



**NUS**  
National University  
of Singapore

# Protocol Conformance with Choreographic PlusCal

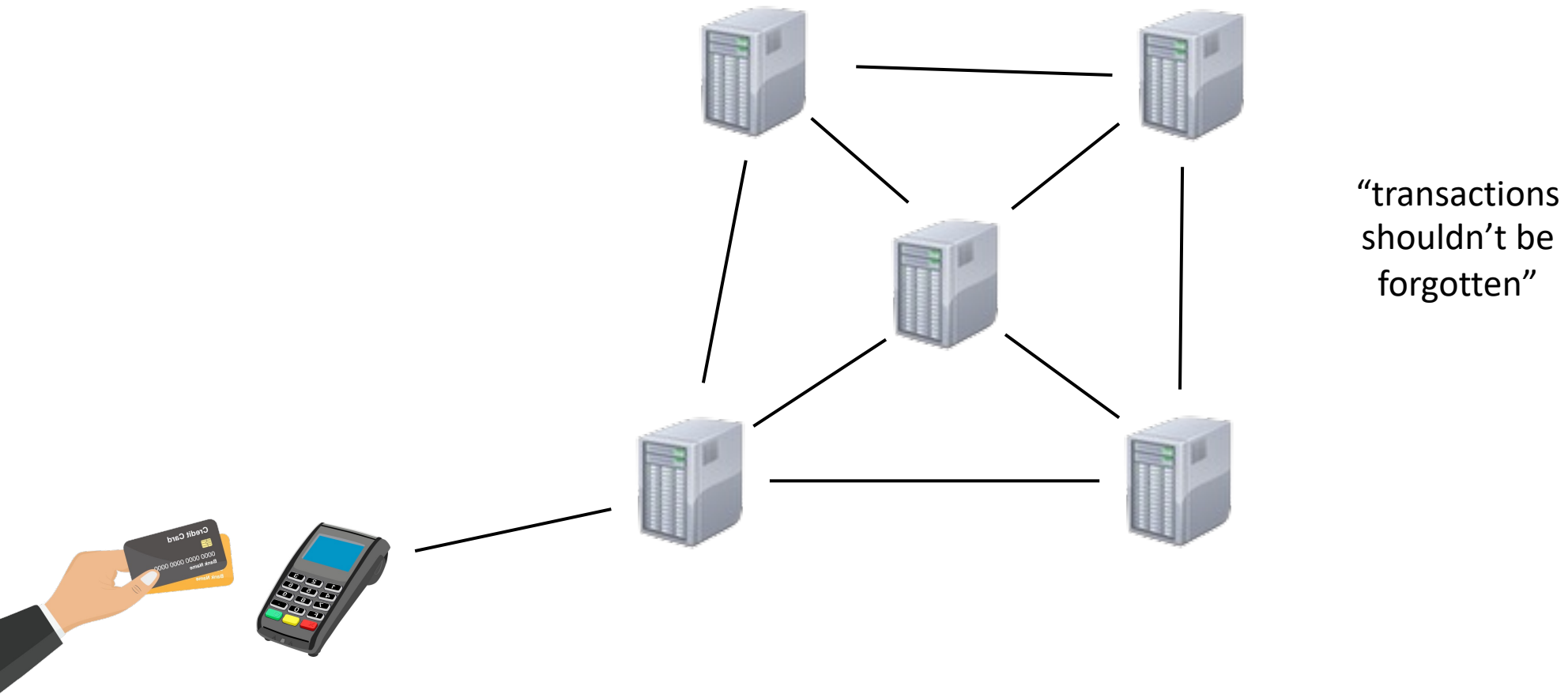
Darius Foo, Andreea Costea, and Wei-Ngan Chin

National University of Singapore

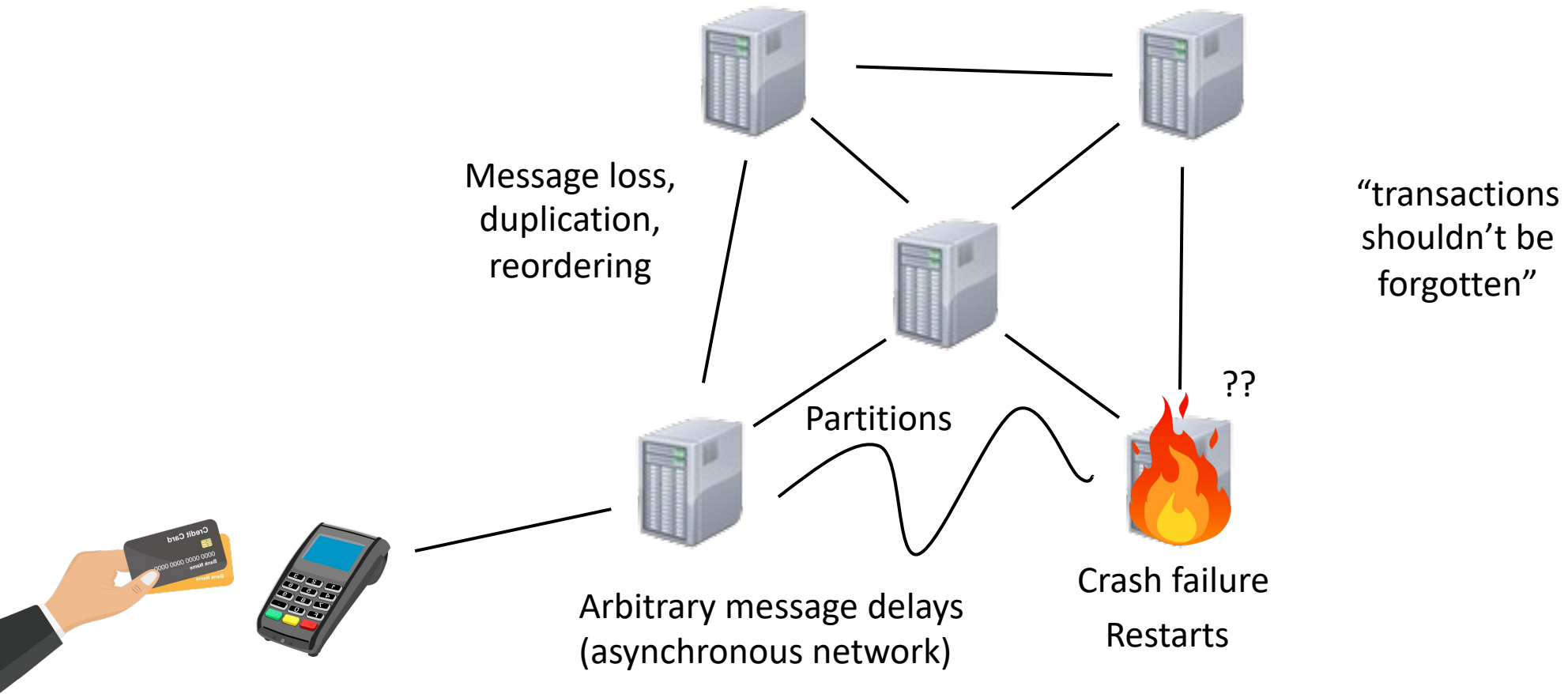
17th International Symposium on Theoretical Aspects of Software Engineering

