```
1 DApp KVLock(lk: DLock, data: DVec[string], meta: DVec[int]) {
2   proc set(k, v) {
3     while (lk.pull() == FAIL) {}
4     ok := lk.invoke(tryAcquire());
5
6
7
8   }
9 }
```

# Distributed Applications

```
1 DApp KVLock(lk: DLock, data: DVec[string], meta: DVec[int]) {
2   proc set(k, v) {
3     while (lk.pull() == FAIL) {}
4     ok := lk.invoke(tryAcquire());
5     while (lk.push() == FAIL) {}
6     if (!ok) { return; }
7     /* ... */
8   }
9 }
```

## Handling Failures

```
1 DApp KVLockAbort(lk: DLock, data: DVec[string], meta: DVec[int]) {
2   proc set(k, v) {
3     if (lk.pull() == FAIL) { return; }
4     ok := lk.invoke(tryAcquire());
5     if (lk.push() == FAIL) { return; }
6     if (!ok) { return; }
7     /* ... */
8   }
9 }
```

## Handling Failures

```
1 DApp KVLockRetry(lk: DLock, data: DVec[string], meta: DVec[int]) {
2   proc set(k, v) {
3     for retry in 0..N {
4       if (lk.pull() == FAIL) { continue; }
5       ok := lk.invoke(tryAcquire());
6       if (lk.push() == FAIL) { continue; }
7       if (!ok) { continue; }
8     }
9     if (retry == N) { return; }
10    /* ... */
11  }
12 }
```

## Handling Failures

```
1 obj.m()! :=
2   while (obj.pull() == FAIL) {}
3   obj.invoke(m());
4   while (obj.push() == FAIL) {}
5
6 DApp KVLock(lk: DLock, data: DVec[string], meta: DVec[int]) {
7   proc set(k, v) {
8     ok := lk.tryAcquire()!;
9     if (!ok) { return; }
10    data.insert(hash(k), v)!;
11    meta.insert(hash(k), len(v))!;
12    lk.release()!;
13  }
14 }
```

Motivation
○○○○○○○○○○○○○○○○

ADO Model
○○○○○○○

Advert
○○○○●○○

Adore
○○○○○○○○○○

AdoB
○○○○○○○○○○

Conclusions
○○○

# End-to-End Verification

# End-to-End Verification

# End-to-End Verification

Motivation
0000000000000

ADO Model
00000000

**Advert**
0000●0

Adore
0000000000

AdoB
0000000000

Conclusions
000

# End-to-End Verification



Request Layer

```
Definition elect (time: nat) : bool :=
  let ok := broadcast(Elect) in ...
```
Coq

```
int elect(time: int) {
  int ok = broadcast(ELECT); ... }
```
C

Broadcast Layer

```
Definition broadcast (msg: Msg) : bool :=
  fold (send msg) replicas
```
Coq

```
int broadcast(msg: Msg) {
  for (int i = 0; i < NUM_REPLICAS; i++) {
    send(msg, replicas[i]); } }
```
C

Send Layer

...

## Proof Effort

|  | **Proof LOC (Coq)** |
|---|---|
| KVLock DApp | ~600 |
| KVLockFree DApp | ~300 |
| 2PC DApp | ~600 |
| Generic Paxos Refinement | ~5k |
| Chain Replication Refinement | ~2k |
| Shared Libraries | ~11k |
| Multi Paxos C Refinement | ~44k |
| Single Paxos | ~80 |
| Multi Paxos | ~90 |
| Vertical Paxos | ~100 |
| CASPaxos | ~80 |

## Roadmap

- Motivation
- ADO Overview
- Case Study: Advert
- **Case Study: Adore**
  - Atomic Distributed Objects with Certified Reconfiguration
  - Prove safety at the ADO level.
  - Support hot reconfiguration.
- Case Study: AdoB
- Conclusions

# Reconfiguration

Hot Reconfiguration

## Reconfiguration



Committed log entries decided by current configuration

Consensus

Current configuration stored in log

Reconfiguration

## Reconfiguration

bug in single-server membership changes    Subscribe ☐    ↕

4784 views

onga...@gmail.com    Jul 10, 2015, 12:58:53 AM    ☆    ↩    ⋮
to raft...@googlegroups.com

Hi raft-dev,

Unfortunately, I need to announce a bug in the dissertation version of membership changes (the single-server changes, not joint consensus). The bug is potentially severe, but the fix I'm proposing is easy to implement.

