

PSync: A Partially Synchronous Language for Fault- Tolerant Distributed Algorithms

Authors: Cezara
Dragoi, Thomas
A. Henzinger,
Damien Zufferey

Presented by:
Abhishek A. Singh

PSYNC

1

Domain Specific
Language

2

Based on Heard-Of
Model

- Which views asynchronous faulty systems as synchronous systems with an adversarial component

3

Implemented as an
embedding in Scala

4

Created for
modelling,
programming &
verification of
distributed fault-
tolerant algorithms

Heard-of Model

- Algorithms structured in communication closed rounds.
- Each round consists of two consecutive operations: Send and Update
- Two components in the model:
 - Set of Processes
 - Adversarial environment: which determines whether messages are received or not
- Each process has a Heard-of set: HO
- Each round is communication closed: all messages sent in a round are either received or dropped

The Heard-of set (HO)

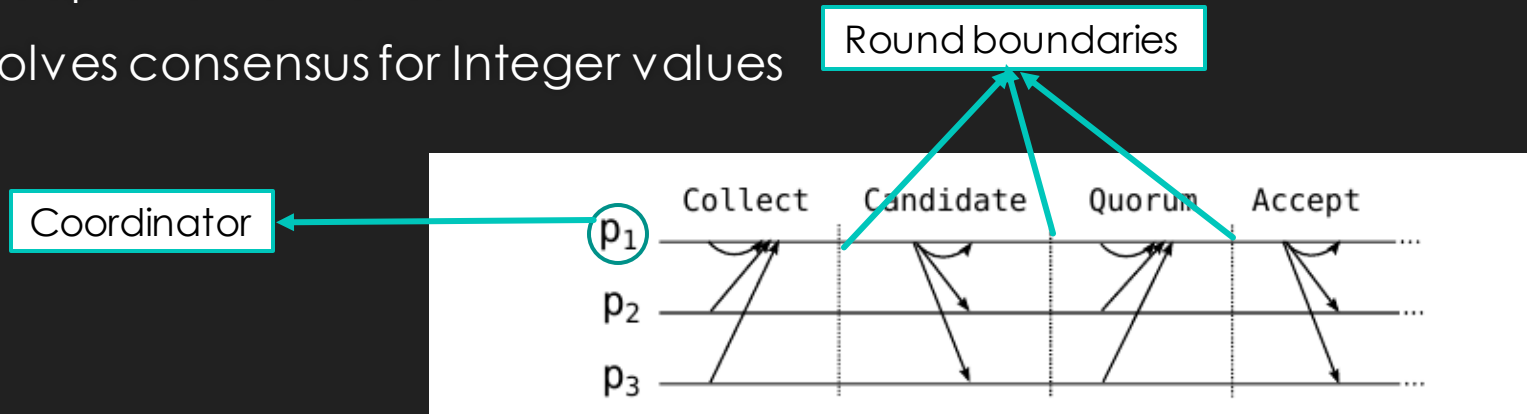
- Abstraction for the asynchronous and faulty behavior of the network.
- HO is a set of processes.
- In a round:
 - Process p receives a message from q if q sent a message to process p and $q \in \text{HO}(p)$
- HO+Rounds \Rightarrow Abstract notion of time + control structure for programmers

Runtime model

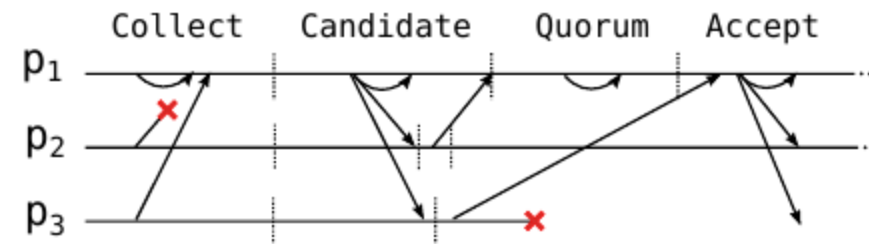
- Based on timeouts
- Send, Wait, Update, Move to next round, Repeat

