# Programming abstractions for automating the verification of replicated state machine

Cezara Drăgoi, INRIA Paris/ENS

joint work with Tom Henzinger, IST Austria Josef Widder, TU Wien Damien Zufferey, MPI Kaiserslautern

#### Replication



#### Replicated state machine



#### Zab Viewstamped

#### Goal

#### Design automated verification for implementations of replicated state machine

## Difficulties

- Impossibility result [FLP'85]: "Consensus cannot be reached in asynchronous networks in the presence of at least one faulty process"
  - protocols have been developed for different network assumptions: Different degrees of synchrony and faults

• asynchronous parallel composition



- asynchronous parallel composition
- communication via message passing: peerto-peer, broadcast
- timer constraints: how long does a process wait for a message? what if two messages are received in the reverse sending order ?



- asynchronous parallel composition
- communication via message passings: peer-to-peer, broadcast
- timer constraints: how long does a process wait for a message? what if two messages are received in the reverse order ?
- **UNDOUNDED DUFFERS:** A process receives a bunch of messages that momentarily it does not need.



- asynchronous parallel composition
- communication via message passings: peer-to-peer, broadcast
- timer constraints: how long does a process wait for a message? what if two messages are received in the reverse order ?
- **UNDOUNDED DUFFERS:** A process receives a bunch of messages that momentarily it does not need.
- faults: message drop, process-crash, messages are corrupted



- asynchronous parallel composition
- communication via message passings: peer-to-peer, broadcast
- timer constraints: how long does a process wait for a message? what if two messages are received in the reverse order ?
- **UNDOUNDED DUFFERS:** A process receives a bunch of messages that momentarily it does not need.
- faults: message drop, process-crash, messages are corrupted



## State of the art in verification of RSM

Mechanized verification:

- Verdi, EventML(2012) (only safety)
- IronFleet(2013)

- TLA+

Automated verification:

- model checking of simple algorithms or protocols using minor coordination

- lvy(2016)
- Psync

#### Goal

#### Programming abstraction

simpler source code
+ specifications

Automated Verification

#### Example



Assumptions: synchronous communication and no faults

there is a bound **b** on the message delay

**Assumptions**: the leader receives the messages from all processes and every process hears from the leader

### **Example: Communication**





## **Example: Communication**



the leader hears from all processes





#### Example



We would like a round based model that separates the *network properties* from the <u>algorithmic computation</u>

the leader hears from all processes, every process hears from the leader

<u>compute the</u> <u>minimum value,</u> <u>a commun prefix, etc.</u>

### Programming model

# 1. Separates the *network properties* from the <u>algorithmic</u> <u>computation</u>

2. Faults and asynchrony simulated by an adversarial environment

3. We want a **simple programming abstraction** that offers a synchronous image of the system and compiles into **executable asynchronous code**.

#### Communication close rounds[Erlad,Francez'83]



A round is a computation step that

- · defines the interaction between the processes that participate in that step,
- · processes synchronize at the boundary between rounds,
- it gives a logical unit of time.

**Communication close** = two rounds do not communicate with each other, i.e., messages are scoped in the round.

#### Communication close rounds[Erlad,Francez'83]



A round is a computation step that

- · defines the interaction between the processes that participate in that step,
- · processes synchronize at the boundary between rounds,
- it gives a logical unit of time.

**Communication close** = two rounds do not communicate with each other, i.e., messages are scoped in the round.

#### HO-model [Charrone-Bost, Schiper'09]





Each process p has a variable, *HO(p)=* the set of processes p hears-from *HO(p)* is non-deterministically chosen by an adversarial environment.

1. Determines the delivered messages



Each process p has a variable, *HO(p)=* the set of processes p hears-from *HO(p)* is non-deterministically chosen by an adversarial environment.

- 1. Determines the delivered messages
- 2. The network assumptions are stated over the HO variables

21



 $\Box$  ( $\forall$ p. p \in HO(leader)  $\land$  leader \in HO(p))

## HO-model



 $\Box$  ( $\forall$  p. p \in HO(Leader)  $\land$  Leader \in HO(p))

 $\square$  ( $\forall p. p \in HO(leader) \land leader \in HO(p)$ )

22

Approach



#### PSync Program Structure [popl'16]



• Round<sub>T</sub>

Send: () 
$$\rightarrow$$
 [Id  $\rightarrow$  T] Comm Pred Update: [Id  $\rightarrow$  T]  $\rightarrow$  ()

#### Example



#### **Example: Last Voting Algorithm**



```
new Round[(Int,Time)]{
```

```
def send(): Map[ProcessID, (Int,Time)] = Map( leader -> (x, ts) )
  def update(mailbox: Map[ProcessID, (Int,Time)]) {
    if (id == leader && mailbox.size > n/2) {
      vote = mailbox.maxBy(_._2._2)._2._1 // value with maximal ts
      commit = true
    }
}
```