6 July 2023

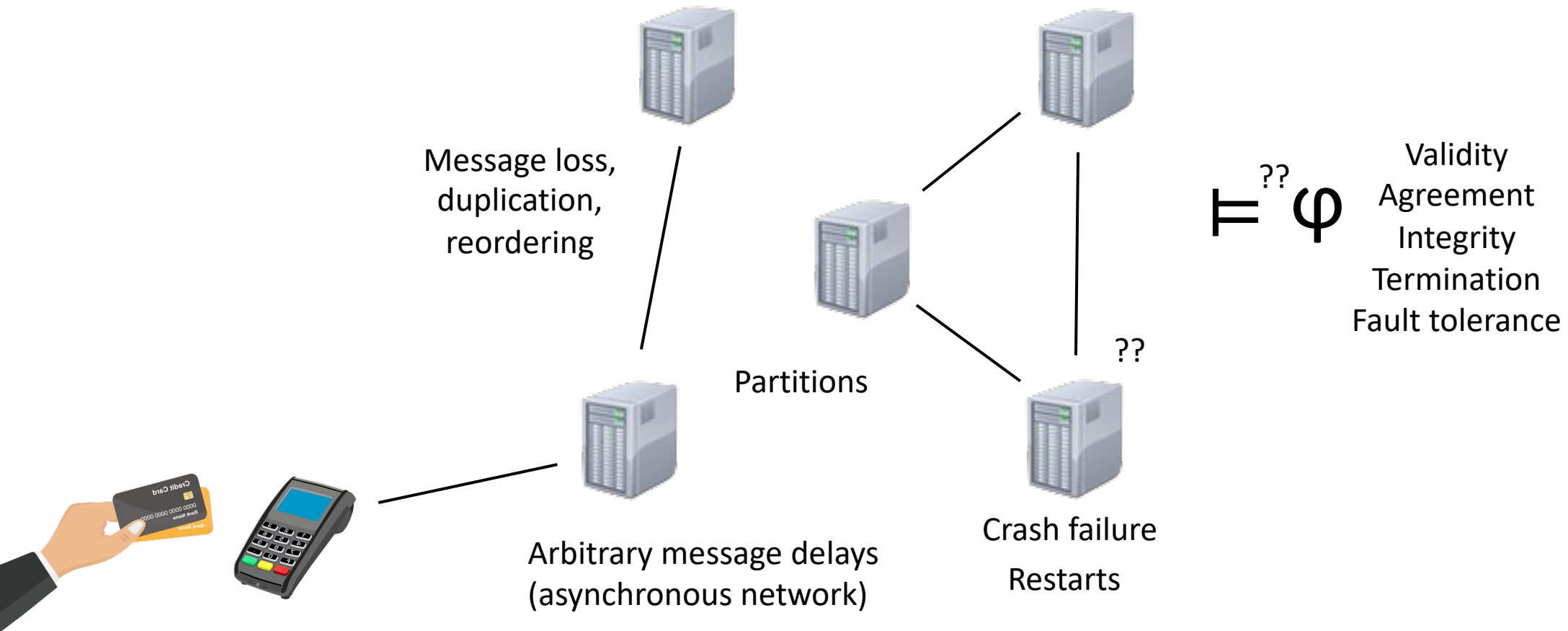
# Distributed Systems



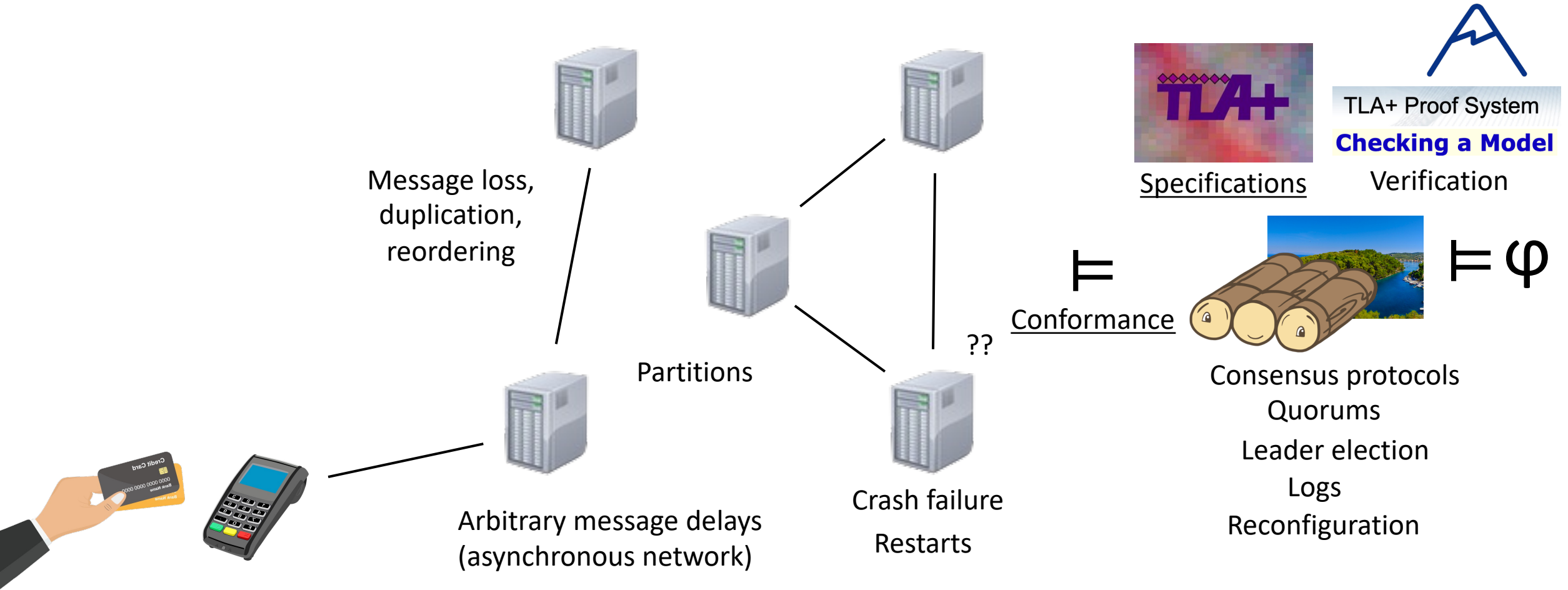
# Distributed Systems



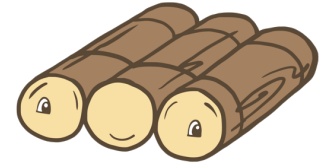
# Distributed Systems



# Distributed Systems



# How hard is it to implement a protocol correctly?

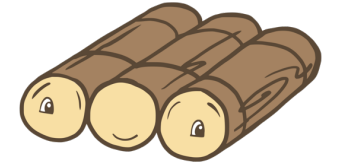


## **In Search of an Understandable Consensus Algorithm**

Diego Ongaro and John Ousterhout  
Stanford University

(2014)

# How hard is it to implement a protocol correctly?



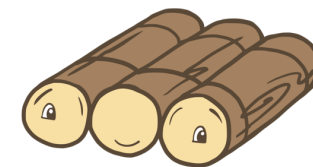
## Where can I get Raft?