Diego Ongaro
(Raft designer)

| Raft hot reconfiguration published | Critical bug discovered | Fix proposed | No complete safety proofs | Safety proved in Adore |
|---|---|---|---|---|
| August 2014 | May 2015 | July 2015 | | June 2021 |

# Safety in Adore

# Reconfiguration in Adore

Created by **pull**
(election)

Created by **push**
(commit)

**ECache**
voters={...}
ldr=ID
time=t

**CCache**
voters={...}
time=t

Created by **invoke**
(local log update)

Created by **reconfig**
(local log update)

**MCache**
method=M
time=t

**RCache**
config={...}
time=t

# Reconfiguration in Adore



ADO

Config is {**S1**,**S2**,**S3**}

**S3**'s config is now
{**S1**,**S2**,**S3**,**S4**}

**ECache**
voters={**S1**,**S2**}
ldr=**S1**

**MCache**
method=A

**CCache**
voters={**S1**,**S3**}

**ECache**
voters={**S1**,**S3**}
ldr=**S3**

**RCache**
config=
{**S1**,**S2**,**S3**,**S4**}

Raft

**S1**
config={**S1**,**S2**,**S3**}  | A

**S2**
config={**S1**,**S2**,**S3**}  | Empty

**S3**
config={**S1**,**S2**,**S3**,**S4**}  | A | **+S4**

Motivation
0000000000000

ADO Model
00000000

Advert
000000

**Adore**
0000●00000

AdoB
0000000000

Conclusions
000

# Reconfiguration Rules

# Reconfiguration Rules

# Reconfiguration Rules

Original config = C     New config = C'

Rule 1

... → **RCache** config=C'

$$\forall Q \subseteq C. \; \forall Q' \subseteq C'. \\ is\_quorum(Q) \land \\ is\_quorum(Q') \Rightarrow \\ Q \cap Q' \neq \emptyset$$

Must be a CCache between

Rule 2

**RCache** config=C → ... → **RCache** config=C'

Must be a CCache between

Rule 3

**ECache** → ... → **RCache** config=C

Motivation
000000000000
ADO Model
00000000
Advert
000000
**Adore**
0000●00000
AdoB
0000000000
Conclusions
000

