

Analyzing Paxos with Fault-Tolerant Multiparty Session Types

Bachelor thesis by Nicolas Daniel Torres
Date of submission: January 15, 2022

1. Review: Prof. Dr. Kirstin Peters
2. Review: M.Sc. Anna Schmitt
Darmstadt



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Computer Science
Department

<institute>

<working group>

Contents

1	Introduction	3
2	Technical Preliminaries	4
2.1	Paxos	4
2.2	Multiparty Session Types	4
2.3	Fault-Tolerant Multiparty Session Types	4
2.3.1	Fault-Tolerance in Distributed Algorithms	4
2.3.2	Fault-Tolerant Types and Processes	6
2.3.3	A Semantics with Failure Patterns	11
2.3.4	Typing Fault-Tolerant Processes	14
2.4	Additional Notation	19
3	Model	20
3.1	Sorts	20
3.2	Global Type	21
3.3	Functions	21
3.4	Processes	23
3.4.1	System Initialization	23
3.4.2	Proposer	24
3.4.3	Acceptor	25
3.5	Failure Patterns	26
3.6	Example	26
3.6.1	Scenario	28
3.6.2	Formulae	28
4	Analysis	34
4.1	Local Types	34
4.2	Type Check	35
4.2.1	System Initialization	35
4.2.2	Proposer	36
4.2.3	Acceptor	38
4.3	Termination, Agreement, Validity	40
4.3.1	Termination	40
4.3.2	Agreement	41
4.3.3	Validity	41

1 Introduction

In distributed systems components on different computers coordinate and communicate via message passing to achieve a common goal. Sometimes, to achieve this goal, the individual components need to reach consensus, i.e., agree on the value of some data using a consensus algorithm. For example in state machine replication or when deciding which database transactions should be committed in what order. For such a distributed system to behave correctly the consensus algorithm needs to be correct. Thus, analyzing consensus algorithms is important.

To achieve consensus, consensus algorithms must satisfy the following properties: termination, validity, and agreement [7]. Proving these properties can be complicated. Model checking tools lead to big state-spaces so static analysis is preferable. For static analysis Multiparty Session Types are particularly interesting because session typing can ensure protocol conformance and the absence of communication errors and deadlocks [14].

Due to the presence of faulty processes and unreliable communication consensus algorithms are designed to be fault-tolerant. Modelling fault-tolerance is not possible using Multiparty Session Types, thus a fault-tolerant extension is necessary. Peters, Nestmann, and Wagner developed such an extension called Fault-Tolerant Multiparty Session Types.

In this work we will use Fault-Tolerant Multiparty Session Types to analyze the consensus algorithm Paxos, as described in [11].

2 Technical Preliminaries

In this chapter we will introduce the Paxos consensus algorithm, Multiparty Session Types (MPST), Fault-Tolerant Multiparty Session Types (FTMPST), and some notation.

2.1 Paxos

2.2 Multiparty Session Types

Multiparty Session Types (MPST) are used to statically ensure correctly coordinated behavior in systems without global control ([9, 6]). One important such property is progress, i.e., the absence of deadlock. Like with every other static typing approach, the main advantage is their efficiency, i.e., they avoid the problem of state space explosion.

MPST are designed to abstractly capture the structure of communication protocols. They describe global behaviors as *sessions*, i.e., units of conversations [9, 2, 3]. The participants of such sessions are called *roles*. *Global types* specify protocols from a global point of view. These types are used to reason about processes formulated in a *session calculus*. Most of the existing session calculi are variants of the π -calculus [12].

2.3 Fault-Tolerant Multiparty Session Types

Our model of the Paxos algorithm uses a fault-tolerant extension of Multiparty Session Types introduced by Peters, Nestmann, and Wagner in [13]. The following explanation of Fault-Tolerant Multiparty Session Types is from that same paper.

2.3.1 Fault-Tolerance in Distributed Algorithms

We consider three sources of failure in an unreliable communication (Fig. 2.1(a)): (1) the sender may crash before it releases the message, (2) the receiver may crash before it can consume the message, or (3) the communication medium may lose the message. The design of a distributed algorithm may allow it to handle some kinds of failures better than others. Failures are unpredictable events that occur at runtime. Since types consider only static and predictable information, we do not distinguish between different kinds of failure or model their source in types. Instead we only allow types, i.e., the specifications of systems, to distinguish between potentially faulty and reliable interactions.

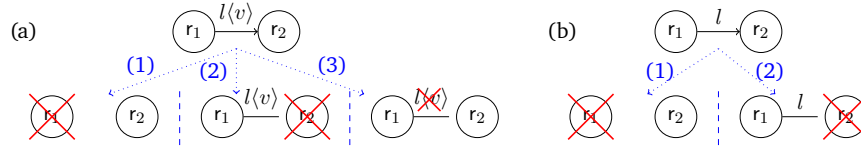


Figure 2.1: Unreliable Communication (a) and Weakly Reliable Branching (b).

A fault-tolerant algorithm has to solve its task despite such failures. Remember that MPST analyse the communication structure. Accordingly, we need a mechanism to tolerate faults in the communication structure. We want our type system to ensure that a faulty interaction neither blocks the overall protocol nor influences the communication structure of the system after this fault. We consider an unreliable communication as fault-tolerant if a failure does not influence the guarantees that our type system provides for the overall communication structure except for this particular communication. Moreover, if a potentially unreliable communication is executed successfully, then our type system ensures the same guarantees as for reliable communication such as e.g. the absence of communication mismatches.