There are many implementations of Raft available in various stages of development. This table lists the implementations we know about with source code available. The most popular and/or recently updated implementations are towards the top. This information will inevitably get out of date; please submit a [pull request](#) or an issue to update it.

Stars	Name	Primary Authors	Language	License	Leader Election + Log Replication?	Persistence?	Membership Changes?	Log Compaction?
13,312★	<a href="#">TiKV</a>	<a href="#">Jay</a> , <a href="#">ngaut</a> , <a href="#">siddontang</a> , <a href="#">tiancaiamao</a>	Rust	Apache-2.0	Yes	Yes	Yes	Yes
9,211★	<a href="#">nebula-graph-storage</a>	<a href="#">Sherman Ye</a> , <a href="#">Doodle Wang</a>	C++	Apache-2.0	Yes	Yes	Yes	Yes
26,166★	<a href="#">RethinkDB</a>		C++	Apache-2.0	Yes	Yes	Yes	Yes
10,501★	<a href="#">Seastar Raft</a>	<a href="#">Gleb Natapov</a> , <a href="#">Konstantin Osipov</a> , <a href="#">Pavel Solodovnikov</a> , <a href="#">Alejo Sanchez</a> , <a href="#">Kamil Braun</a> , <a href="#">Tomash Grabiec</a>	C++20	AGPL	Yes	Yes	Yes	Yes
5,439★	<a href="#">hazelcast-raft</a>	<a href="#">Mehmet Dogan</a> , <a href="#">Ensar Basri Kahveci</a>	Java	Apache-2.0	Yes	Yes	Yes	Yes
7,220★	<a href="#">hashicorp/raft</a>	<a href="#">Armon Dadgar</a>	Go	MPL-2.0	Yes	Yes	Yes	Yes
3,560★	<a href="#">braft</a>	<a href="#">Zhangyi Chen</a> , <a href="#">Yao Wang</a>	C++	Apache-2.0	Yes	Yes	Yes	Yes

... 137 more implementations

# How hard is it to implement a protocol correctly?



## Minimizing Faulty Executions of Distributed Systems (2015)

Colin Scott<sup>\*</sup>      Aurojit Panda<sup>\*</sup>      Vjekoslav Brajkovic<sup>◇</sup>      George Necula<sup>\*</sup>  
Arvind Krishnamurthy<sup>†</sup>      Scott Shenker<sup>\*◇</sup>  
<sup>\*</sup>*UC Berkeley*      <sup>◇</sup>*ICSI*      <sup>†</sup>*University of Washington*

### Abstract

When troubleshooting buggy executions of distributed systems, developers typically start by manually separating events that are responsible for triggering the bug from those that are extraneous (noise). We present DEMi, a tool for automatically performing this task. We apply DEMi to buggy executions of two prominent distributed systems, Raft and Spark, and find that it produces minimized executions that are between 1X and 4.6X the size of optimal executions.

8 bugs

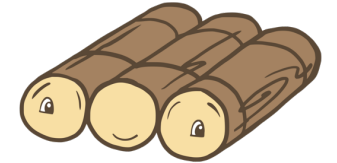
much more costly than machine time, automated minimization tools for *sequential* test cases [24, 86, 94] have already proven themselves valuable, and are routinely applied to bug reports for software projects such as Firefox [1], LLVM [7], and GCC [6].

In this paper we address the problem of automatically minimizing executions of distributed systems. We focus on executions generated by fuzz testing, but we also illustrate how one might minimize production traces.

Distributed executions have two distinguishing fea-



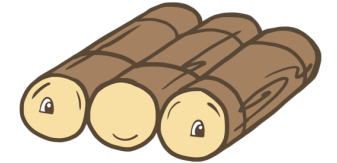
# How hard is it to implement a protocol correctly?



## Fuzz testing distributed systems with QuickCheck (2016)



# How hard is it to implement a protocol correctly?



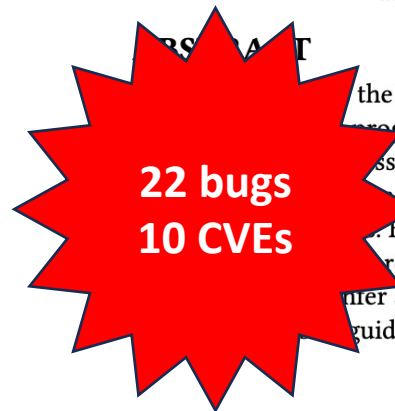
## Distributed System Fuzzing (2023)

Ruijie Meng<sup>\*†</sup>  
National University of Singapore  
Singapore  
ruijie\_meng@u.nus.edu

George Pirlea<sup>\*</sup>  
National University of Singapore  
Singapore  
gpirlea@comp.nus.edu.sg

Abhik Roychoudhury<sup>‡</sup>  
National University of Singapore  
Singapore  
abhik@comp.nus.edu.sg

Ilya Sergey  
National University of Singapore  
Singapore  
ilya@nus.edu.sg



the lightweight approach of choice for finding programs. It provides a balance between efficiency by conducting a biased random search over inputs using a feedback function from observations. For distributed system testing, however, the presented today by only black-box tools that enter and exploit any knowledge of the system's guide the search for bugs.

are generated in a purely random fashion, or it can be guided by knowledge of the program's internal structure (white-box). The most popular fuzzers are *grey-box*, where the search is guided by run-time observations of program behaviour, collected, as tests execute, for artefacts instrumented at compile time. Thanks to the ease of its deployment and use, grey-box fuzzing is the state-of-the-practice for automatically discovering bugs in sequential programs.

A common approach to finding bugs in distributed systems in practice is *stress-testing*, in which the system is subjected to

# Why is conformance hard?



Distributed systems are notoriously hard to get right. Protocol designers struggle to reason about concurrent execution on multiple machines, which leads to subtle errors. Engineers implementing such protocols face the same subtleties and, worse, must improvise to fill in gaps between abstract protocol descriptions and practical constraints, e.g., that real logs cannot grow without bound. Thorough testing is considered best practice, but its efficacy is limited by distributed systems' combinatorially large state spaces.