## Reconfiguration Rules

Original config =
{**S1**,**S2**,**S3**,**S4**}

No Rule 3 leads to a
safety bug

**ECache**
voters=
{**S1**,**S2**,**S3**}
ldr=**S1**
time=1

**RCache**
config=
{**S1**,**S2**,**S3**}
time=1

# Reconfiguration Rules

Original config =
{**S1**,**S2**,**S3**,**S4**}

No Rule 3 leads to a
safety bug

**ECache**
voters=
{**S1**,**S2**,**S3**}
ldr=**S1**
time=1

**RCache**
config=
{**S1**,**S2**,**S3**}
time=1

**ECache**
voters=
{**S2**,**S3**,**S4**}
ldr=**S2**
time=2

# Reconfiguration Rules

Original config =
{**S1**,**S2**,**S3**,**S4**}

No Rule 3 leads to a
safety bug

**ECache**
voters=
{**S1**,**S2**,**S3**}
ldr=**S1**
time=1

**RCache**
config=
{**S1**,**S2**,**S3**}
time=1

**ECache**
voters=
{**S2**,**S3**,**S4**}
ldr=**S2**
time=2

**RCache**
config=
{**S1**,**S2**,**S4**}
time=2

# Reconfiguration Rules

# Reconfiguration Rules



Original config = {**S1**,**S2**,**S3**,**S4**}

No Rule 3 leads to a safety bug

**ECache**
voters=
{**S1**,**S2**,**S3**}
ldr=**S1**
time=1

**RCache**
config=
{**S1**,**S2**,**S3**}
time=1

**ECache**
voters={**S1**,**S3**}
ldr=**S1**
time=3

{**S1**,**S3**} is a quorum of {**S1**,**S2**,**S3**}

**ECache**
voters=
{**S2**,**S3**,**S4**}
ldr=**S2**
time=2

**RCache**
config=
{**S1**,**S2**,**S4**}
time=2

**CCache**
voters={**S2**,**S4**}
time=2

# Reconfiguration Rules



Original config =
{**S1,S2,S3,S4**}

No Rule 3 leads to a
safety bug

**ECache**
voters=
{**S1,S2,S3**}
ldr=**S1**
time=1

**RCache**
config=
{**S1,S2,S3**}
time=1

**ECache**
voters={**S1,S3**}
ldr=**S1**
time=3

**MCache**
method=A
time=3

**CCache**
voters={**S1,S3**}
time=3

**ECache**
voters=
{**S2,S3,S4**}
ldr=**S2**
time=2

**RCache**
config=
{**S1,S2,S4**}
time=2

**CCache**
voters={**S2,S4**}
time=2

CCaches on different branches
= consistency is broken

# Reconfiguration Rules

Original config = {**S1**,**S2**,**S3**,**S4**}

Rule 3 prevents the safety bug

**ECache** voters= {**S1**,**S2**,**S3**} ldr=**S1** time=1

**RCache** config= {**S1**,**S2**,**S3**} time=1

Required by Rule 3

**ECache** voters= {**S2**,**S3**,**S4**} ldr=**S2** time=2

**MCache** method=A time=2

**CCache** voters= {**S2**,**S3**,**S4**} time=2

**RCache** config= {**S1**,**S2**,**S4**} time=2

**CCache** voters={**S2**,**S4**} time=2

# Reconfiguration Rules



Original config =
{**S1,S2,S3,S4**}

Rule 3 prevents the
safety bug

**ECache**
voters=
{**S1,S2,S3**}
ldr=**S1**
time=1

**RCache**
config=
{**S1,S2,S3**}
time=1

**ECache**
voters={**S1,S3**}
ldr=**S1**
time=3

Now impossible

**S3** voted here already

**ECache**
voters=
{**S2,S3,S4**}
ldr=**S2**
time=2

**MCache**
method=A
time=2

**CCache**
voters=
{**S2,S3,S4**}
time=2

**RCache**
config=
{**S1,S2,S4**}
time=2

**CCache**
voters={**S2,S4**}
time=2

## Proving Safety



Nearest common
CCache ancestor

**CCache**

...

**ECache**
voters=Q
time=t

t ≠ t' because ?

...

**ECache**
voters=Q'
time=t'

# Proving Safety



Diagram contents:

- No RCaches
- Nearest common CCache ancestor
- CCache
- ... → **ECache** voters=Q time=t
- t ≠ t' because configs are equal so Q ∩ Q' ≠ ∅
- ... → **ECache** voters=Q' time=t'
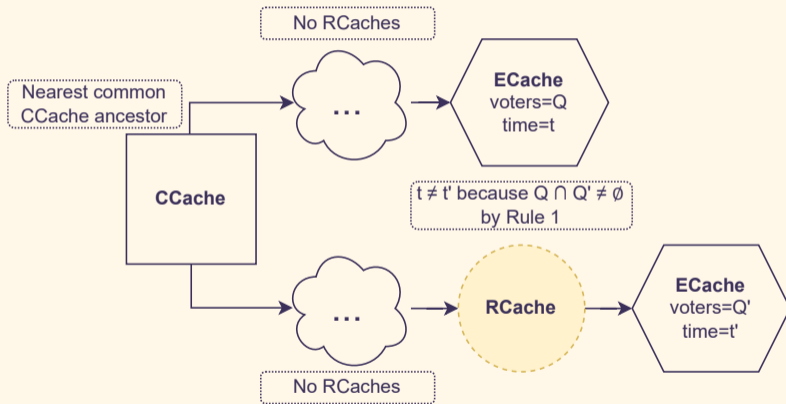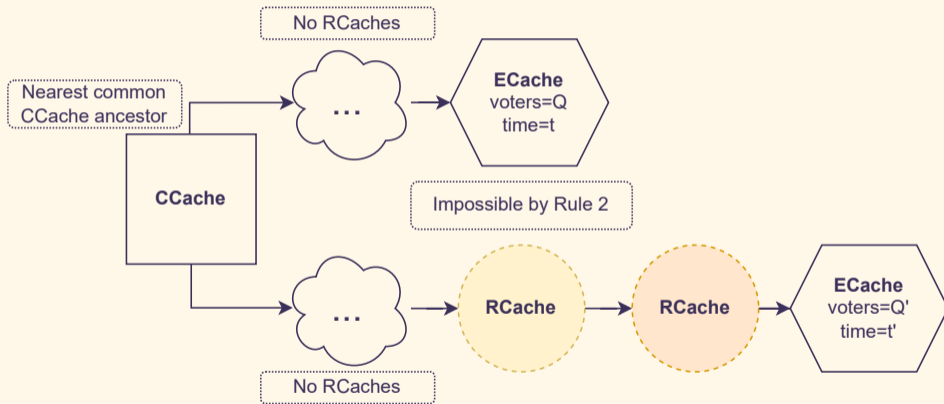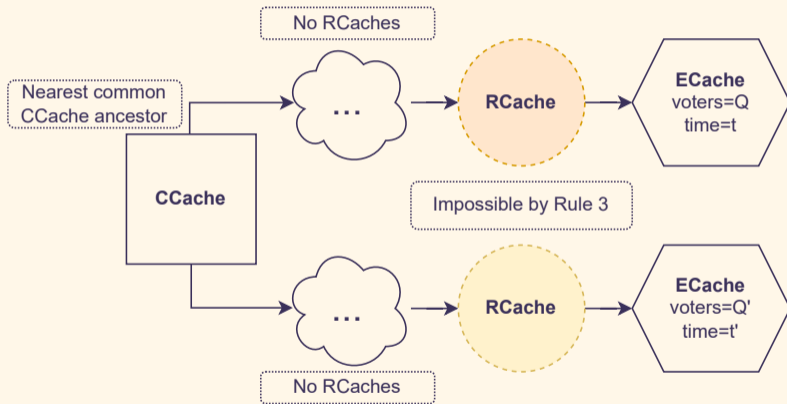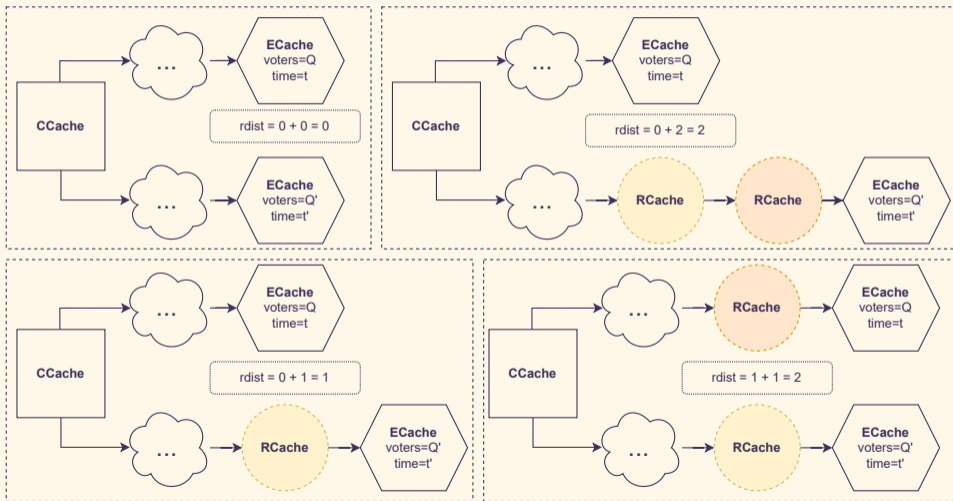- No RCaches

