

PSYNC: A partially synchronous language for fault-tolerant distributed algorithms

Cezara Drăgoi, Thomas A. Henzinger, Damien Zufferey

(POPL '16)

Overview

PSYNC

- is a domain-specific language (DSL)
- for programming fault-tolerant distributed algorithms
- in a high-level, round-based model with lockstep semantics

- which compiles into ***equivalent***, efficient asynchronous programs
- and is suitable for semi-automated verification

PSYNC is **descriptive** enough
to specify lots of dist. algos



PSYNC is **sufficiently restricted** so programs
can be semi-automatically proven

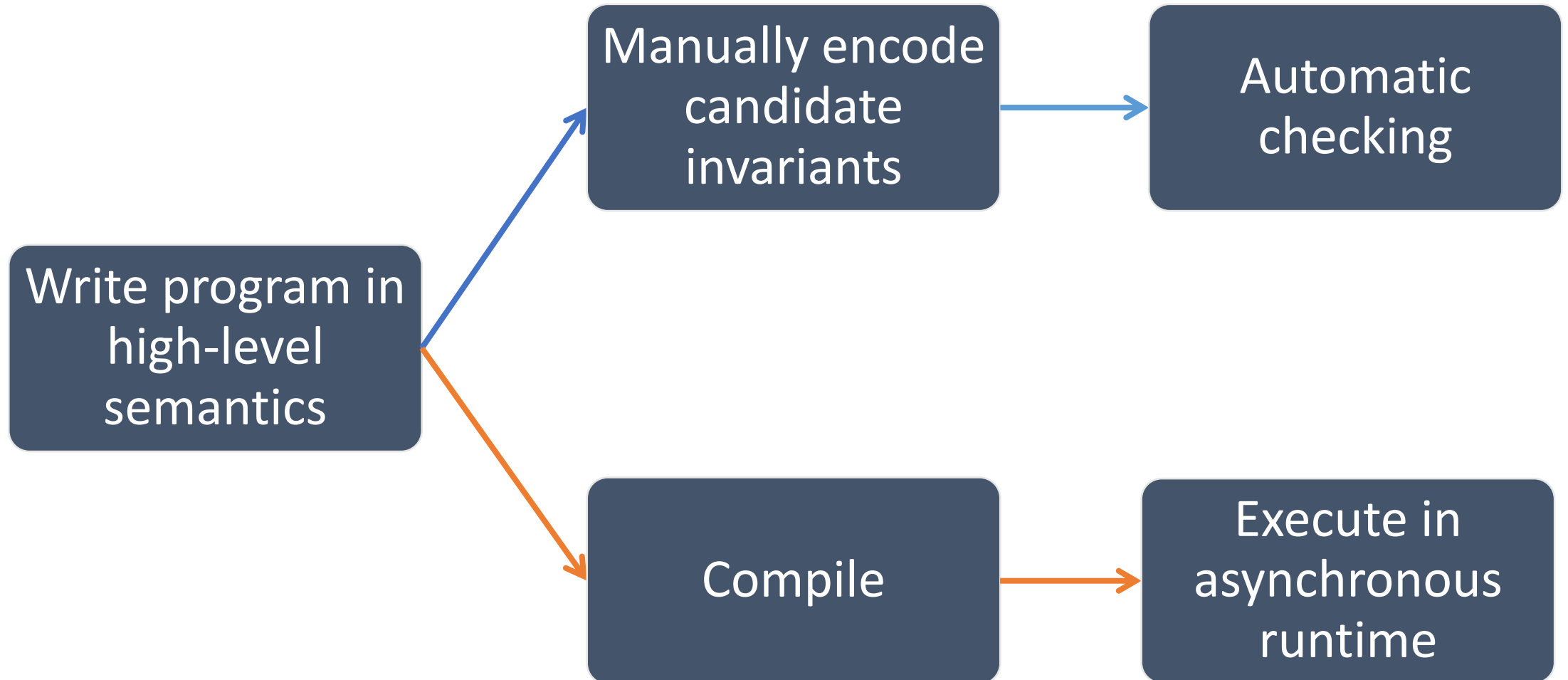


Unify **modelling, programming** and **verification** of fault-tolerant distributed algorithms



PSYNC can be **compiled**
into asynchronous programs

Workflow



Why does it matter?