## **IronFleet: Proving Practical Distributed Systems Correct**

Chris Hawblitzel, Jon Howell, Manos Kapritsos, Jacob R. Lorch,  
Bryan Parno, Michael L. Roberts, Srinath Setty, Brian Zill

Microsoft Research

# Why is conformance hard?

- Medium of specification
  - Fully-fledged protocols are large and complex
    - [Basic Raft](#): 485 LoC
    - [TLC-optimized Raft](#): 653 LoC
    - [Raft with reconfiguration](#): 1083 LoC

```
\* Defines how the variables may transition.
Next == /\ \/\ \E i \in Server : Restart(i)
        \/\ \E i \in Server : Timeout(i)
        \/\ \E i,j \in Server : RequestVote(i, j)
        \/\ \E i \in Server : BecomeLeader(i)
        \/\ \E i \in Server, v \in Value : ClientRequest(i, v)
        \/\ \E i \in Server : AdvanceCommitIndex(i)
        \/\ \E i,j \in Server : AppendEntries(i, j)
        \/\ \E m \in DOMAIN messages : Receive(m)
        \/\ \E m \in DOMAIN messages : DuplicateMessage(m)
        \/\ \E m \in DOMAIN messages : DropMessage(m)
\* History variable that tracks every log ever:
/\ allLogs' = allLogs \cup {log[i] : i \in Server}
```

- State machines (TLA<sup>+</sup>) become harder to extend as they grow

# Why is conformance hard?

- Medium of specification
  - Fully-fledged protocols are large and complex
    - [Basic Raft](#): 485 LoC
    - [TLC-optimized Raft](#): 653 LoC
    - [Raft with reconfiguration](#): 1083 LoC
  - State machines (TLA<sup>+</sup>) become harder to extend as they grow
  - PlusCal solves some problems, but is not used much in practice
    - 25% of protocols in [github.com/tlaplus/Examples](https://github.com/tlaplus/Examples)

## The PlusCal Algorithm Language

# Why is conformance hard?

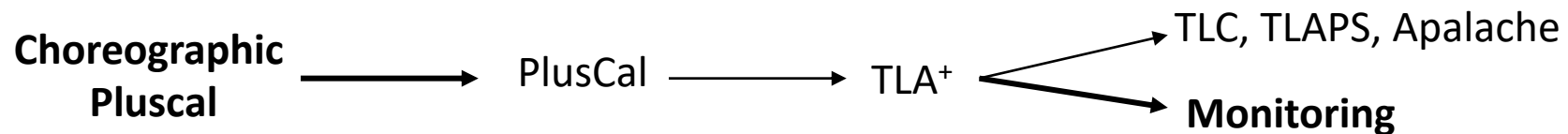
- *Implementations* are large and complex
  - Real-world Raft: etcd, 20k LoC, with concurrency, I/O, etc.
  - Implementation bugs can compromise protocol guarantees
  - Lack of lightweight tools for justifying parts of the implementation and supporting automated checks

# Challenges

1. Underspecification due to difficulty of extending large specifications
2. Conformance of real-world consensus implementations

# Contributions

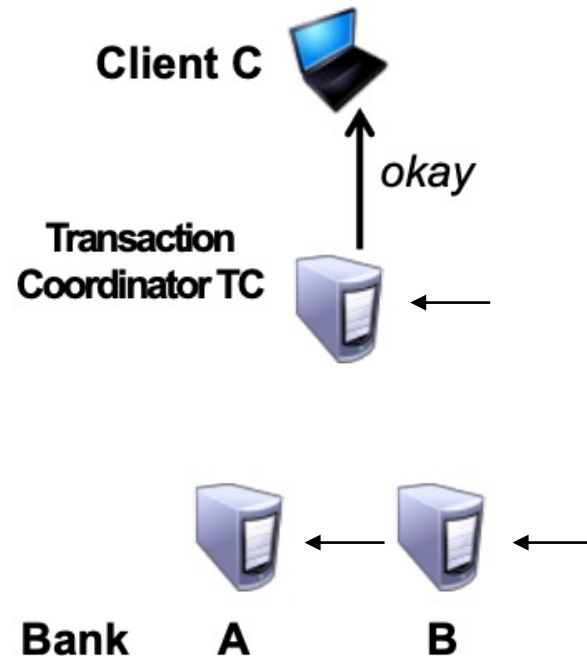
1. Choreographic PlusCal
2. Practical monitoring using *existing* TLA+ specifications



# Two-phase commit

## A correct atomic commit protocol

---



1.  $C \rightarrow TC$ : "go!"
  2.  $TC \rightarrow A, B$ : "prepare!"
  3.  $A, B \rightarrow TC$ : vote "yes" or "no"
  4.  $TC \rightarrow A, B$ : "commit!" or "abort!"
    - TC sends **commit** if both say yes
    - TC sends **abort** if either say no
  5.  $TC \rightarrow C$ : "okay" or "failed"
- A, B commit on receipt of commit message



# Two-phase commit in PlusCal

```
process (C \in coordinators)
  variables temp = participants,
            aborted = FALSE; {
  while (temp /= {}) {
    with (r \in temp) {
      Send(self, r, "prepare");
      temp := temp \ {r};
    };
  };
  temp := participants;
  while (temp /= {} \ / aborted) {
    with (r \in temp) {
      either {
        Receive(r, self, "prepared");
      } or {
        Receive(r, self, "abort");
        aborted := TRUE;
      };
      temp := temp \ {r};
    };
  };
  if (aborted) {
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
        Send(coord, r, "abort");
        temp := temp \ {r};
      };
    };
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
        Receive(r, coord, "aborted");
        temp := temp \ {r};
      };
    };
  } else {
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
        Send(coord, r, "commit");
        temp := temp \ {r};
      };
    };
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
        Receive(r, coord, "committed");
        temp := temp \ {r};
      };
    };
  }
} } }

process (P \in participants) {
  Receive(coord, self, "prepare");
  either {
    psend:
      Send(self, coord, "prepared");
  } or {
    Send(self, coord, "abort");
  };
  either {
    Receive(coord, self, "commit");
    Send(self, coord, "committed");
  } or {
    Receive(coord, self, "abort");
    Send(self, coord, "aborted");
  }
}
```