To ensure that a failure does not block the algorithm, both the receiver and the sender need to be allowed to proceed without their unreliable communication partner. Therefore, the receiver of an unreliable communication is required to specify a default value that, in the case of failure, is used instead of the value the process was supposed to receive. The type system ensures the existence of such default values and checks their sort. Moreover, we augment unreliable communication with labels that help us to avoid communication mismatches. This is enough to ensure that the communication structure of a distributed algorithm is fault-tolerant.

Branching in the context of failures is more difficult, because a branch marks a decision point in a specification, i.e., the participants of the session are supposed to behave differently w.r.t. this decision. In an unreliable setting it is difficult to ensure that all participants are informed consistently about such a decision and adapt their behaviour accordingly.

Consider a reliable branching that is decided by a process r_1 and transmitted to r_2 . If we try to execute such a branching despite failures, we observe that there are again three ways in that this branching can go wrong (Fig. 2.1(b)): (1) The sender may crash before it releases its decision. This will block r_2 , because it is missing the information about the branch it should move to. (2) The receiver might crash. (3) The message of r_1 about the decided branch is lost. Then again r_2 is blocked.

Case (2) can be dealt with similar to unreliable communication, i.e., by marking the branching as potentially faulty and by ensuring that a crash of r_2 will not block another process. For Case (1) we declare one of the offered branches as default to that r_2 moves if r_1 has crashed. Then r_2 will not necessarily move to the branch that r_1 had in mind before it crashed, but to a valid/specified branch and, since r_1 is crashed, no two processes move to different branches. The main problem is in Case (3). Let r_1 move to a non-default branch and transmit its decision to r_2 , this message gets lost, and r_2 moves to the default branch. Now both processes did move to branches that are described by their types; but they are in different branches. Accordingly, this case violates the specification in the type and we want to reject it. More precisely, we consider three levels of failures in interactions:

Strongly Reliable (r): Neither the sender nor the receiver can crash as long as they are involved in this interaction. The message cannot be lost by the communication medium. This form corresponds to reliable communication as it was described in [1] in the context of distributed algorithms. This is the standard, failure-free case.

Global Types	Local Types	Processes
$G ::= r_1 \rightarrow_r r_2 : \langle S \rangle . G$ $ \quad r_1 \rightarrow_u r_2 : l \langle S \rangle . G$ $ \quad r_1 \rightarrow_r r_2 : \{l_i . G_i\}_{i \in I}$ $ \quad r \rightarrow_w R : \{l_i . G_i\}_{i \in I, l_d}$ $ \quad G_1 \parallel G_2$ $ \quad (\mu t)G \quad \quad t \quad \quad \text{end}$ $ \quad r_1 \rightarrow r_2 : \langle s'[r] : T \rangle . G$	$T ::= [r_2]!_r \langle S \rangle . T$ $ \quad [r_1]?_r \langle S \rangle . T$ $ \quad [r_2]!_u l \langle S \rangle . T$ $ \quad [r_1]?_u l \langle S \rangle . T$ $ \quad [r_2]!_r \{l_i . T_i\}_{i \in I}$ $ \quad [r_1]?_r \{l_i . T_i\}_{i \in I}$ $ \quad [R]!_w \{l_i . T_i\}_{i \in I}$ $ \quad [r]?_w \{l_i . T_i\}_{i \in I, l_d}$ $ \quad (\mu t)T \quad \quad t \quad \quad \text{end}$ $ \quad [r_2]! \langle s'[r] : T \rangle . T'$ $ \quad [r_1]? \langle s'[r] : T \rangle . T'$	$P ::= \bar{a}[n](s).P$ $ \quad a[r](s).P$ $ \quad s[r_1, r_2]!_r \langle e \rangle . P$ $ \quad s[r_2, r_1]?_r(x).P$ $ \quad s[r_1, r_2]!_u l \langle e \rangle . P$ $ \quad s[r_2, r_1]?_u l \langle v \rangle (x).P$ $ \quad s[r_1, r_2]!_r l . P$ $ \quad s[r_2, r_1]?_r \{l_i . P_i\}_{i \in I}$ $ \quad s[r, R]!_w l . P$ $ \quad s[r_j, r]?_w \{l_i . P_i\}_{i \in I, l_d}$ $ \quad P_1 \mid P_2$ $ \quad (\mu X)P \quad \quad X \quad \quad \mathbf{0}$ $ \quad \text{if } b \text{ then } P_1 \text{ else } P_2$ $ \quad (\nu x)P \quad \quad \perp$ $ \quad s[r_1, r_2]! \langle s'[r] \rangle . P$ $ \quad s[r_2, r_1]? \langle s'[r] \rangle . P$ $ \quad s_{r_1 \rightarrow r_2} : M$
Message Types		Messages
$MT ::= \langle S \rangle^r \quad \quad l \langle S \rangle^u \quad \quad l^r \quad \quad l^w \quad \quad s[r]$		$M ::= \langle v \rangle^r \quad \quad l \langle v \rangle^u \quad \quad l^r$ $ \quad l^w \quad \quad s[r]$

Figure 2.2: Syntax of Fault-Tolerant MPST

Weakly Reliable (w): Both the sender and the receiver might crash at every possible point during this interaction. But the communication medium cannot lose the message.