- Fault-tolerant distributed systems are complex:
 - packets get re-ordered, duplicated or dropped
 - machines fail & potentially restart
 - concurrency is hard
- For critical applications, we want assured *correctness*
 - we need **formal verification**
 - but want development and verification to be easy
- PSYNC provides a solution

Internals

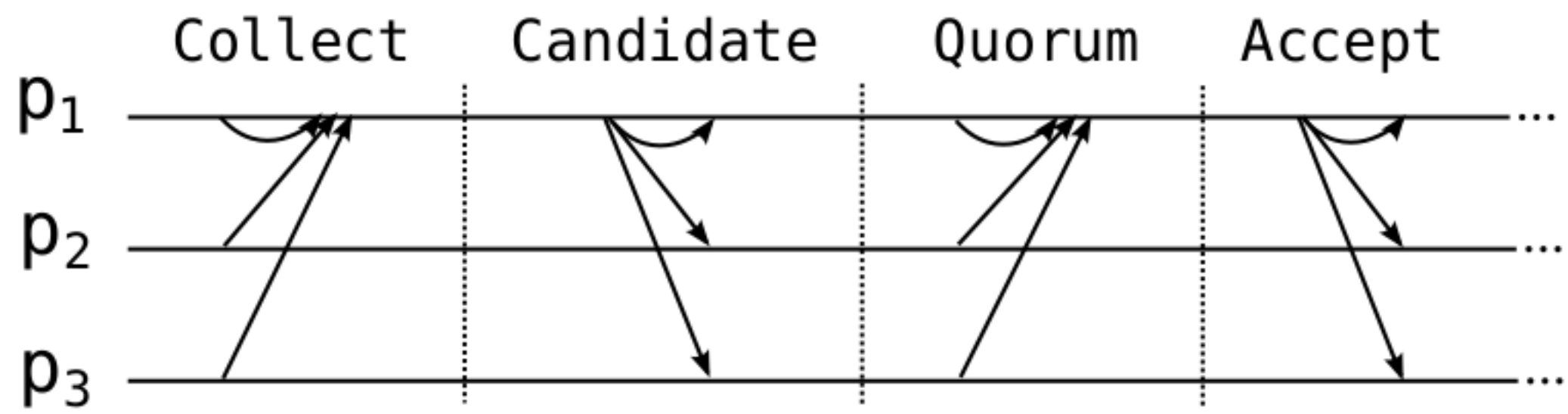
How does it work?

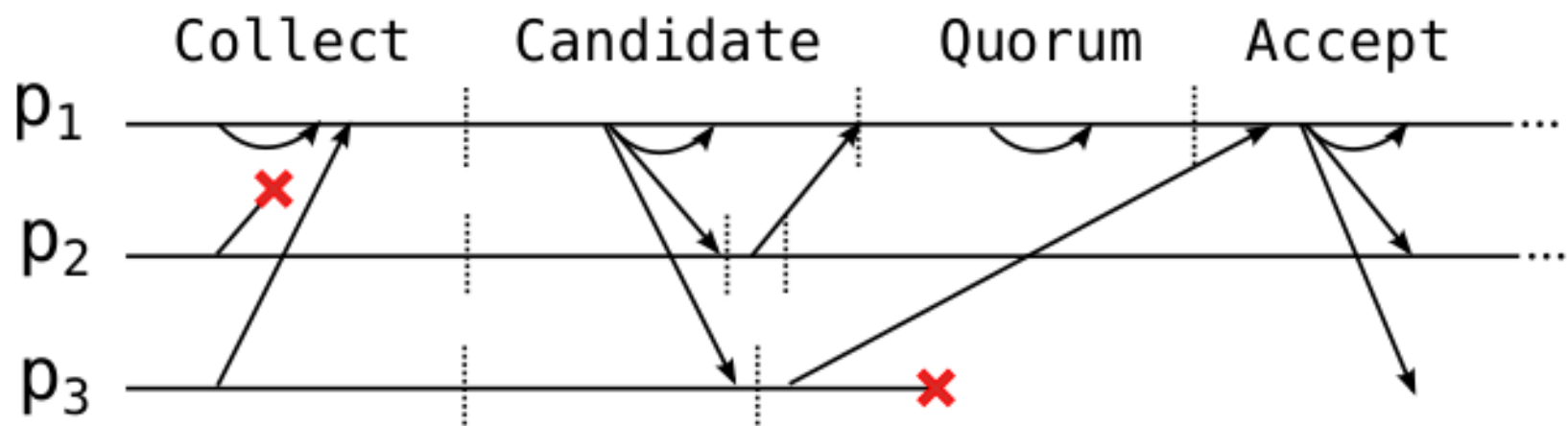
- **Communication-closed rounds with lockstep semantics:**
 - a PSYNC program is defined as a sequence of rounds
 - all processes execute the same round
 - messages sent within a round are either delivered in that round or dropped forever
- Each round consists of two operations, executed in this order:
 1. Send – send messages
 2. Update – update local state based on messages received in this round