## Proving Safety

# Proving Safety

# Proving Safety

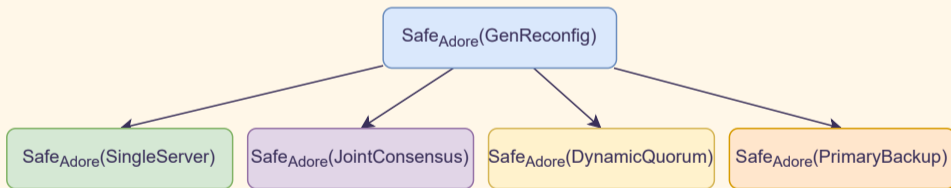## Proving Safety

## Generalized Quorums

- Safety proved once for generic reconfiguration scheme.
- A quorum is any set that guarantees overlap.
- Can be instantiated many times with minimal proof effort.

```
                    Safe_Adore(GenReconfig)

Safe_Adore(SingleServer)  Safe_Adore(JointConsensus)  Safe_Adore(DynamicQuorum)  Safe_Adore(PrimaryBackup)
```

## Generalized Quorums

### Single-Server

$$Config \triangleq Set(\mathbb{N}_{nid})$$

$$canReconfig(C, C') \triangleq C = C' \lor$$
$$\exists s.\, C = C' \cup \{s\} \lor C' = C \cup \{s\}$$

$$isQuorum(S, C) \triangleq |C| < 2 * |S \cap C|$$

Motivation
000000000000

ADO Model
00000000

Advert
000000

Adore
0000000●000

AdoB
0000000000

Conclusions
000

## Generalized Quorums

### Joint Consensus

$$Config \triangleq Set(\mathbb{N}_{nid}) * Option(Set(\mathbb{N}_{nid}))$$
$$canReconfig(C, C') \triangleq \exists\, old.\, (C = (old, \bot) \wedge C' = (old, \_)) \vee$$
$$\exists\, new.\, (C = (\_, new) \wedge C' = (new, \bot))$$
$$isQuorum(S, (old, new)) \triangleq |old| < 2 * |S \cap old| \wedge$$
$$(new = \bot \vee |new| < 2 * |S \cap new|)$$