# Two-phase commit in PlusCal

```
process (C \in coordinators)
variables temp = participants,
          aborted = FALSE; {
while (temp /= {}) {
  with (r \in temp) {
    Send(self, r, "prepare");
    temp := temp \ {r};
  };
temp := participants;
while (temp /= {} \ / aborted) {
  with (r \in temp) {
    either {
      Receive(r, self, "prepared");
    } or {
      Receive(r, self, "abort");
      aborted := TRUE;
    };
    temp := temp \ {r};
  };
if (aborted) {
  temp := participants;
  while (temp /= {}) {
    with (r \in temp) {
      Send(coord, r, "abort");
      temp := temp \ {r};
    };
temp := participants;
while (temp /= {}) {
  with (r \in temp) {
    Receive(r, coord, "aborted");
    temp := temp \ {r};
  }
} else {
  temp := participants;
  while (temp /= {}) {
    with (r \in temp) {
      Send(coord, r, "commit");
      temp := temp \ {r};
    };
temp := participants;
while (temp /= {}) {
  with (r \in temp) {
    Receive(r, coord, "committed");
    temp := temp \ {r};
  }
}
} } }
```

```
process (P \in participants) {
  Receive(coord, self, "prepare");
  either {
    psend:
      Send(self, coord, "prepared");
  } or {
    Send(self, coord, "abort");
  };
  either {
    Receive(coord, self, "commit");
    Send(self, coord, "committed");
  } or {
    Receive(coord, self, "abort");
    Send(self, coord, "aborted");
  }
}
```

# Two-phase commit in PlusCal

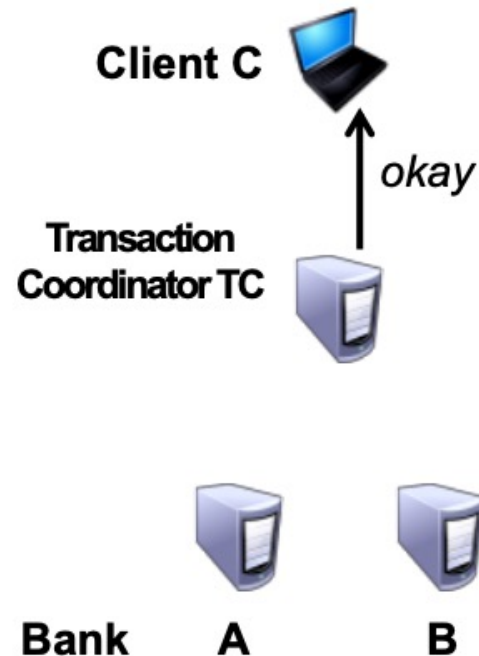
```
process (C \in coordinators)
  variables temp = participants,
             aborted = FALSE; {
  while (temp /= {}) {
    with (r \in temp) {
      Send(self, r, "prepare");
      temp := temp \ {r};
    } };
  temp := participants;
  while (temp /= {} \ / aborted) {
    with (r \in temp) {
      either {
        Receive(r, self, "prepared");
      } or {
        Receive(r, self, "abort");
        aborted := TRUE;
      };
      temp := temp \ {r};
    } };
  if (aborted) {
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
        Send(coord, r, "abort");
        temp := temp \ {r};
        ...
      }
    }
  }
}

process (P \in participants) {
  Receive(coord, self, "prepare");
  either {
    psend:
    Send(self, coord, "prepared");
  } or {
    Send(self, coord, "abort");
  };
  either {
    Receive(coord, self, "commit");
    Send(self, coord, "committed");
  } or {
    Receive(coord, self, "abort");
    Send(self, coord, "aborted");
  }
}
```

# Two-phase commit

## A correct atomic commit protocol

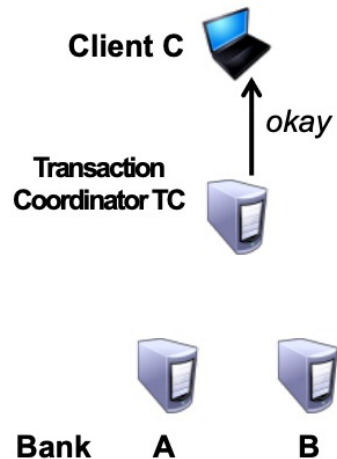
---



1.  $C \rightarrow TC$ : "go!"
  2.  $TC \rightarrow A, B$ : "prepare!"
  3.  $A, B \rightarrow TC$ : vote "yes" or "no"
  4.  $TC \rightarrow A, B$ : "commit!" or "abort!"
    - TC sends **commit** if both say yes
    - TC sends **abort** if either say no
  5.  $TC \rightarrow C$ : "okay" or "failed"
- A, B commit on receipt of commit message

# Choreographic PlusCal: choreography

## A correct atomic commit protocol



1. C → TC: "go!"
  2. TC → A, B: "prepare!"
  3. A, B → TC: vote "yes" or "no"
  4. TC → A, B: "commit!" or "abort!"
    - TC sends **commit** if both say yes
    - TC sends **abort** if either say no
  5. TC → C: "okay" or "failed"
- A, B commit on receipt of commit message

```
choreography
(P \in participants),
(C \in coordinators) {

  all (p \in participants) {
    Transmit(coord, p, "prepare");
    either {
      Transmit(p, coord, "prepared");
    } or {
      Transmit(p, coord, "aborted");
    } } };
  if (aborted) {
    all (p \in participants) {
      Transmit(coord, p, "abort");
      Transmit(p, coord, "aborted");
    }
  } else {
    all (p \in participants) {
      Transmit(coord, p, "commit");
      Transmit(p, coord, "committed");
    } } }
}
```

# Choreographic PlusCal: all

```
process (C \in coordinators)
  variables temp = participants,
             aborted = FALSE; {
    while (temp /= {}) {
      with (r \in temp) {
        Send(self, r, "prepare");
        temp := temp \ {r};
      };
    }
  }
temp := participants;
while (temp /= {} \ / aborted) {
  with (r \in temp) {
    either {
      Receive(r, self, "prepared");
    } or {
      Receive(r, self, "abort");
      aborted := TRUE;
    };
    temp := temp \ {r};
  };
}
if (aborted) {
  temp := participants;
  while (temp /= {}) {
    with (r \in temp) {
```

Encoding a  
multicast

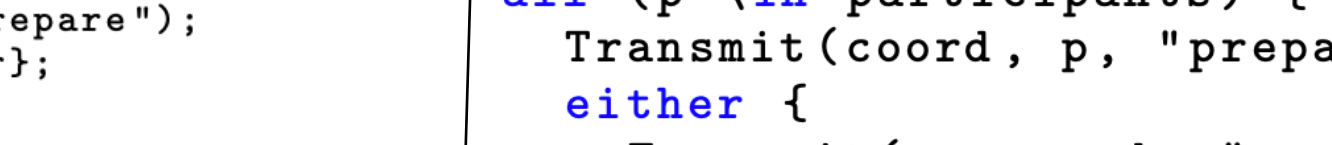
```
all (p \in participants) {
  Transmit(coord, p, "prepare");
  either {
    Transmit(p, coord, "prepared");
  } or {
    Transmit(p, coord, "aborted");
    cancel "phase1";
  }
}
```