Unreliable (u): Both the sender and the receiver might crash at every possible point during this interaction and the communication medium might lose the message. There are no guarantees that this interaction—or any part of it—takes place. In this case, it is difficult for the type system to ensure interesting properties in branching.

2.3.2 Fault-Tolerant Types and Processes

We assume that the sets \mathcal{N} of names a, s, x, \dots ; \mathcal{R} of roles n, r, \dots ; \mathcal{L} of labels l, l_d, \dots ; \mathcal{V}_T of type variables t ; and \mathcal{V}_P of process variables X are pairwise distinct. For clarity, we often distinguish names into *values*, i.e., the payload of messages, *shared channels*, or *session channels* according to their usage; there is, however, no need to formally distinguish between different kinds of names. To simplify the reduction semantics of our session calculus, we use natural numbers as roles (compare to [9]). Sorts S range over $\mathbb{B}, \mathbb{N}, \dots$. The set \mathcal{E} of expressions e, v, b, \dots is constructed from the standard Boolean operations, natural numbers, names, and (in)equalities.

Global types specify the desired communication structure of systems from a global point of view. In local types this global view is projected to the specification of a single role/participant. We use standard MPST ([8, 9]) extended by operators for unreliable communication and weakly reliable branching that are highlighted in blue colour in Fig. 2.2.

The processes $\bar{a}[n](s).P$ and $a[r](s).P$ initialise a new session s with n roles via the shared channel a and then proceed as P . We identify sessions with their unique session channel.

The type $r_1 \rightarrow_r r_2 : \langle S \rangle . G$ specifies a strongly reliable communication from role r_1 to role r_2 to transmit a value of the sort S and then continues with G . A system with this type will be guaranteed to perform a corresponding action. In a session s this communication is implemented by the sender $s[r_1, r_2]!_r \langle e \rangle . P_1$ (specified as $[r_2]!_r \langle S \rangle . T_1$) and the receiver $s[r_2, r_1]?_r \langle x \rangle . P_2$ (specified as $[r_1]?_r \langle S \rangle . T_2$). As result of the communication, the receiver instantiates x in its continuation P_2 with the received value.

The type $r_1 \rightarrow_u r_2 : l \langle S \rangle . G$ specifies an unreliable communication from r_1 to r_2 transmitting (if successful) a label l and a value of type S and then continues (regardless of the success of this communication) with G . The unreliable counterparts of senders and receivers are $s[r_1, r_2]!_u l \langle e \rangle . P_1$ (specified as $[r_2]!_u l \langle S \rangle . T_1$) and $s[r_2, r_1]?_u l \langle v \rangle \langle x \rangle . P_2$ (specified as $[r_1]?_u l \langle S \rangle . T_2$). The receiver $s[r_2, r_1]?_u l \langle v \rangle \langle x \rangle . P_2$ declares a default value v that is used instead of a received value to instantiate x after a failure. Moreover, a label is communicated that helps us to ensure that a faulty unreliable communication has no influence on later actions.

The strongly reliable branching $r_1 \rightarrow_r r_2 : \{l_i . G_i\}_{i \in I}$ allows r_1 to pick one of the branches offered by r_2 . We identify the branches with their respective label. Selection of a branch is implemented by $s[r_1, r_2]!_r l . P$ (specified as $[r_2]!_r \{l_i . T_i\}_{i \in I}$). Upon receiving branch l_j from r_1 the process $s[r_2, r_1]?_r \{l_i . P_i\}_{i \in I}$ (specified as $[r_1]?_r \{l_i . T_i\}_{i \in I}$) continues with P_j .

The weakly reliable counterpart of branching is $r \rightarrow_w R : \{l_i . G_i\}_{i \in I, l_d}$, where $R \subseteq \mathcal{R}$ and l_d with $d \in I$ is the default branch. We use a broadcast from r to all roles in R to ensure that the sender can influence several participants with its decision consistently as it is the case for strongly reliable branching. Note that splitting this action to inform the roles in R separately does not work, because strongly reliable branching does not allow participants to crash and subsequent weakly reliable branchings cannot ensure that all receivers get the message if the sender crashes while performing these subsequent actions. The type system will ensure that this branching construct is weakly reliable, i.e., the involved participants might crash but no message is lost. Because of that, all processes that are not crashed will move to the same branch. We often abbreviate branching w.r.t. to a small set of branches by omitting the set brackets and instead separating the branches by \oplus , where the last branch is always the default branch. In contrast to the strongly reliable cases, the weakly reliable selection $s[r, R]!_w l . P$ (specified as $[R]!_w \{l_i . T_i\}_{i \in I}$) allows to broadcast its decision to R and $s[r_j, r]?_w \{l_i . P_i\}_{i \in I, l_d}$ (specified as $[r]?_w \{l_i . T_i\}_{i \in I, l_d}$) defines a default label l_d .

The \perp denotes a process that crashed. Similar to [9], we use message queues to implement asynchrony in sessions. Therefore, session initialisation introduces a directed and initially empty message queue $s_{r_1 \rightarrow r_2} : []$ for each pair of roles $r_1 \neq r_2$ of the session s . The separate message queues ensure that messages with different sources or destinations are not ordered, but each message queue is FIFO. Since strongly reliable, weakly reliable, and unreliable forms of interaction might be implemented differently (e.g. by TCP or UDP), it make sense to further split the message queues into three message queues for each pair $r_1 \neq r_2$ such that different kinds of messages do not need to be ordered. To simplify the presentation of examples in this paper and not to blow up the number of message queues, we stick to a single message queue for each pair $r_1 \neq r_2$, but the correctness of our type system does not depend on this decision. We have five kinds of messages and corresponding message types in Fig. 2.2—one for each kind of interaction.