## Generalized Quorums

### Dynamic Quorum Size

$$Config \triangleq \mathbb{N} * Set(\mathbb{N}_{nid})$$

$$canReconfig((q, C), (q', C')) \triangleq (C \subseteq C' \wedge |C'| < q + q') \vee$$
$$(C' \subseteq C \wedge |C| < q + q')$$

$$isQuorum(S, (q, C)) \triangleq q \leq |S \cap C|$$

## Generalized Quorums

### Primary Backup

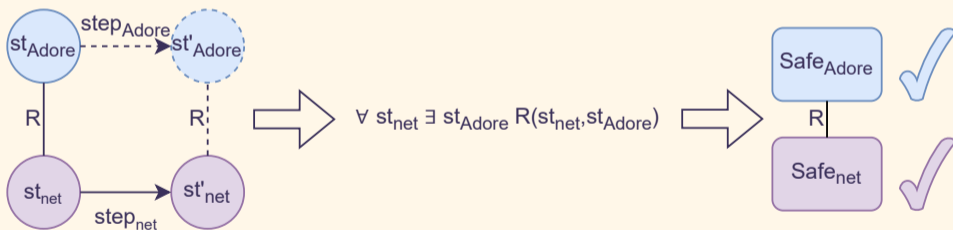$$Config \triangleq \mathbb{N}_{nid} * Set(\mathbb{N}_{nid})$$

$$canReconfig((P, \_), (P', \_)) \triangleq P = P'$$
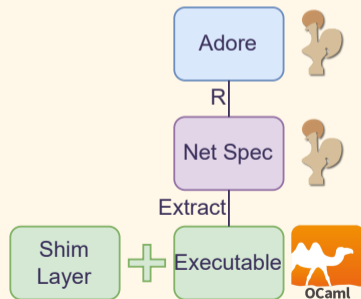
$$isQuorum(S, (P, \_)) \triangleq P \in S$$

## Refinement

- ▶ Refinement between Raft network-based specification and Adore.
- ▶ Also generic with respect to reconfiguration scheme.

# Extraction

- ▶ Automated extraction from Coq specification to executable OCaml.
- ▶ Extracted code contains core logic, unverified shim layer handles network communication.
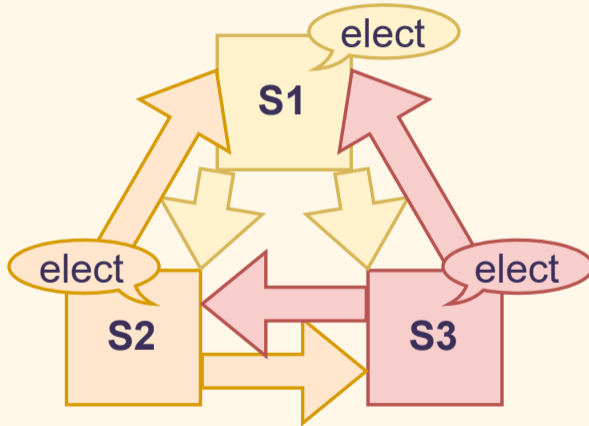- ▶ Safety guaranteed through Adore and refinement.

Motivation
000000000000
ADO Model
00000000
Advert
000000
**Adore**
000000000●
AdoB
0000000000
Conclusions
000

## Proof Effort

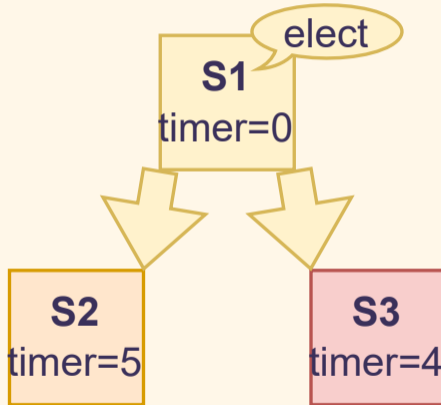|                               | **Proof LOC (Coq)** | **Proof Time**   |
| ----------------------------- | ------------------- | ---------------- |
| Cache Tree Library/Properties | ∼6k                 | 2 person-weeks   |
| Safety Proof                  | ∼4k                 | 3 person-weeks   |
| Refinement Proof              | ∼13k                | 9 person-weeks   |
| Reconfiguration Schemes (6)   | ∼300                | <1 person-week   |