# Choreographic PlusCal: task and cancel

```
process (C \in coordinators)
  variables temp = participants,
            aborted = FALSE; {
  while (temp /= {}) {
    with (r \in temp) {
      Send(self, r, "prepare");
      temp := temp \ {r};
    };
  temp := participants;
  while (temp /= {} \ / aborted) {
    with (r \in temp) {
      either {
        Receive(r, self, "prepared");
      } or {
        Receive(r, self, "abort");
        aborted := TRUE;
      };
      temp := temp \ {r};
    };
  };
  if (aborted) {
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
```

task coordinators "phase1" {  
 all (p \in participants) {  
 Transmit(coord, p, "prepare");  
 either {  
 Transmit(p, coord, "prepared");  
 } or {  
 Transmit(p, coord, "aborted");  
 cancel "phase1";  
 } } };

Stop if  
participant  
aborts



# Choreographic PlusCal

```
process (C \in coordinators)
  variables temp = participants,
             aborted = FALSE; {
  while (temp /= {}) {
    with (r \in temp) {
      Send(self, r, "prepare");
      temp := temp \ {r};
    };
  };
  temp := participants;
  while (temp /= {} \ / aborted) {
    with (r \in temp) {
      either {
        Receive(r, self, "prepared");
      } or {
        Receive(r, self, "abort");
        aborted := TRUE;
      };
      temp := temp \ {r};
    };
  };
  if (aborted) {
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
        Send(coord, r, "abort");
        temp := temp \ {r};
      };
    };
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
        Receive(r, coord, "aborted");
        temp := temp \ {r};
      };
    };
  } else {
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
        Send(coord, r, "commit");
        temp := temp \ {r};
      };
    };
    temp := participants;
    while (temp /= {}) {
      with (r \in temp) {
        Receive(r, coord, "committed");
        temp := temp \ {r};
      };
    };
  } } }
```

```
process (P \in participants) {
  Receive(coord, self, "prepare");
  either {
    psend:
    Send(self, coord, "prepared");
  } or {
    Send(self, coord, "abort");
  };
  either {
    Receive(coord, self, "commit");
    Send(self, coord, "committed");
  } or {
    Receive(coord, self, "abort");
    Send(self, coord, "aborted");
  } } }
```

## choreography

```
(P \in participants),
(C \in coordinators) {
  task coordinators "phase1" {
    all (p \in participants) {
      Transmit(coord, p, "prepare");
      either {
        Transmit(p, coord, "prepared");
      } or {
        Transmit(p, coord, "aborted");
        cancel "phase1";
      } } };
  if (aborted) {
    all (p \in participants) {
      Transmit(coord, p, "abort");
      Transmit(p, coord, "aborted");
    };
  }
} else {
  all (p \in participants) {
    Transmit(coord, p, "commit");
    Transmit(p, coord, "committed");
  } } }
```



# Choreographic PlusCal

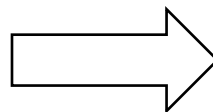
<b>Protocol</b>	<b>Ch. PlusCal</b>	<b>TLA<sup>+</sup></b>
Two-phase commit [23]	23	66
Non-blocking atomic commit [35]	36	96
Raft leader election [32]	46	186

Table 1: Relative specification sizes (LoC)

```
choreography
(P \in participants),
(C \in coordinators) {
  task coordinators "phase1" {
    all (p \in participants) {
      Transmit(coord, p, "prepare");
      either {
        Transmit(p, coord, "prepared");
      } or {
        Transmit(p, coord, "aborted");
        cancel "phase1";
      } } };
  if (aborted) {
    all (p \in participants) {
      Transmit(coord, p, "abort");
      Transmit(p, coord, "aborted");
    }
  } else {
    all (p \in participants) {
      Transmit(coord, p, "commit");
      Transmit(p, coord, "committed");
    } } }
```

# Projection & monitoring

```
choreography
(P \in participants),
(C \in coordinators) {
task coordinators "phase1" {
  all (p \in participants) {
    Transmit(coord, p, "prepare");
  } either {
    Transmit(p, coord, "prepared");
  } or {
    Transmit(p, coord, "aborted");
    cancel "phase1";
  } } };
if (aborted) {
  all (p \in participants) {
    Transmit(coord, p, "abort");
    Transmit(p, coord, "aborted");
  }
} else {
  all (p \in participants) {
    Transmit(coord, p, "commit");
    Transmit(p, coord, "committed");
  } } }
```



```
process (C \in coordinators)
variables temp = participants,
          aborted = FALSE; {
while (temp /= {}) {
  with (r \in temp) {
    Send(self, r, "prepare");
    temp := temp \ {r};
  } };
temp := participants;
while (temp /= {} \/\ aborted) {
  with (r \in temp) {
    either {
      Receive(r, self, "prepared");
    } or {
      Receive(r, self, "abort");
      aborted := TRUE;
    };
    temp := temp \ {r};
  } };
if (aborted) {
  temp := participants;
  while (temp /= {}) {
    with (r \in temp) {
      Send(coord, r, "abort");
      temp := temp \ {r};
    } };
temp := participants;
while (temp /= {}) {
  with (r \in temp) {
    Receive(r, coord, "aborted");
    temp := temp \ {r};
  } }
} else {
  temp := participants;
  while (temp /= {}) {
    with (r \in temp) {
      Send(coord, r, "commit");
      temp := temp \ {r};
    } }
temp := participants;
while (temp /= {}) {
  with (r \in temp) {
    Receive(r, coord, "committed");
    temp := temp \ {r};
  } } } }
```

```
process (P \in participants) {
  Receive(coord, self, "prepare");
  either {
    psend:
      Send(self, coord, "prepared");
  } or {
    Send(self, coord, "abort");
  };
  either {
    Receive(coord, self, "commit");
    Send(self, coord, "committed");
  } or {
    Receive(coord, self, "abort");
    Send(self, coord, "aborted");
  } } }
```

# Projection & monitoring

- Choreographic languages/logics (e.g. session types) have a *projection* operation to derive local programs for verification, monitoring, and/or code generation

$$\mathit{project}(a \rightarrow b) = \{! b, ? a\}$$

- We define projection across *both* Choreographic PlusCal and TLA<sup>+</sup>
  - Integrates with existing toolchain
  - Monitoring works for vanilla TLA<sup>+</sup> as well (assuming some syntactic conditions)

# Projection & monitoring

Ch. PlusCal	<b>VARIABLES</b> v, inbox, outbox	<b>choreography</b> (a \in A, b \in B) { Transmit(a, b, v, "msg") }	
PlusCal	<b>VARIABLES</b> v, inbox, outbox	<b>process</b> (a \in A) { Send(b, "msg") }	<b>process</b> (b \in B) { v = Receive(a) }
TLA+	<b>VARIABLES</b> v, inbox, outbox	A_send(self, b) == /\ Send(b, "msg") /\ UNCHANGED <<inbox, v>>	B_send(self, a) == /\ v'[self] := Receive(a) /\ UNCHANGED <<outbox>>
Multiple TLA+ models	<b>VARIABLES</b> inbox, outbox A_send(b) == /\ Send(b, "msg") /\ UNCHANGED <<inbox>>		<b>VARIABLES</b> v, inbox, outbox B_send(a) == /\ v := Receive(a) /\ UNCHANGED <<outbox>>