#### Outline



#### Runtime: Round switch

![](_page_27_Figure_1.jpeg)

#### Runtime: Round switch

![](_page_28_Figure_1.jpeg)

#### Partial synchrony

T

![](_page_29_Figure_1.jpeg)

#### Network

		Respects liveness assumptions		າຣ	Respects liveness assum	
	Asynchronous	Synch	ironous	Asynchronous	Synchronous	
Runtime Execution						
	• ·					
	Asynchronous	Asynchronous	Synchronous			

#### **Runtime Correctness**

![](_page_30_Figure_1.jpeg)

For any Psync program P, the runtime semantics of P observationally refines the HO semantics of P, if the client is commutative.

#### Clients II Runtime(P) $\subseteq$ Clients II HO(P)

#### Outline

![](_page_31_Figure_1.jpeg)

#### **Psync: Benefits for Verification**

Reason about rounds in isolation. Lockstep semantics, no interleaving.

![](_page_32_Figure_2.jpeg)

Simple invariants that reason at the round boundaries, no messages are in flight, only the local states matter.

#### Psync: semi-automated verification

![](_page_33_Figure_1.jpeg)

#### Hoare-style Verification [vmcai'14]

![](_page_34_Figure_1.jpeg)

Safety: Inductive invariant checking Liveness: Variant functions

#### **Invariant for Agreement**

$$\begin{array}{l} \forall i. \neg decided(i) \land \neg ready(i) \\ \forall \exists v, t, A. \ A = \{ i. \ ts(i) > t \} \land |A| > n/2 \\ \land \quad \forall i. i \in A \Rightarrow x(i) = v \\ \land \quad \forall i. decided(i) \Rightarrow x(i) = v \\ \land \quad \forall i. commit(i) \lor ready(i) \Rightarrow vote(i) = v \\ \land \quad t \leq \Phi \\ \land \quad \forall i. \ ts(i) = \Phi \Rightarrow commit(coord) \end{array}$$

## **Specification logic**

- is able to express:
  - -- properties of sets of processes in the network
  - -- cardinality constraints
  - -- properties of the data values stored by each process
  - -- **JV** quantifier alternation
  - captures the transition relation of algorithms in the HO-model
- has a semi-decision procedure for checking entailment
- has a decidable satisfiability problem for a fragment Cl<sub>dec</sub>

#### **Psync: verification correctness**

![](_page_37_Figure_1.jpeg)

Given a specification S closed under indistinguishability, if a Psync program P satisfies S then the asynchronous semantics of P refines S.

#### Clients II Runtime(P) $\subseteq$ Clients II Spec(P)

#### Do Algorithms use Rounds?

Algorithm	LOC	Use rounds	Asynchronous
One third rule	52	$\checkmark$	$\checkmark$
Last Voting (Paxos)	89	$\checkmark$	$\checkmark$
Flood min consensus	24	$\checkmark$	×
Ben-Or randomized consensus	56	$\checkmark$	$\checkmark$
K-set agreement	42	$\checkmark$	$\checkmark$
K-set agreement early stopping	33	$\checkmark$	×
Lattice agreement	34	×	$\checkmark$
	54	$\checkmark$	$\checkmark$
Two phases commit	53	$\checkmark$	×
Eager reliable broadcast	36	×	×
Chandra Toueg	121	$\checkmark$	$\checkmark$
Generalized Paxos	152	$\checkmark$	$\checkmark$
ViewStampted	91	×	$\checkmark$

#### Code Size (Easy to Implement)

Paxos in	LOC	Executable	Verification
PSync	89	$\checkmark$	Semi-
DistAlgo	43	$\checkmark$	×
Distal	157	$\checkmark$	×
Overlog	107	$\checkmark$	×
TLA+	53	×	Interactive
IO Automata	142	×	Interactive
EventML	1729 N	$\checkmark$	Interactive
Verdi (Raft)	520	$\checkmark$	Interactive

#### Performance and Verification

Implementation	Year	Language	Throughput (x 1000 req./s)
Last Voting in PSync	2015	Scala	170
Egalitarian Paxos	2013	Go	450
Paxos in Distal	2013	Scala	150
JPaxos / SPaxos	2012	Java	75 / 300
Paxos for system builder	2008	С	40

Verification of	# Invariants (LOC)	# VCs	Solving time in
One third rule	4 (23)	27	5
Last Voting	8 (35)	45	16

### Conclusions

PSync uses a simple programming abstraction: Communication-closed rounds Separates the algorithm from the network requirements Asynchrony and faults are modelled by an adversary that drops messages

Runtime:

Asynchronous semantics refines the lockstep semantics

Preserves strong consistency

Can be implemented efficiently

Automated verification becomes possible

![](_page_41_Figure_7.jpeg)