# Roadmap

- ▶ Motivation
- ▶ ADO Overview
- ▶ Case Study: Advert
- ▶ Case Study: Adore
- ▶ **Case Study: AdoB**
  - ▶ <u>A</u>tomic <u>D</u>istributed <u>O</u>bjects for <u>B</u>enign/<u>B</u>yzantine Consensus
  - ▶ Prove liveness at the ADO level.
  - ▶ Support benign and byzantine failures in a generic abstraction.
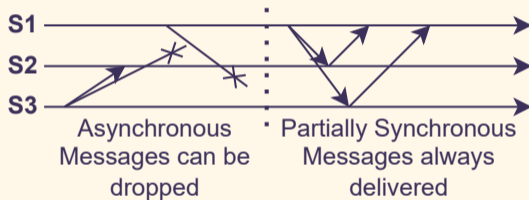- ▶ Conclusions

# Liveness

# Liveness

# Liveness Assumptions

▶ Partial synchrony

S1 ——————————————————→

S2 ——————————————————→

S3 ——————————————————→

Asynchronous
Messages can be
dropped

Partially Synchronous
Messages always
delivered

Liveness Assumptions
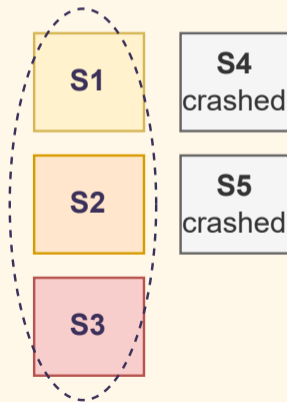
- Partial synchrony
- Productive strategy

```
if not isLeader() and timer() == 0:
  startElection()
else if isLeader() and hasUncommitted():
  startCommit()
else if timer() == 0:
  sendTimeout()
```
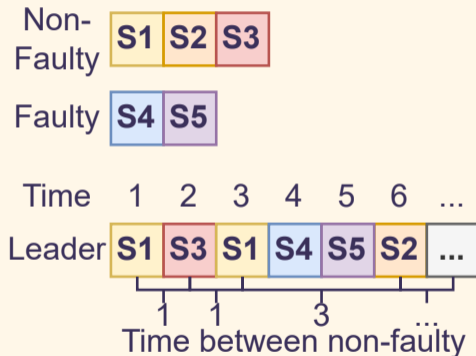
## Liveness Assumptions

Quorum = majority = 3/5

- ▶ Partial synchrony
- ▶ Productive strategy
- ▶ Non-faulty quorum
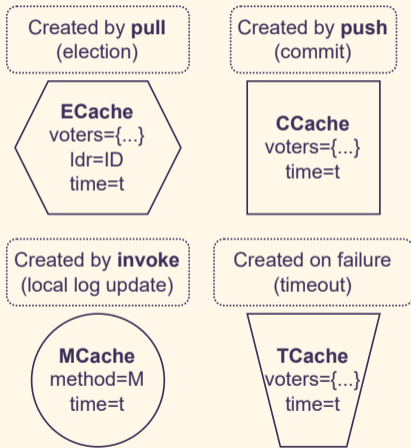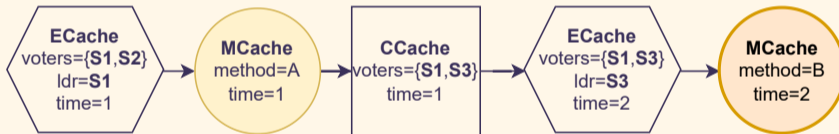
**S1**

**S4**
crashed

**S2**

**S5**
crashed

**S3**

## Liveness Assumptions

- ▶ Partial synchrony
- ▶ Productive strategy
- ▶ Non-faulty quorum
- ▶ Fair election rotation

Non-Faulty   **S1 S2 S3**

Faulty   **S4 S5**

Time   1   2   3   4   5   6   ...

Leader   **S1 S3 S1 S4 S5 S2** ...

1   1    3    ...