# Monitoring



Choreographic PlusCal

## The PlusCal Algorithm Language

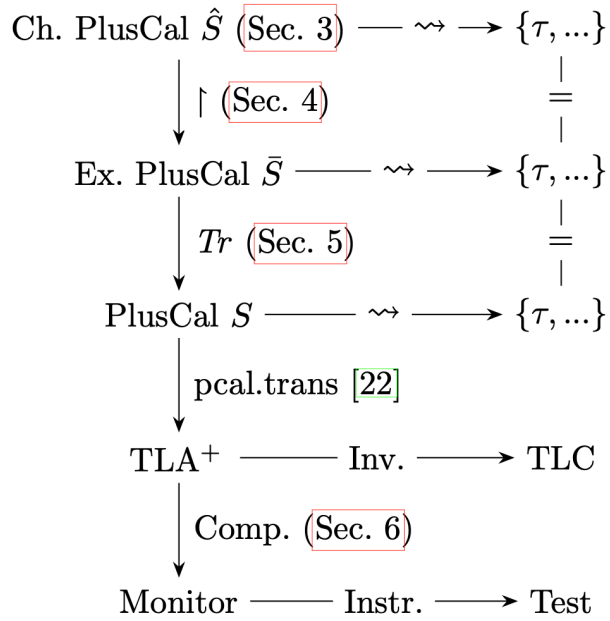


Fig. 4: Overview

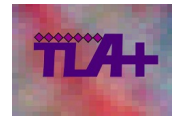


$\models \varphi$



Checking a Model

TLA+ Proof System



$\equiv$



$\equiv$



# Monitoring

- Instrument system to collect traces
  - Refinement mapping
    - Function from concrete to abstract state
    - Abstracts away details, reinterprets behavior in terms of the model's
    - May require auxiliary state to define
    - Deep embedding of TLA+ formulae in Go

```
type TLA interface {  
    String() string  
    MarshalJSON() ([]byte, error)  
}
```

# Monitoring

- Instrument system to collect traces
  - Refinement mapping
  - Linearization points
    - Program locations where state changes become visible
    - Can vary significantly between implementations
    - May require auxiliary state to define

# Monitoring

- Instrument system to collect traces
- Validate behaviors
  - Model-based trace checking [Pressler 18, Davis 20]
  - Compile model into monitor and validate on the fly
    - Offline, also online
    - Scalable, possible to enable in production/fuzzing

<b>Project</b>	<b>Protocol</b>	<b>LoC</b>	<b>Overhead</b>
vadiminshakov/committer	2PC	3032	19% (5 ms)
etcd-io/raft	Raft leader election	21,064	2% (4 ms)

Table 2: Monitor overhead



# Conclusion

- Choreographic PlusCal + monitoring
- What's in the paper?
  - Details, formalization, soundness of new features and projection
- Future work
  - Liveness: runtime verification
  - New classes of protocols, e.g. role-parametric
  - User-provided refinement mapping and linearization points are all trusted – statically check

# Thank you!

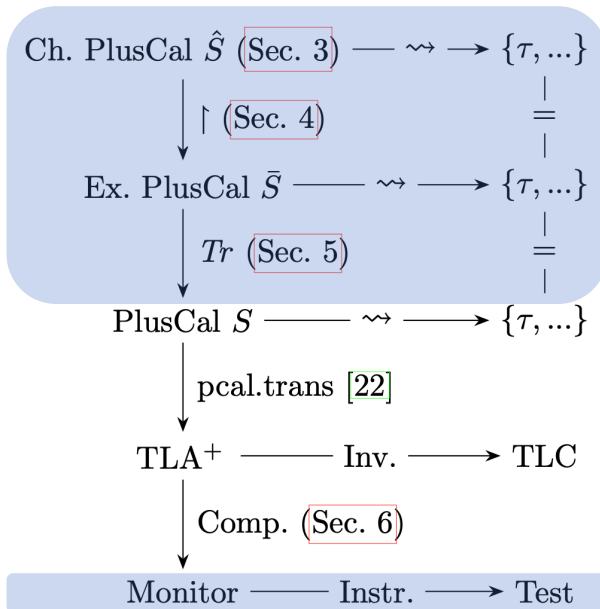
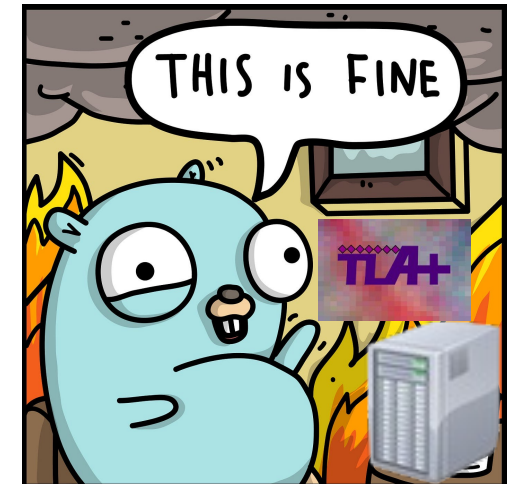


Fig. 4: Overview

## Protocol Conformance with Choreographic PlusCal

Darius Foo, Andreea Costea, and Wei-Ngan Chin

National University of Singapore  
 {dariusf, andreeac, chinwn}@comp.nus.edu.sg



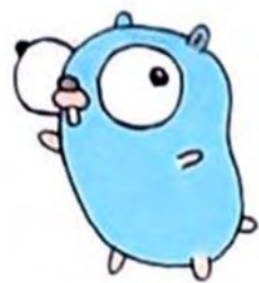
<https://github.com/dariusf/tlaplus/tree/cpcal>



# Monitoring

```
func psend(prev state, this state, self TLA) bool {
    if !(reflect.DeepEqual(prev.pc, Str("psend"))) {
        return false
    }
    // ... outbox check elided
    if !(reflect.DeepEqual(this.pc, Str("Lb1_2"))) {
        return false
    }
    return true
}
```

Fig. 6: Go rendering of psend in generated monitor



# Specification with TLA<sup>+</sup>

```
\* Defines how the variables may transition.
Next == /\ \E i \in Server : Restart(i)
        \E i \in Server : Timeout(i)
        \E i,j \in Server : RequestVote(i, j)
        \E i \in Server : BecomeLeader(i)
        \E i \in Server, v \in Value : ClientRequest(i, v)
        \E i \in Server : AdvanceCommitIndex(i)
        \E i,j \in Server : AppendEntries(i, j)
        \E m \in DOMAIN messages : Receive(m)
        \E m \in DOMAIN messages : DuplicateMessage(m)
        \E m \in DOMAIN messages : DropMessage(m)
\* History variable that tracks every log ever:
/\ allLogs' = allLogs \cup {log[i] : i \in Server}
```