The remaining operators for independence $G \parallel G'$; parallel composition $P \mid P'$; recursion $(\mu t)G$, $(\mu X)P$; inaction end , $\mathbf{0}$; conditionals $\text{if } b \text{ then } P_1 \text{ else } P_2$; session delegation $r_1 \rightarrow r_2 : \langle s'[r] : T \rangle . G$, $s[r_1, r_2]! \langle \langle s'[r] \rangle \rangle . P$, $s[r_2, r_1]? \langle \langle s'[r] \rangle \rangle . P$; and restriction $(\nu x)P$ are all standard.

Consider the specification $G_{\text{dice},r}$ of a simple dice game in a bar

$$(\mu t)3 \rightarrow_r 1:\langle \mathbb{N} \rangle.3 \rightarrow_r 2:\langle \mathbb{N} \rangle.3 \rightarrow_r 1:\{roll.3 \rightarrow_r 2:roll.t, exit.3 \rightarrow_r 2:exit.end\} \quad (1)$$

where the dealer Role 3 continues to *roll* a dice and tell its value to player 1 and then to *roll* another time for player 2 until the dealer decides to *exit* the game.

We can combine strongly reliable communication/branching and unreliable communication, e.g. by ordering a drink before each round in $G_{\text{dice},r}$.

$$\begin{aligned} (\mu t)3 \rightarrow_u 4:drink\langle \mathbb{N} \rangle.3 \rightarrow_r 1:\langle \mathbb{N} \rangle.3 \rightarrow_r 2:\langle \mathbb{N} \rangle. \\ 3 \rightarrow_r 1:\{roll.3 \rightarrow_r 2:roll.t, \quad exit.3 \rightarrow_r 2:exit.end\} \end{aligned}$$

where role 4 represents the bar tender and the noise of the bar may swallow these orders. Moreover, we can remove the branching and specify a variant of the dice game in that 3 keeps on rolling the dice forever, but, e.g. due to a bar fight, one of our three players might get knocked out at some point or the noise of this fight might swallow the announcements of role 3:

$$G_{\text{dice},u} = (\mu t)3 \rightarrow_u 1:roll\langle \mathbb{N} \rangle.3 \rightarrow_u 2:roll\langle \mathbb{N} \rangle.t \quad (2)$$

To restore the branching despite the bar fight that causes failures, we need the weakly reliable branching mechanism.

$$\begin{aligned} G_{\text{dice}} = (\mu t)3 \rightarrow_w \{1, 2\}: play.3 \rightarrow_u 1:roll\langle \mathbb{N} \rangle.3 \rightarrow_u 2:roll\langle \mathbb{N} \rangle.t, \\ \oplus end.3 \rightarrow_u 1:win\langle \mathbb{B} \rangle.3 \rightarrow_u 2:win\langle \mathbb{B} \rangle.end \end{aligned} \quad (3)$$

If 3 is knocked out by the fight, i.e., crashes, the game cannot continue. Then 1 and 2 move to the default branch *end*, have to skip the respective unreliable communications, and terminate. But the game can continue as long as 3 and at least one of the players 1, 2 participate.

An implementation of G_{dice} is $P_{\text{dice}} = P_3 \mid P_1 \mid P_2$, where for $i \in \{1, 2\}$:

$$\begin{aligned} P_3 &= \bar{a}[3](s).(\mu X)\text{if } x_1 \leq 21 \wedge x_2 \leq 21 \\ &\quad \text{then } s[3, \{1, 2\}]!_w play.s[3, 1]!_u roll\langle roll(x_1) \rangle.s[3, 2]!_u roll\langle roll(x_2) \rangle.X \\ &\quad \text{else } s[3, \{1, 2\}]!_w end.s[3, 1]!_u win\langle x_1 \leq 21 \rangle.s[3, 2]!_u win\langle x_2 \leq 21 \rangle.\mathbf{0} \\ P_i &= a[i](s).(\mu X)s[i, 3]?_w play.s[i, 3]?_u roll\langle x \rangle(x).X \oplus end.s[i, 3]?_u win\langle f \rangle(w).\mathbf{0} \end{aligned}$$

Role 3 stores the sums of former dice rolls for the two players in its local variables x_1 and x_2 , and $roll(x_i)$ rolls a dice and adds its value to the respective x_i . Role 3 keeps rolling dice until the sum x_i for one of the players exceeds 21. If both sums x_1 and x_2 exceed 21 in the same round, then 3 wins, i.e., both players receive f ; else, the player that stayed below 21 wins and receives t . The players 1 and 2 use their respective last known sum that is stored in x as default value for the unreliable communication in the branch *play* and f as default value in the branch *end*. The last branch, i.e., *end*, is the default branch.

Our type system verifies processes, i.e., implementations, against a specification that is a global type. Since processes implement local views, local types are used as a mediator between the global specification and the respective local end points. To ensure that the local types correspond to the global type, they are derived by *projection*. Instead of the projection function described in [9] we use a more relaxed variant of projection as introduced in [17].

