

Describing Asynchronous Communications in Distributed Algorithms

FRIDA 2017

Philippe QUÉINNEC

(joint work with Florent Chevrou, Aurélie Hurault, Adam Shimi)

Université de Toulouse – IRIT ACADIE – Pardi ANR-16-CE25-0006



Université
de Toulouse



October 16, 2017

Description of Interactions

Description of communication in a distributed algorithm

- Explicit communication actions or shared variables?
- Messages in transit ?
- Set/multiset of messages in transit, input queues?
- Embedded in the chosen formalism?
- ...

→ Common knowledge or folklore ?

Unexpected properties hidden in the communication models.

- 1 Several Models of Interaction
 - Asynchronous Consensus
 - Several Models
 - Comparison
- 2 Modelisation of the Communication
 - Messages in transit
 - Order of delivery
 - Fairness
 - Atomicity
 - Latency
- 3 A Framework for the Verification of Asynchronous Communicating Systems
 - Framework Architecture
 - Communication Models
 - Correctness & Completeness
 - Conclusion

- 1 Several Models of Interaction
 - Asynchronous Consensus
 - Several Models
 - Comparison
- 2 Modelisation of the Communication
- 3 A Framework for the Verification of Asynchronous Communicating Systems

Consensus

Consensus

A set of processes each propose a value. All processes must agree on a single one.

Properties

Agreement The correct processes all decide the same value.

Integrity A process decides at most once.

Validity The decision value is one of the proposed values.

Termination A correct process eventually decides.

An algorithm in an asynchronous world

Each process compute the minimum of the proposed values: it broadcasts its value, waits for all the values and chooses the min.

Questions: broadcast? waits? messages or values?

VARIABLES *proposed*, for each process, its proposed value
hasdecided, the set of processes that have decided on a value
decision for each process, the decided value

TypeInvariant \triangleq \wedge *proposed* \in [*PROC* \rightarrow *VALUE*]
 \wedge *hasdecided* \in SUBSET *PROC*
 \wedge *decision* \in [*PROC* \rightarrow *VALUE*]

The correct processes all decide the same value.

Agreement \triangleq $\square(\forall i, j \in \textit{hasdecided} : \textit{decision}[i] = \textit{decision}[j])$

Integrity \triangleq A process decides at most once.

$\forall i \in \textit{PROC}, k \in \textit{VALUE} : \square(i \in \textit{hasdecided} \wedge (\textit{decision}[i] = k) \Rightarrow \square(\textit{decision}[i] = k))$

Validity \triangleq The decision value is one of the proposed values.

$\square(\forall i \in \textit{hasdecided} : \exists j \in \textit{PROC} : \textit{decision}[i] = \textit{proposed}[j])$

A correct process eventually decides.

Termination \triangleq $\forall i \in \textit{PROC} : \diamond(i \in \textit{hasdecided})$

Not a Distributed Algorithm

VARIABLES *proposed*, for each process, its proposed value
hasdecided, the set of processes that have decided on a value
decision for each process, the decided value

Init \triangleq \wedge *hasdecided* = {}
 \wedge *proposed* \in [*PROC* \rightarrow *VALUE*]
 \wedge *decision* = *proposed*

decide(*i*) \triangleq
 \wedge *i* \notin *hasdecided*
 \wedge *hasdecided'* = *hasdecided* \cup {*i*}
 \wedge *decision'* = [*decision* EXCEPT
!*i* = $\min(\{\textit{proposed}[j] : j \in \textit{PROC}\})$]
 \wedge UNCHANGED \langle *proposed* \rangle

Next \triangleq \exists *i* \in *PROC* : *decide*(*i*)

What is proved ?

Simple Model

net = bag of values, atomic broadcast, atomic reception

$broadcast(i) \triangleq$

$\wedge i \notin hasbroadcasted$

$\wedge hasbroadcasted' = hasbroadcasted \cup \{i\}$

$\wedge net' = net \oplus SetToBag(\{proposed[i]\})$

$\wedge UNCHANGED \langle proposed, hasdecided, decision \rangle$

$decide(i) \triangleq$

$\wedge i \notin hasdecided$

$\wedge \exists M \in SubBag(net) :$

$\wedge BagCardinality(M) = N$

$\wedge hasdecided' = hasdecided \cup \{i\}$

$\wedge decision' = [decision \text{ EXCEPT } ![i] = min(M)]$

$\wedge UNCHANGED \langle proposed, hasbroadcasted, net \rangle$

$Next \triangleq \exists i \in PROC : broadcast(i) \vee decide(i)$

Messages?