# Specification with TLA<sup>+</sup>

Figuring out how actions are related is tedious, e.g. sequentially

```
\* Server i times out and starts a new election.
Timeout(i) == /\ state[i] \in {Follower, Candidate}
              /\ state' = [state EXCEPT ![i] = Candidate]
              /\ currentTerm' = [currentTerm EXCEPT ![i] = currentTerm[i] + 1]
              \* Most implementations would probably just set the local vote
              \* atomically, but messaging localhost for it is weaker.
              /\ votedFor' = [votedFor EXCEPT ![i] = Nil]
              /\ votesResponded' = [votesResponded EXCEPT ![i] = {}]
              /\ votesGranted' = [votesGranted EXCEPT ![i] = {}]
              /\ voterLog' = [voterLog EXCEPT ![i] = [j \in {} |-> <<>>]]
              /\ UNCHANGED <<messages, leaderVars, logVars>>
```

```
\* Candidate i sends j a RequestVote request.
RequestVote(i, j) ==
  /\ state[i] = Candidate
  /\ j \notin votesResponded[i]
  /\ Send([mtype |-> RequestVoteRequest,
          mterm |-> currentTerm[i],
          mlastLogTerm |-> LastTerm(log[i]),
          mlastLogIndex |-> Len(log[i]),
          msource |-> i,
          mdest |-> j])
  /\ UNCHANGED <<serverVars, candidateVars, leaderVars, logVars>>
```

... must also check if other actions are enabled in Candidate state, else nondeterminism

# Specification with TLA<sup>+</sup>

Figuring out how actions are related is tedious, e.g. send-receive

```

/* Server i receives a RequestVote request from server j with
/* m.mterm <= currentTerm[i].
HandleRequestVoteRequest(i, j, m) ==
  LET logOk == \V m.mlastLogTerm > LastTerm(log[i])
              \V \A m.mlastLogTerm = LastTerm(log[i])
              \A m.mlastLogIndex >= Len(log[i])
      grant == \A m.mterm = currentTerm[i]
              \A logOk
              \A votedFor[i] \in {Nil, j}
  IN \A m.mterm <= currentTerm[i]
      \A \V grant \A votedFor' = [votedFor EXCEPT ![i] = j]
      \V ~grant \A UNCHANGED votedFor

  \A Reply([mtype      |-> RequestVoteResponse,
           mterm       |-> currentTerm[i],
           mvoteGranted |-> grant,
           /* mlog is used just for the `elections' history variable for
           /* the proof. It would not exist in a real implementation.
           mlog         |-> log[i],
           msource     |-> i,
           mdest       |-> j],
           m)

  \A UNCHANGED <<state, currentTerm, candidateVars, leaderVars, logVars>>
```

```

/* Candidate i sends j a RequestVote request.
RequestVote(i, j) ==
  \A state[i] = Candidate
  \A j \notin votesResponded[i]
  \A Send([mtype      |-> RequestVoteRequest,
          mterm       |-> currentTerm[i],
          mlastLogTerm |-> LastTerm(log[i]),
          mlastLogIndex |-> Len(log[i]),
          msource      |-> i,
          mdest        |-> j])
  \A UNCHANGED <<serverVars, candidateVars, leaderVars, logVars>>
```

Must do this repeatedly to get a sense of the flow of the protocol



# Specification with TLA<sup>+</sup>

## Non-compositionality

```
\* Add a message to the bag of messages.
Send(m) == messages' = WithMessage(m, messages)

\* Candidate i sends j a RequestVote request.
RequestVote(i, j) ==
  /\ state[i] = Candidate
  /\ j \notin votesResponded[i]
  /\ Send([mtype      |-> RequestVoteRequest,
          mterm       |-> currentTerm[i],
          mlastLogTerm |-> LastTerm(log[i]),
          mlastLogIndex |-> Len(log[i]),
          msource      |-> i,
          mdest        |-> j])
  /\ UNCHANGED <<serverVars, candidateVars, leaderVars, logVars>>
```

Must thread state through functions manually

# Linking specification to implementation

Many “industrial-grade” unverified protocol implementations...

<b>etcd-io/etcd</b>	Distributed reliable key-value store for the most critical data of a distributed system		
<b>W</b>	<b>baidu/braft</b>	An industrial-grade C++ implementation of RAFT consensus algorithm based on braft.	
<b>ht</b>			
<b>Ad</b>	<b>wid</b>	<b>Tencent/phxpaxos</b>	
	<b>Sta</b>	The Paxos library implemented in C++ that has been used in the WeChat production environment.	
	<b>290</b>		
	<b>Add</b>	<b>Stars</b>	<b>Language</b>
		3045	C++
		Added by GitHub	

# Linking specification to implementation

... many specifications as well, but unrelated

## List of Examples

No	Name	Short description	Spec's authors	TLAPS Proof	TLC Check	
39	MultiPaxos	<a href="#">The abstract specification of Generalized Paxos (Lamport, 2004)</a>	Giuliano Losa		✓	
45	Paxos	<a href="#">Paxos consensus algorithm (Lamport, 1998)</a>	Leslie Lamport		✓	
47	raft	<a href="#">Raft consensus algorithm (Ongaro, 2014)</a>	Diego Ongaro		✓	
57	transaction_commit	<a href="#">Consensus on transaction commit (Gray &amp; Lamport, 2006)</a>	Leslie Lamport		✓	
67	Tencent-Paxos	<a href="#">PaxosStore: high-availability storage made practical in WeChat. Proceedings of the VLDB Endowment (Zheng et al., 2017)</a>	Xingchen Yi, Hengfeng Wei	✓	✓	
59	TwoPhase	<a href="#">Two-phase handshaking</a>	Leslie Lamport, Stephan Merz		✓	Nat
62	Misra Reachability Algorithm	<a href="#">Misra Reachability Algorithm</a>	Leslie Lamport	✓	✓	Int, Seq, FiniteS TLC, TLAPS, NaturalsInducti
63	Loop Invariance	<a href="#">Loop Invariance</a>	Leslie Lamport	✓	✓	Int, Seq, FiniteS TLC, TLAPS, SequenceTheor NaturalsInducti
69	Paxos	<a href="#">Paxos</a>			✓	Int, FiniteSets
75	Lock-Free Set	<a href="#">PlusCal spec of a lock-free set used by TLC</a>	Markus Kuppe		✓	Sequences, FiniteSets, Integ TLC
77	ParallelRaft	<a href="#">A variant of Raft</a>	Xiaosong Gu, Hengfeng Wei, Yu Huang		✓	Integers, FiniteS Sequences, Naturals
83	Raft (with cluster changes)	<a href="#">Raft with cluster changes, and a version with Apache type annotations but no cluster changes</a>	George Pirlea, Darius Foo, Brandon Amos, Huanchen Zhang, Daniel Ricketts		✓	Functions, SequencesExt, FiniteSetsExt, TypedBags