Time between non-faulty

# Time in AdoB



Created by **pull**
(election)

**ECache**
voters={...}
ldr=ID
time=t

Created by **push**
(commit)

**CCache**
voters={...}
time=t

Created by **invoke**
(local log update)

**MCache**
method=M
time=t

Created on failure
(timeout)

**TCache**
voters={...}
time=t

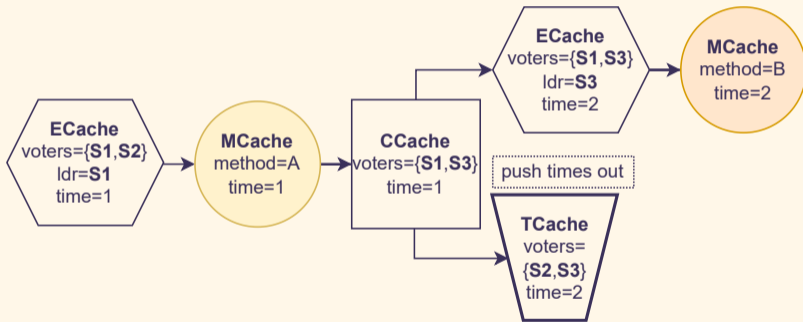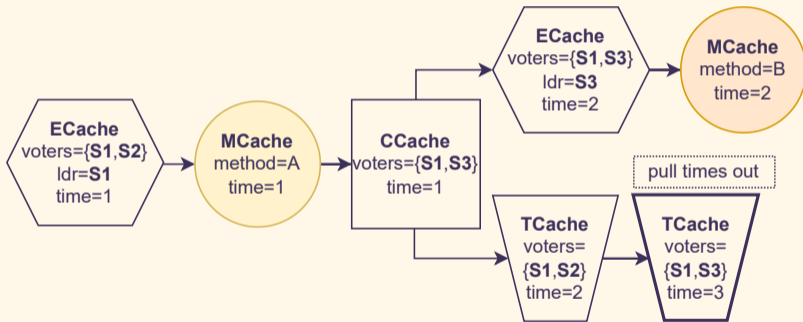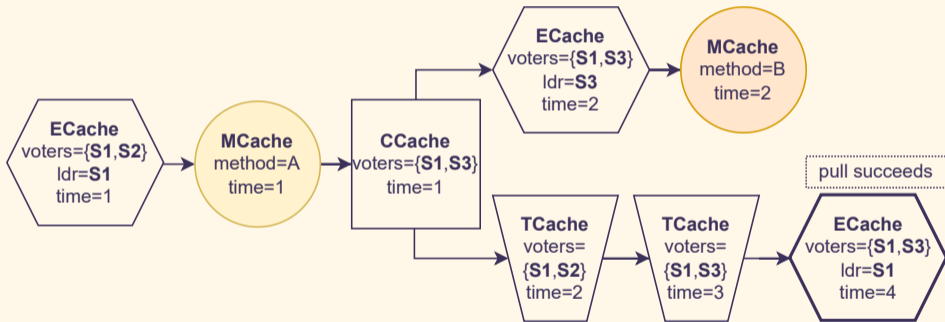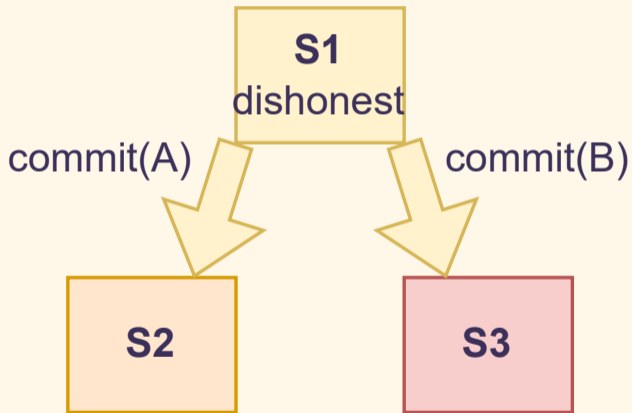# Liveness in AdoB

# Liveness in AdoB

# Liveness in AdoB

# Liveness in AdoB

# Byzantine Failures

# Byzantine Failures

Honest Dishonest

Configuration   S1 S2 S3 S4 S5

Quorums
(3/5)

S1 S3 S5

S2 S4 S5

Always overlap, but
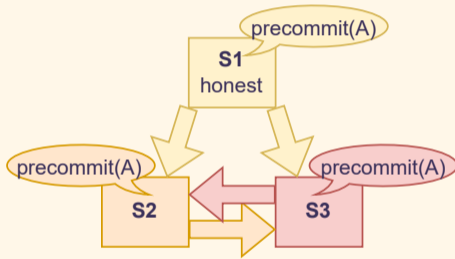may not be honest

Super Quorums
(4/5)

S1 S2 S3 S5

S4 S2 S3 S5

Overlap is always honest

# Byzantine Failures

# Byzantine Failures



precommit(A)  HotStuff/Jolteon

**S1** honest

**S2**    **S3**

commit(A)

**S1** honest

precommit(A) signed by **S2**, **S3**

**S2**    **S3**

# Byzantine Failures in AdoB



Created by **pull**
(election)