net = set of messages, atomic broadcast, atomic reception

$broadcast(i) \triangleq$

$\wedge i \notin hasbroadcasted$

$\wedge hasbroadcasted' = hasbroadcasted \cup \{i\}$

$\wedge net' = net \cup \{[src \mapsto i, dest \mapsto j, value \mapsto proposed[i]] : j \in PROC\}$

$\wedge UNCHANGED \langle proposed, hasdecided, decision \rangle$

$decide(i) \triangleq$

$\wedge i \notin hasdecided$

$\wedge \exists M \in SUBSET \ net :$

$\wedge Cardinality(M) = N$

$\wedge \forall m \in M : m.dest = i$

$\wedge hasdecided' = hasdecided \cup \{i\}$

$\wedge decision' = [decision \text{ EXCEPT } ![i] = \min(\{m.value : m \in M\})]$

$\wedge net' = net \setminus M$

$\wedge UNCHANGED \langle proposed, hasbroadcasted \rangle$

$Next \triangleq \exists i \in PROC : broadcast(i) \vee decide(i)$

Input Queues?

Input queues, atomic broadcast, atomic reception

$broadcast(i) \triangleq$

$\wedge i \notin hasbroadcasted$

$\wedge hasbroadcasted' = hasbroadcasted \cup \{i\}$

$\wedge net' = [j \in PROC \mapsto Append(net[j], proposed[i])]$

$\wedge UNCHANGED \langle proposed, hasdecided, decision \rangle$

$decide(i) \triangleq$

$\wedge i \notin hasdecided$

$\wedge LET M \triangleq net[i] IN$

$\wedge Len(M) = N$

$\wedge hasdecided' = hasdecided \cup \{i\}$

$\wedge decision' = [decision EXCEPT ![i] = min(M)]$

$\wedge net' = [net EXCEPT ![i] = \langle \rangle]$

$\wedge UNCHANGED \langle proposed, hasbroadcasted \rangle$

$Next \triangleq \exists i \in PROC : broadcast(i) \vee decide(i)$

Non atomic reception

A set of messages, atomic broadcast, non-atomic reception

$broadcast(i) \triangleq$

$\wedge i \notin hasbroadcasted$

$\wedge hasbroadcasted' = hasbroadcasted \cup \{i\}$

$\wedge net' = net \cup \{[src \mapsto i, dest \mapsto j, value \mapsto proposed[i]] : j \in PROC\}$

$\wedge UNCHANGED \langle proposed, hasdecided, decision, received \rangle$

$receive(i) \triangleq$

$\exists m \in net :$

$\wedge m.dest = i$

$\wedge received' = [received \text{ EXCEPT } ![i] = @ \oplus SetToBag(\{m.value\})]$

$\wedge net' = net \setminus \{m\}$

$\wedge UNCHANGED \langle proposed, hasdecided, decision, hasbroadcasted \rangle$

$decide(i) \triangleq$

$\wedge i \notin hasdecided$

$\wedge BagCardinality(received[i]) = N$

$\wedge hasdecided' = hasdecided \cup \{i\}$

$\wedge decision' = [decision \text{ EXCEPT } ![i] = \min(received[i])]$

$\wedge UNCHANGED \langle proposed, hasbroadcasted, received, net \rangle$

Non atomic broadcast and reception

A set of messages, non-atomic broadcast, non-atomic reception

$send(i) \triangleq$

$\exists j \in PROC :$

$\wedge j \notin sent[i]$

$\wedge sent' = [sent \text{ EXCEPT } ![i] = @ \cup \{j\}]$

$\wedge net' = net \cup \{[src \mapsto i, dest \mapsto j, value \mapsto proposed[i]]\}$

$\wedge \text{UNCHANGED } \langle proposed, hasdecided, decision, received \rangle$

$receive(i) \triangleq$

$\exists m \in net :$

$\wedge m.dest = i$

$\wedge received' = [received \text{ EXCEPT } ![i] = @ \oplus SetToBag(\{m.value\})]$

$\wedge net' = net \setminus \{m\}$

$\wedge \text{UNCHANGED } \langle proposed, hasdecided, decision, sent \rangle$

$decide(i) \triangleq \dots$

$Next \triangleq \exists i \in PROC : send(i) \vee receive(i) \vee decide(i)$

Comparison

net	broadcast	reception	number of distinct states	
			3 sites, 2 values	4 / 2
direct access	–	–	64	256
bag of values	atomic	atomic	120	496
set of messages	atomic	atomic	120	496
input queues	atomic	atomic	216	1450
set of messages	atomic	non atomic	7568	1624096
synchronous-like	non atomic	non atomic	23814	8214536
set of messages	non atomic	non atomic	175616	?