Projection maps global types onto the respective local type for a given role p . The projections of the new global types are obtained straightforwardly from the projection of their respective strongly reliable counterparts:

$$(r_1 \rightarrow_{\diamond} r_2 : \mathfrak{S}.G) \upharpoonright_p \triangleq \begin{cases} [r_2]!_{\diamond} \mathfrak{S}.G \upharpoonright_p & \text{if } p = r_1 \\ [r_1]?_{\diamond} \mathfrak{S}.G \upharpoonright_p & \text{if } p = r_2 \\ G \upharpoonright_p & \text{otherwise} \end{cases}$$

where either $\diamond = r$, $\mathfrak{S} = \langle S \rangle$ or $\diamond = u$, $\mathfrak{S} = l \langle S \rangle$ and

$$(r_1 \rightarrow_{\diamond} \mathfrak{R} : \{l_i.G_i\}_{i \in I \mathfrak{D}}) \upharpoonright_p \triangleq \begin{cases} [\mathfrak{R}]!_{\diamond} \{l_i.G_i \upharpoonright_p\}_{i \in I} & \text{if } p = r_1 \\ [r_1]?_{\diamond} \{l_i.G_i \upharpoonright_p\}_{i \in I \mathfrak{D}} & \text{if } \mathfrak{B} \\ \bigsqcup_{i \in I} (G_i \upharpoonright_p) & \text{otherwise} \end{cases}$$

where either $\diamond = r$, $\mathfrak{R} = r_2$, \mathfrak{B} is $p = r_2$, \mathfrak{D} is empty or $\diamond = w$, $\mathfrak{R} = R$, \mathfrak{B} is $p \in R$, \mathfrak{D} is l_d . In the last case of strongly reliable or weakly reliable branching—when projecting onto a role that does not participate in this branching—we map to $\bigsqcup_{i \in I} (G_i \upharpoonright_p) = (G_1 \upharpoonright_p) \sqcup \dots \sqcup (G_n \upharpoonright_p)$. The operation \sqcup is (similar to [17]) inductively defined as:

$$\begin{aligned} T \sqcup T &= T \\ ([r]?_r I_1) \sqcup ([r]?_r I_2) &= [r]?_r (I_1 \sqcup I_2) \\ ([r]?_w I_1) \sqcup ([r]?_w I_2) &= [r]?_w (I_1 \sqcup I_2) \quad \text{if } I_1 \text{ and } I_2 \text{ have the same default branch} \\ I \sqcup \emptyset &= I \\ I \sqcup (\{l.T\} \cup J) &= \begin{cases} \{l.(T' \sqcup T)\} \cup ((I \setminus \{l.T'\}) \sqcup J) & \text{if } l.T' \in I \\ \{l.T\} \cup (I \sqcup J) & \text{if } l \notin I \end{cases} \end{aligned}$$

where $T, T' \in \mathcal{T}$, $l \notin I$ is short hand for $\nexists T'$. $l.T' \in I$, and is undefined in all other cases. The mergeability relation \sqcup states that two types are identical up to their branching types, where only branches with distinct labels are allowed to be different. This ensures that if the sender r_1 in $r_1 \rightarrow_r r_2 : \{l_i.G_i\}_{i \in I}$ decides to branch then only processes that are informed about this decision can adapt their behaviour accordingly; else projection is **not** defined.

The remaining global types are projected as follows:

$$\begin{aligned} (G_1 \parallel G_2) \upharpoonright_p &\triangleq \begin{cases} G_1 \upharpoonright_p & \text{if } p \notin R(G_2) \\ G_2 \upharpoonright_p & \text{if } p \notin R(G_1) \end{cases} \quad ((\mu t)G) \upharpoonright_p \triangleq \begin{cases} (\mu t)G \upharpoonright_p & \text{if } p \in R(G) \\ \text{end} & \text{otherwise} \end{cases} \\ t \upharpoonright_p &\triangleq t \quad \text{end} \upharpoonright_p \triangleq \text{end} \end{aligned}$$

The projection of $G_1 \parallel G_2$ on p is **not** defined if p occurs on both sides of this parallel composition; it is $G_i \upharpoonright_p$ if p occurs in exactly one side $i \in \{1, 2\}$; or it is $(G_1 \parallel G_2) \upharpoonright_p = G_1 \upharpoonright_p = G_2 \upharpoonright_p = \text{end}$ if p does not occur at all. Projecting a recursive global type results in a recursive local type if p occurs in the body of the recursion or else in successful termination. Type variables and successful termination are mapped onto themselves. We denote a global type G as *projectable* if for all $r \in R(G)$ the projection $G \upharpoonright_r$ is defined. We restrict our attention to projectable global types.

Projecting the global type $G_{\text{dice},r}$ in (1) results in the local types

$$\begin{aligned} T_{3:\text{dice},r} &= (\mu t)[1]!_r \langle \mathbb{N} \rangle . [2]!_r \langle \mathbb{N} \rangle . [1]!_r \{ \text{roll}. [2]!_r \text{roll}. t, \quad \text{exit}. [2]!_r \text{exit}. \text{end} \} \\ T_{i:\text{dice},r} &= (\mu t)[3]?_r \langle \mathbb{N} \rangle . [3]?_r \{ \text{roll}. t, \quad \text{exit}. \text{end} \} \end{aligned}$$

where the types of the two players $T_{1:\text{dice},r} = T_{2:\text{dice},r} = T_{i:\text{dice},r}$ are identical. The projection of the outer branching in $G_{\text{dice},r}$ on 2 results in $[3]?_r \text{roll}.t$ for the first branch and $[3]?_r \text{exit}.end$ for the second branch. These two $[3]?_r$ types are unified by \sqcup into a single $[3]?_r$ type with two branches.