PSYNC is based on the Heard-Of (\mathcal{HO}) model

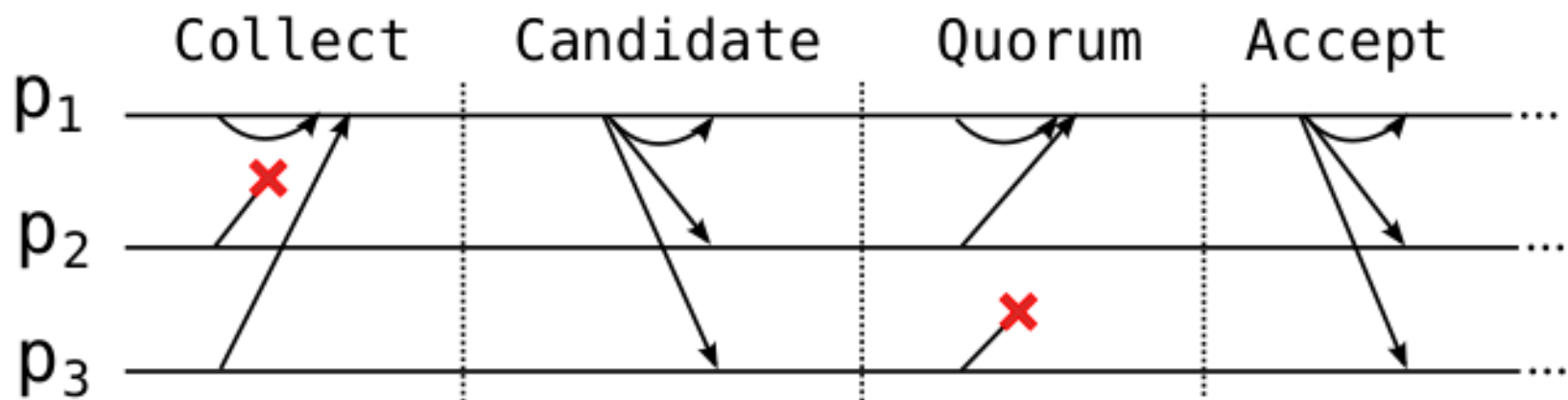
- a distributed system is a set of processes & an **adversarial environment**
- every round, the environment decides which messages processes receive
- each process p has a *Heard-Of set*, $\mathcal{HO}(p)$ = the set of processes p hears from
- in a given round, process p receives a message from process q if q sends a message to p and $q \in \mathcal{HO}(p)$
- \mathcal{HO} uniformly models *asynchronous* behaviour and *faults* while providing the illusion of a lockstep semantics

- A dropped message sent from q to p can be modelled as:
 - $q \notin HO(p)$
- A crashed process q can be modelled as:
 - $\forall p, q \notin HO(p)$
- Generally, assumptions regarding the HO sets can be given as linear temporal logic (LTL) formulas





(a) An asynchronous, faulty execution of the *LastVoting*



(b) Corresponding indistinguishable lockstep execution

LastVoting: Paxos adapted to \mathcal{HO} model

```
1 interface
2     init(v: Int); out(v: Int)
3
4 variable
5     x: Int; ts: Int; vote: Int
6     ready: Boolean; commit: Boolean
7     decided: Boolean; decision: Int
8
9 //auxiliary function: rotating coordinator
10 def coord(phi: Int): ProcessID =
11     new ProcessID((phi/phase.length) % n)
12
13 //initialization
14 def init(v: Int) =
15     x := v
16     ts := -1
17     ready := false
18     commit := false
19     decided := false
```

```
1 val phase = Array[Round]( //the rounds
2   Round /* Collect */ {
3     def send(): Map[ProcessID, (Int,Int)] =
4       return MapOf(coord(r) → (x, ts))
5     def update(mbox: Map[ProcessID, (Int,Int)]) =
6       if (id = coord(r) ∧ mbox.size > n/2)
7         vote := mbox.valWithMaxTS
8         commit := true },
9   Round /* Candidate */ {
10    def send(): Map[ProcessID, Int] =
11      if (id = coord(r) ∧ commit) return
12        broadcast(vote)
13      else return ∅
14    def update(mbox: Map[ProcessID, Int]) =
15      if (mbox contains coord(r))
16        x := mbox(coord(r))
17        ts := r/4 },
```

```
17 Round /* Quorum */ {
18     def send(): Map[ProcessID, Int] =
19         if ( ts = r/4 ) return MapOf(coord(r) → x)
20         else return ∅
21     def update(mbox: Map[ProcessID, Int]) =
22         if (id = coord(r) ∧ mbox.size > n/2)
23             ready := true },
24 Round /* Accept */ {
25     def send(): Map[ProcessID, Int] =
26         if (id = coord(r) ∧ ready) return broadcast(vote)
27         else return ∅
28     def update(mbox: Map[ProcessID, Int]) =
29         if (mbox contains coord(r) ∧ ¬decided)
30             decision := mbox(coord(r))
31             out(decision)
32             decided := true
33     ready := false
34     commit := false }
```