Unexpected properties:

- input queues \Rightarrow delivery order stricter than causal delivery (so-called fifo n-1)
- atomic broadcast with input queues \Rightarrow total order broadcast!
- atomic multiple reception = transaction

- 1 Several Models of Interaction
- 2 **Modelisation of the Communication**
 - Messages in transit
 - Order of delivery
 - Fairness
 - Atomicity
 - Latency
- 3 A Framework for the Verification of Asynchronous Communicating Systems

Messages in transit

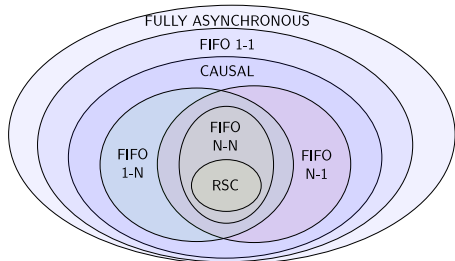
How can we model the messages in transit?

- Shared variables: no explicit message
- A set or multiset of messages (asynchronous communication)
- A sequence (a global queue)
- An array of sets (unordered input ports)
- An array of sequences (ordered input ports)
- A bi-dimensional array of sequences (fifo channels between couple of sites)
- Buffers
- ...

Is it important?

A data structure imposes a delivery order \Rightarrow different communication models

- Set = asynchronous
- Sequence = Fifo N-N
- Input queues = Fifo N-1
- Channels = Fifo 1-1
- 1-bounded = RSC



- Lots of different delivery orders: fifo11, causal, etc
- Hidden in the data structures
- Almost all are stricter than causal delivery (i.e. they generate less behaviors)

On the Diversity of Asynchronous Communication. F. Chevou, A. Hurault, and P. Quéinnec. *Formal Aspects of Computing*, 28(5):2016.

Fairness

- On each message \Rightarrow no message can stay forever in transit.
- Per receptor \Rightarrow no receptor can refuse to receive some message.
(however, by over-flooding it, one can lose message progress if no delivery order exists)
- Per sender \Rightarrow a site cannot be permanently ignored?
- Per couple (sender,receptor) \Rightarrow communication occurs between all sites.

Atomicity of communication actions

- Receive-and-reply

$$net' = net \setminus \{m1\} \cup \{m2\}$$

Analogous to a test-and-set or compare-and-swap

Assume we have *zero* which tests if there is a message in transit

<i>P1</i>	<i>P2</i>	<i>P3</i>	\perp reachable?
$a! \cdot b!$	$a? \cdot c!$	$b? \cdot (c? + zero \cdot \perp)$	yes
$a! \cdot b!$	$[a?c!]$	$b? \cdot (c? + zero \cdot \perp)$	no

- Multiple receptions

$$net' = net \setminus \{m_1, m_2, \dots\}$$

Looks like a transaction

$a! \cdot b! \parallel \{a, b\}? \parallel \{a, c\}?$ is different from $a! \cdot b! \parallel a? \cdot b? \parallel a? \cdot c?$
(in the second case, the third process can get *a* and deadlock)

Atomicity of communication actions

- Multiple emissions: $net' = net \cup \{m \mid \dots\}$

Chandy-Misra algorithm for mutual exclusion

A requesting site must send a request message to all the sites to which it has previously given its permission.

Pitfall: a site must not handle a request message while it is itself sending request messages (or it must be careful).

- Broadcast = multiple emissions?

Atomicity of multiple emissions + queues (be it per site or per channel) = total order broadcast : if two broadcasts concurrently occur, the messages to the common destinations are in the same order everywhere.

Latency

- Is a message immediately available, in the transition which follows the emission? Or do we need buffers to delay? It should be irrelevant (interleaving, non-determinism).
- Application priority on channels:

$c! \cdot a! \cdot b! \parallel b? \cdot (a? + c?)$ with $a > c$.

The a branch is mandatory without explicit latency.