Projection maps G_{dice} in (3) to:

$$\begin{aligned} T_{3:\text{dice}} &= (\mu t)[\{1, 2\}]!_w \text{ play}.[1]!_u \text{roll}\langle\mathbb{N}\rangle.[2]!_u \text{roll}\langle\mathbb{N}\rangle.t \\ &\quad \oplus \text{end}.[1]!_u \text{win}\langle\mathbb{B}\rangle.[2]!_u \text{win}\langle\mathbb{B}\rangle.\text{end} \\ T_{i:\text{dice}} &= (\mu t)[3]?_w(\text{play}.[3]?_u \text{roll}\langle\mathbb{N}\rangle.t \oplus \text{end}.[3]?_u \text{win}\langle\mathbb{B}\rangle.\text{end}) \end{aligned}$$

where $i \in \{1, 2\}$ and both $T_{i:\text{dice}}$ are obtained by the second case of projection. The type system will ensure that either 3 transmits the request to branch to both players 1, 2 simultaneously and, since these messages cannot be lost, all players that are not crashed move to same branch or 3 crashes and all remaining players move to the default branch.

Assume instead that 3 can only inform one of the players 1, 2 at once. The type

$$\begin{aligned} (\mu t)3 \rightarrow_w \{1\} : \text{play}.3 \rightarrow_u 1:\text{roll}\langle\mathbb{N}\rangle.3 \rightarrow_u 2:\text{roll}\langle\mathbb{N}\rangle.t \\ \oplus \text{end}.3 \rightarrow_u 1:\text{win}\langle\mathbb{B}\rangle.3 \rightarrow_u 2:\text{win}\langle\mathbb{B}\rangle.\text{end} \end{aligned}$$

is not projectable, because \sqcup does not allow to unify the projections $[3]?_u \text{roll}\langle\mathbb{N}\rangle.t$ and $[3]?_u \text{win}\langle\mathbb{B}\rangle.\text{end}$ of the two branches of 2. Replacing the two unreliable communications with 2 by strongly reliable communications implies that neither 3 nor 2 fail. The type

$$\begin{aligned} (\mu t)3 \rightarrow_w \{1\} : \\ \text{play}.3 \rightarrow_u 1:\text{roll}\langle\mathbb{N}\rangle.3 \rightarrow_w \{2\} : (\text{play}.3 \rightarrow_u 2:\text{roll}\langle\mathbb{N}\rangle.t \oplus \text{end}.\text{end}) \\ \oplus \text{end}.3 \rightarrow_u 1:\text{win}\langle\mathbb{B}\rangle.3 \rightarrow_w \{2\} : (\text{play}.\text{end} \oplus \text{end}.3 \rightarrow_u 2:\text{win}\langle\mathbb{B}\rangle.\text{end}) \end{aligned}$$

where 3 informs its two players subsequently about the chosen branch is projectable. But it introduces the two additional branches $\text{end}.\text{end}$ and $\text{play}.\text{end}$, i.e., 3 is allowed to choose the branches for the players 1, 2 separately and differently, whereas in (1) as well as in (3) the players 1, 2 are always in the same branch. Because of that, we allow for broadcast in weakly reliable branching such that 3 can inform both players consistently without introducing additional and not-intended branches.

In types $(\mu t)G$ and $(\mu t)T$ the type variable t is *bound*. In processes $(\mu X)P$ the process variable X is bound. Similarly, all names in round brackets are bound in the remainder of the respective process, e.g. s is bound in P by $\bar{a}[n](s).P$ and x is bound in P by $s[r_1, r_2]?_r(x).P$. A variable or name is *free* if it is not bound. Let $\text{FN}(P)$ return the free names of P .

Let *subterm* denote a (type or process) expression that syntactically occurs within another (type or process) term. We use $' '$ (as e.g. in $\bar{a}[r](s).P$) to denote sequential composition. In all operators the *prefix* before $' '$ guards the *continuation* after the $' '$. Let $\prod_{1 \leq i \leq n} P_i$ abbreviate the parallel composition $P_1 \mid \dots \mid P_n$.

We write $\text{nsr}(G)$, $\text{nsr}(T)$, and $\text{nsr}(P)$, if none of the prefixes in G , T , and P is strongly reliable or for delegation and if P does not contain message queues. Let $\text{R}(G)$ return all roles that occur in G . A global type is *well-formed* if (1) it neither contains free nor unguarded type variables, (2) $\text{R}(G) = \{1, \dots, |\text{R}(G)|\}$, (3) for all its subterms of the form $r_1 \rightarrow_r r_2:\langle S \rangle.G$ and $r_1 \rightarrow_u r_2:l\langle S \rangle.G$, we have $r_1 \neq r_2$, (4) for all its subterms of the form $r_1 \rightarrow_r r_2:\{l_i.G_i\}_{i \in I}$ and $r \rightarrow_w R:\{l_i.G_i\}_{i \in I, l_i}$, we have $r_1 \neq r_2$, $r \notin R$, $d \in I$, and the labels l_i are pairwise distinct, and (5) for all its subterms of the form $G_1 \parallel G_2$, we have $\text{R}(G_1) \cap \text{R}(G_2) = \emptyset$. We restrict our attention to well-formed global types for that projection is defined on all its roles.