**ECache**
voters={...}
ldr=ID
time=t

Created by **push**
(commit)

**CCache**
voters={...}
time=t

Created by **invoke**
(pre-commit)

**MCache**
voters={...}
method=M
time=t

Created on failure
(timeout)

**TCache**
voters={...}
time=t

## Generalizing Benign and Byzantine Failures



honest = {**S1**,**S2**,**S3**}
dishonest = {**S4**}

common honest = {**S1**}

common honest = {**S1**,**S3**}

**ECache**
voters=
{**S1**,**S2**,**S3**}
ldr=**S1**

**MCache**
voters=
{**S1**,**S3**,**S4**}
method=A

**CCache**
voters=
{**S1**,**S2**,**S4**}

## Generalizing Benign and Byzantine Failures

honest = {**S1**,**S2**,**S3**,**S4**}
dishonest = {}

common honest = {**S1**}

common honest = {**S1**}

**ECache**
voters=
{**S1**,**S2**,**S3**}
ldr=**S1**

→

**MCache**
voters={**S1**}
method=A

→

**CCache**
voters=
{**S1**,**S2**,**S4**}

# Generalizing Benign and Byzantine Failures

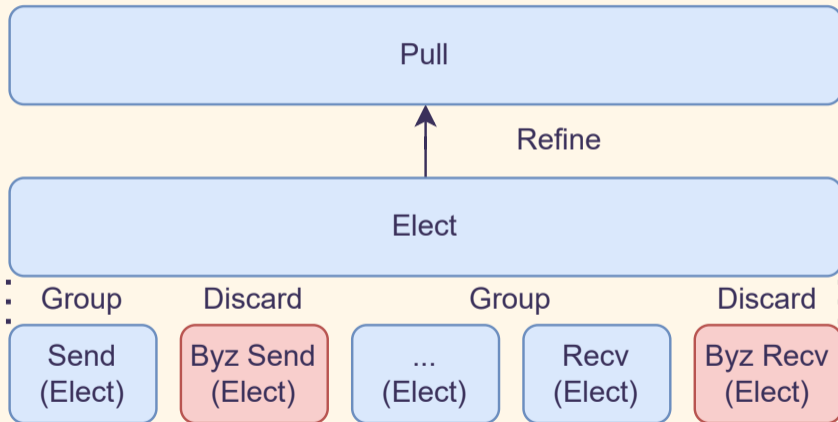| Failure Model | Required Number of Votes | | |
|---------------|---------|---------|---------|
| | pull | invoke | push |
| Benign | Quorum | Only leader | Quorum |
| Byzantine | Super Quorum | Super Quorum | Super Quorum |
| Generalized | Super Quorum | MQuorum | Super Quorum |

### Definition

Two quorums have a common voter (e.g., $> 1/2$ of configuration).
Super quorums have a common honest voter (e.g., $> 2/3$ of configuration).
An MQuorum and super quorum with the same leader have a common honest voter.

## Refinement

## Proof Effort

|  | **Proof LOC (Coq)** | **Proof Time** |
| --- | --- | --- |
| Safety Proof | $\sim$3k | 2 person-weeks |
| Liveness Proof | $\sim$3k | 2 person-weeks |
| Refinement Proof | $\sim$4k | 6 person-weeks |

# Roadmap

- ▶ Motivation
- ▶ ADO Overview
- ▶ Case Study: Advert
- ▶ Case Study: Adore
- ▶ Case Study: AdoB
- ▶ **Conclusions**
  - ▶ Summary of results.
  - ▶ Future work.

## Summary

It facilitates formal verification by hiding network-level details behind a global tree-based state representation and atomic interface.

- ▶ ADO model: novel protocol-level abstraction for consensus.
- ▶ Atomic tree-based representation of replicated state.
- ▶ Exposes partial failures to distributed applications (Advert).
- ▶ Enables safety and liveness reasoning (Adore, AdoB).
- ▶ Correctly models a wide range of consensus protocols both benign (Advert, Adore) and byzantine (AdoB).
- ▶ Supports practical extensions like reconfiguration (Adore).

## Future Work

- ► Automate refinement.
  - ► Verdi verified system transformers (PLDI '15).
  - ► CSPEC (OSDI '18), pretend synchrony (POPL '19), inductive sequentialization (PLDI '20).
- ► Generate code from ADO specification.
  - ► DeepSEA (OOPSLA '19).
- ► Expand beyond consensus.
  - ► Conflict-free replicated data types.
  - ► Causal consistency.