- Assume Test-and-Action: $zero_a!$ sends a message on a if there is no pending message for the peer.

	both can die?
$(\tau \cdot kill_2! + kill_1?) \parallel (\tau \cdot kill_1! + kill_2?)$	yes
$(zero_kill_2! + kill_1?) \parallel (zero_kill_1! + kill_2?)$ without buffers	no
$(zero_kill_2! + kill_1?) \parallel (zero_kill_1! + kill_2?)$ with buffers	yes

Emptiness Test / Capacity / Boundness

Emptiness test

- Is it allowed to test $net = \emptyset$?
- Is it allowed to test if there is a pending message for a process?

→ Carefully (see previous examples in atomicity and latency)

Global predicates

- Example: maximal capacity (maximal number of messages in transit)
- Are global predicates meaningful? Enforceable?

- 1 Several Models of Interaction
- 2 Modelisation of the Communication
- 3 A Framework for the Verification of Asynchronous Communicating Systems
 - Framework Architecture
 - Communication Models
 - Correctness & Completeness
 - Conclusion

Verification of Asynchronous Communicating Systems

Modeling

- *Asynchronously* communicating peers
- Multi-senders multi-receivers channels
- Group of channels associated to different communication models

Method

- Peers and communication specified using transition systems
- Compatibility checking of LTL properties
- Modular (new models)
- Fully automatic: model checking
- TLA⁺

Automated Verification of Asynchronous Communicating Systems with TLA⁺.
F. Chevrou, A. Hurault, and P. Quéinnec. Electronic Communications of the
EASST (special issue: AVOCS'15), 2015.

Explicit hypotheses

- Several communication models with specified delivery order
- Weak fairness on communication actions
- One communication action per transition (no multiple emissions, no receive-and-reply)
- ... but point-to-point and broadcast communication
- No latency, but no test-and-action

Communication models

A communication model is built by assembling μ -models

Two point-to-point channels, with fifo11 ordering, and at most two messages in transit on channel Request.

```
MODELS == {  
  [name |-> "p2p",  
    params |-> [chan |-> {"Request", "Token"} ]],  
  [name |-> "fifo11",  
    params |-> [chan |-> {"Request", "Token"} ]],  
  [name |-> "message_cap",  
    params |-> [chan |-> {"Request"}, bound |-> 2]]  
}
```

A communication action (send/receive) on a channel is enabled if it is enabled in all concerned models.

Sending & Receiving

```
Request(i) ==  
  ∧ ~ requesting[i]  
  ∧ requesting' = [ requesting EXCEPT ![i] = TRUE ]  
  ∧ COM!send(i, {father[i]}, "Request", i)  
    (* sender, destination, channel, content *)  
  ∧ ...
```

```
ReceiveToken(i) ==  
  ∧ ∃ j ∈ Sites : COM!receive(j, i, "Token", 0)  
    (* sender, receiver, channel, content *)  
  ∧ token' = [ token EXCEPT ![i] = TRUE ]  
  ∧ ...
```

Correctness & Completeness

Valid Execution

An execution respects e.g. the causal delivery order iff

$$\forall m_1, m_2 : \text{peer}(d(m_1)) = \text{peer}(d(m_2)) \wedge e(m_1) \prec e(m_2) \\ \Rightarrow d(m_1) \prec d(m_2)$$

where $e(m)$ = emission event of m

$r(m)$ = delivery event of m

\prec = causal order of events

Correctness

Any execution generated by the framework is valid.

Mechanized proofs, by refinement.

Completeness

All the valid executions are generated by the framework.

(with restrictions on the form of the system)

Manual proofs.

Advantages/Limitations

Advantages

- No re-re-re-coding of the communication actions.
- Easy reconfiguration, easy tests.
- No hidden properties.
- Experimenting platform available online (<http://vacs.enseeiht.fr>).

Limitations

- Generic action vs ad-hoc implementation:
 $\exists i, j \in Site : \exists d \in Data : COM!receive(i, j, "chan", d) \wedge \dots$
 $LET m \triangleq Head(chan[i][j]) IN \dots$
→ same number of states/transitions but slower
- Broadcast is still experimental
- Failures are missing (should they be in a communication model?)

Conclusion

In a communication model, beware of

- Hidden properties
- Incomplete model
- Model too close to an implementation. . .
- . . . model too abstract / unclear w.r.t. reality

Need for

- Generic models (of computation, of communication),
- with modularity and parameterization,
- with known properties (and no hidden ones),
- with mechanization.