$$\begin{array}{ll}
(\text{Init}) & \bar{a}[n](s).P_n \mid \prod_{1 \leq i \leq n-1} a[i](s).P_i \mapsto (\nu s) \left(\prod_{1 \leq i \leq n} P_i \mid \prod_{1 \leq i, j \leq n, i \neq j} s_i \rightarrow j : [] \right) \\
& \text{if } a \neq s \\
(\text{RSend}) & s[r_1, r_2]!_r \langle y \rangle . P \mid s_{r_1 \rightarrow r_2} : M \mapsto P \mid s_{r_1 \rightarrow r_2} : M \# \langle v \rangle^r \quad \text{if } \text{eval}(y) = v \\
(\text{RGet}) & s[r_1, r_2]?_r \langle x \rangle . P \mid s_{r_2 \rightarrow r_1} : \langle v \rangle^r \# M \mapsto P\{v/x\} \mid s_{r_2 \rightarrow r_1} : M \\
(\text{USend}) & s[r_1, r_2]!_u \langle y \rangle . P \mid s_{r_1 \rightarrow r_2} : M \mapsto P \mid s_{r_1 \rightarrow r_2} : M \# \langle v \rangle^u \quad \text{if } \text{eval}(y) = v \\
(\text{UGet}) & s[r_1, r_2]?_u \langle dv \rangle (x) . P \mid s_{r_2 \rightarrow r_1} : l' \langle v \rangle^u \# M \mapsto P\{v/x\} \mid s_{r_2 \rightarrow r_1} : M \\
& \text{if } l \doteq l', \text{FP}_{\text{uget}}(s, r_1, r_2, l') \\
(\text{USkip}) & s[r_1, r_2]?_u \langle dv \rangle (x) . P \mapsto P\{dv/x\} \quad \text{if } \text{FP}_{\text{uskip}}(s, r_1, r_2, l) \\
(\text{ML}) & s_{r_1 \rightarrow r_2} : l \langle v \rangle^u \# M \mapsto s_{r_1 \rightarrow r_2} : M \quad \text{if } \text{FP}_{\text{ml}}(s, r_1, r_2, l) \\
(\text{RSel}) & s[r_1, r_2]!_r l . P \mid s_{r_1 \rightarrow r_2} : M \mapsto P \mid s_{r_1 \rightarrow r_2} : M \# l^r \\
(\text{RBran}) & s[r_1, r_2]?_r \{l_i . P_i\}_{i \in I} \mid s_{r_2 \rightarrow r_1} : l^r \# M \mapsto P_j \mid s_{r_2 \rightarrow r_1} : M \quad \text{if } l \doteq l_j, j \in I \\
(\text{WSel}) & s[r, R]!_w l . P \mid \prod_{r_i \in R} s_{r \rightarrow r_i} : M_i \mapsto P \mid \prod_{r_i \in R} s_{r \rightarrow r_i} : M_i \# l^w \\
(\text{WBran}) & s[r_1, r_2]?_w \{l_i . P_i\}_{i \in I, l_d} \mid s_{r_2 \rightarrow r_1} : l^w \# M \mapsto P_j \mid s_{r_2 \rightarrow r_1} : M \quad \text{if } l \doteq l_j, j \in I \\
(\text{WSkip}) & s[r_1, r_2]?_w \{l_i . P_i\}_{i \in I, l_d} \mapsto P_d \quad \text{if } \text{FP}_{\text{wskip}}(s, r_1, r_2) \\
(\text{Crash}) & P \mapsto \perp \quad \text{if } \text{FP}_{\text{crash}}(P)
\end{array}$$

Figure 2.3: Reduction Rules (\mapsto) of Fault-Tolerant Processes (Part 1).

The combination of a session channel and a role uniquely identifies a participant of a session, called an *actor*. A process has an actor $s[r]$ if it has an action prefix on s , where r is the first role mentioned in the prefix. Let $A(P)$ be the set of actors of P .

2.3.3 A Semantics with Failure Patterns

The application of a substitution $\{y/x\}$ on a term A , denoted as $A\{y/x\}$, is defined as the result of replacing all free occurrences of x in A by y , possibly applying alpha-conversion to avoid capture or name clashes. For all names $n \in \mathcal{N} \setminus \{x\}$ the substitution behaves as the identity mapping. We use substitution on types as well as processes and naturally extend substitution to the substitution of variables by terms (to unfold recursions) and names by expressions (to instantiate a bound name with a received value).

We use labels for two purposes: they allow us to distinguish between different branches, as usual in MPST-frameworks, and we assume that they may carry additional runtime information such as timestamps. Of course, the presented type system remains valid if we use labels without additional information. In contrast to standard MPST (as e.g. in [9]) and to support unreliable communication, our MPST variant will ensure that all occurrences of the same label are associated with the same sort. We assume a predicate \doteq that compares two labels and is valid if the parts of the labels that do not refer to runtime information correspond. If runtime information are irrelevant, \doteq can be instantiated with equality. We require that \doteq is unambiguous on labels used in types, i.e., given two labels of processes l_P, l'_P and two labels of types l_T, l'_T then $l_P \doteq l'_P \wedge l_P \doteq l_T \Rightarrow l'_P \doteq l'_T$ and $l_P \doteq l_T \wedge l_T \doteq l'_T \Rightarrow l_P \doteq l'_T$.