Example: LastVoting

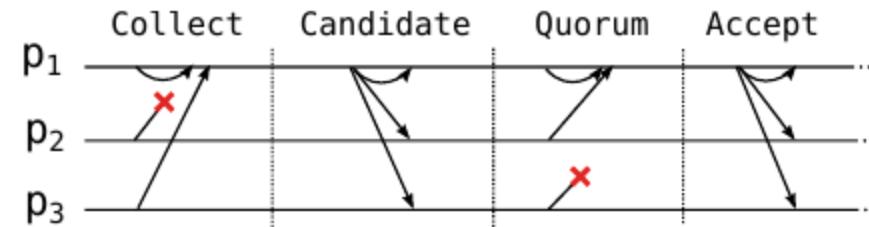
- Adaptation of Paxos
- Solves consensus for Integer values



LastVoting: Execution



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



(b) Corresponding indistinguishable lockstep execution

LastVoting: in PSync

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


PSync Syntax

```
program ::= interface variable* init phase  
interface ::= init: type  $\rightarrow ()$  (name: type  $\rightarrow ()$ )*  
variable ::= name: type  
    init ::= init: type  $\rightarrow ()$   
    phase ::= round+  
roundT ::= send: ()  $\rightarrow [P \mapsto T]$  update: [P  $\mapsto T]$   $\rightarrow ()$ 
```

Lockstep execution

Definition 5 (Lockstep execution). *Given a PSYNC program \mathcal{P} and a non-empty set of processes P , a lockstep execution of \mathcal{P} is the sequence $*A_0s_1A_1s_2 \dots$ such that*

- *$*A_0s_1$ is the result of the INIT rule;*
- *$\forall i. s_iA_isi+1$ satisfy the SEND or the UPDATE rule;*
- *the environment assumptions on HO-sets are satisfied.*

The set of lockstep executions of \mathcal{P} is denoted by $\llbracket \mathcal{P} \rrbracket_{ls}$.

Indistinguishability

- Defined in terms of a transition system.
- Transition system is intended to reflect an instance of execution

Definition 1 (Indistinguishability). *Given two executions π and π' of a transition system TS , a process p cannot distinguish locally between π and π' , denoted $\pi \simeq_p \pi'$, iff the projection of both executions on p agree up to finite stuttering, i.e., $\pi|_p \equiv \pi'|_p$.*

Two executions π and π' are indistinguishable, denoted $\pi \simeq \pi'$, iff no process can distinguish between them, i.e., $\forall p \in P. \pi \simeq_p \pi'$.

Definition 2 (Indistinguishable systems). *A system TS_1 is indistinguishable from a system TS_2 denoted $TS_1 \supseteq TS_2$ iff they are defined over the same set of processes and for any execution $\pi \in \llbracket TS_1 \rrbracket$ there exists an execution $\pi' \in \llbracket TS_2 \rrbracket$ such that $\pi \simeq_{W,L} \pi'$ where $W = V_1 \cap V_2$ and $L = A_1 \cap A_2$.*

Distributed clients

Definition 3 (Distributed client). ¹ Let $TS_i = (\{p_i\}, V_i, A_i, s_0^i, T_i)$ be the transition system associated with a client process, with $A_i \cap A_j = \emptyset$ for all $1 \leq i \neq j \leq n$. Formally, the transitions system associated with the client is $C_{TS} = (P, V, A, s_0, T)$, where $P = \{p_1, p_2, \dots, p_n\}$, $V = \biguplus_i V_i$, $A = \bigcup_i A_i$, $s_0 = (s_0^1, \dots, s_0^n)$, and $T \subseteq \Sigma \times A \times \Sigma$, with $\Sigma = [P \rightarrow V \rightarrow \mathcal{D}]$, such that $\Sigma(p_i) \in [V_i \rightarrow \mathcal{D}]$, and $(s, B, s') \in T$ iff for every $b \in B \cap A_i$ $(s(p_i), b, s'(p_i)) \in T_i$ and each processes takes at most one transition.

- All distributed clients as commutative by definition

Observational refinement

Definition 4 (Observational Refinement). *Let TS_1 and TS_2 be two transition systems and a common interface I . Then, TS_1 refines TS_2 w.r.t. I denoted $TS_1 \sqsubseteq_I TS_2$, if for any client C ,*

$$\text{Runs}(C(TS_1)) \subseteq \text{Runs}(C(TS_2)).$$

- TS_1 observationally refines TS_2 if every run of a client that uses TS_1 is also a run of the same client using TS_2 .

Theorem 1. *Let TS_1 and TS_2 be two systems with a common interface I . If $TS_1 \supseteq TS_2$ then $TS_1 \sqsubseteq_I TS_2$.*

- Indistinguishability is equivalent with sequential consistency



Comments or Questions