Properties for *LastVoting*

- Safety (agreement)

$$\Box(\forall p, p'. p.decided \wedge p'.decided \implies p.decision = p'.decision)$$

with the (simplified) invariant:

$$\begin{aligned} & \forall p. \quad p.decided = false \\ \vee & \quad \exists v, t, A. \quad A = \{p \mid p.ts \geq t\} \wedge |A| > n/2 \\ & \quad \wedge \forall p. p \in A \implies p.x = v. \end{aligned}$$

- Liveness

$$\Diamond(\forall p. p.decided)$$

```
21     val noDecision: Formula = P.forall( i => !i.decided && !i.ready)
22
23     val majority: Formula =
24         V.exists( v => V.exists( t => {
25             val A = P.filter( i => i.ts >= t )
26             A.size > n/2 &&
27             r > 0 &&
28             t <= r/4 &&
29             P.forall( i => (A.contains(i) ==> (i.x == v) ) &&
30                 (i.decided ==> (i.decision == v) ) &&
31                 (i.commit ==> (i.vote == v) ) &&
32                 (i.ready ==> (i.vote == v) ) &&
33                 ((i.ts == r/4) ==> coord.commit ) )
34         } ) )
35
36     val keepInit: Formula = P.forall( i => P.exists( j1 => i.x == init(j1.x) ) )
37
38     val safetyInv = And(keepInit, Or(noDecision, majority))
```

Semantics, runtime & verification

INIT

$$\frac{\forall p \in P. * \xrightarrow{\text{init}(v_p)} s(p)}{* \xrightarrow{\emptyset, \{\text{init}_p(v_p) \mid p \in P\}} \langle \text{Snd}, s, 0, \emptyset, \text{HO} \rangle}$$

SEND

$$\frac{\begin{array}{l} \forall p \in P. s(p) \xrightarrow{\text{phase}[r].\text{send}(m_p)} s(p) \\ \text{msg} = \{(p, t, q) \mid p \in P \wedge (t, q) \in m_p\} \end{array}}{\langle \text{Snd}, s, r, \emptyset, \text{HO} \rangle \xrightarrow{\{\text{send}_p(m_p) \mid p \in P\}, \emptyset} \langle \text{Updt}, s, r, \text{msg}, \text{HO}' \rangle}$$

UPDATE

$$\frac{\begin{array}{l} \forall p \in P. \text{mbox}_p = \{(q, t) \mid (q, t, p) \in \text{msg} \wedge q \in \text{HO}(p)\} \\ \forall p \in P. s(p) \xrightarrow{\text{phase}[r].\text{update}(\text{mbox}_p), o_p} s'(p) \\ r' = r + 1 \quad O = \{o_p \mid p \in p\} \end{array}}{\langle \text{Updt}, s, r, \text{msg}, \text{HO} \rangle \xrightarrow{\{\text{update}_p(\text{mbox}_p) \mid p \in p\}, O} \langle \text{Snd}, s', r', \emptyset, \text{HO} \rangle}$$

Runtime

- The lockstep semantics (LS) can be translated into asynchronous semantics (AS)
- Rounds are implemented using a **timeout**
 - *Paper describes how an appropriate timeout should be determined*
- If the specification is closed under indistinguishability, properties that hold in LS also hold in AS
 - *Consensus, k-set agreement and lattice agreement are closed under indist.*

Verification

- Can verify both *safety* and *liveness*
- Semi-automated, i.e. programmer manually writes invariants that are auto-checked by an SMT solver (**deductive verification**)
- Could potentially support **model-checking**, as lockstep semantics is sufficiently restricted

Conclusion

Take-aways

- *PSYNC strikes a balance between high-level constructs, performance, and automated verification*
- *The \mathcal{HO} model along with communication-closed rounds and lockstep execution, greatly simplifies implementing & verifying dist. algos*
- *For an important class of specifications including consensus, if a PSYNC program satisfies the specification, then its runtime system satisfies it as well*