We use structural congruence to abstract from syntactically different processes with the same meaning, where \equiv is the least congruence that satisfies alpha conversion and the rules:

$$\begin{array}{lll}
P \mid \mathbf{0} \equiv P & P_1 \mid P_2 \equiv P_2 \mid P_1 & P_1 \mid (P_2 \mid P_3) \equiv (P_1 \mid P_2) \mid P_3 \\
(\nu X)\mathbf{0} \equiv \mathbf{0} & (\nu x)\mathbf{0} \equiv \mathbf{0} & (\nu x)(\nu y)P \equiv (\nu y)(\nu x)P \\
(\nu x)(P_1 \mid P_2) \equiv P_1 \mid (\nu x)P_2 & \text{if } x \notin \text{FN}(P_1)
\end{array}$$

(If-T)	$\text{if } e \text{ then } P \text{ else } P' \mapsto P$	if e is true
(If-F)	$\text{if } e \text{ then } P \text{ else } P' \mapsto P'$	if e is false
(Deleg)	$s[r_1, r_2]!\langle\langle s'[r] \rangle\rangle.P \mid s_{r_1 \rightarrow r_2}:M \mapsto P \mid s_{r_1 \rightarrow r_2}:M\#s'[r]$	
(SRecv)	$s[r_1, r_2]?(\langle s'[r] \rangle).P \mid s_{r_2 \rightarrow r_1}:s''[r']\#M \mapsto P\{s''/s'\}\{r'/r\} \mid M$	
(Par)	$P_1 \mid P_2 \mapsto P'_1 \mid P_2$	if $P_1 \mapsto P'_1$
(Res)	$(\nu x)P \mapsto (\nu x)P'$	if $P \mapsto P'$
(Rec)	$(\mu X)P \mapsto P\{(\mu X)P/X\}$	
(Struc)	$P_1 \mapsto P'_1$	if $P_1 \equiv P_2, P_2 \mapsto P'_2, P'_2 \equiv P'_1$

Figure 2.4: Reduction Rules (\mapsto) of Fault-Tolerant Processes (Part 2).

The reduction semantics of the session calculus is defined in Fig. 2.3 and Fig. 2.4, where we follow [9]: we assume that session initialisation is synchronous and communication within a session is asynchronous implemented using message queues.

Rule (Init) initialises a session with n roles. Session initialisation introduces a fresh session channel and unguards the participants of the session. Finally, the message queues of this session are initialised with the empty list under the restriction of the session channel.

Rule (RSend) implements an asynchronous strongly reliable message transmission. As a result the value $\text{eval}(y)$ is wrapped in a message and added to the end of the corresponding message queue and the continuation of the sender is unguarded. Rule (USend) is the counterpart of (RSend) for unreliable senders. (RGet) consumes a message that is marked as strongly reliable with the index r from the head of the respective message queue and replaces in the unguarded continuation of the receiver the bound variable x by the received value y .

There are two rules for the reception of a message in an unreliable communication that are guided by failure patterns. *Failure patterns* are predicates that we deliberately choose not to define here (see below). They allow us to provide information about the underlying communication medium and the reliability of processes. Rule (UGet) is similar to Rule (RGet), but specifies a failure pattern FP_{uget} to decide whether this step is allowed. This failure pattern could, e.g., be used to reject messages that are too old. The Rule (USkip) allows to skip the reception of a message in an unreliable communication using a failure pattern FP_{uskip} and instead substitutes the bound variable x in the continuation with the default value dv . The failure pattern FP_{uskip} tells us whether a reception can be skipped (e.g. via failure detector).

Rule (RSeI) puts the label l selected by r_1 at the end of the message queue towards r_2 . Its weakly reliable counterpart (WSeI) is similar, but puts the label at the end of all relevant message queues. With (RBran) a label is consumed from the top of a message queue and the receiver moves to the indicated branch. There are again two weakly reliable counterparts of (RBran). Rule (WBran) is similar to (RBran), whereas (WSkip) allows r_1 to skip the message and to move to its default branch if the failure pattern FP_{wskip} holds.

The Rules (Crash) for *crash failures* and (ML) for *message loss*, describe failures of a system. With Rule (Crash) P can crash if FP_{crash} , where FP_{crash} can e.g. model immortal processes or global bounds on the number of crashes. (ML) allows to drop an unreliable message if the failure pattern FP_{ml} is valid. FP_{ml} allows, e.g., to implement safe channels that never lose messages or a global bound on the number of lost messages.

The remaining rules for conditionals, session delegation, parallel composition, restriction, recursion, and structural congruence in Fig. 2.4 are standard.

We augmented our reduction semantics in Fig. 2.3 by five different failure patterns that we deliberately do not specify, although we usually assume that the failure patterns FP_{uget} , FP_{uskip} , and FP_{wskip} use only local

information, whereas FP_{ml} and FP_{crash} may use global information of the system in the current run. We provide these predicates to allow for the implementation of system requirements or abstractions like failure detectors that are typical for distributed algorithms. Directly including them in the semantics has the advantage that all traces satisfy the corresponding requirements, i.e., all traces are valid w.r.t. the assumed system requirements. An example for the instantiation of these patterns is given implicitly via the Conditions 1.1–1.6 in Section 2.3.4 and explicitly in Section ?? . If we instantiate the patterns FP_{uget} with true and the patterns FP_{uskip} , FP_{wskip} , FP_{crash} , FP_{ml} with false, then we obtain a system without failures. In contrast, the instantiation of all five patterns with true results in a system where failures can happen completely non-deterministically at any time.

One of the properties that our type system has to ensure even in the case of failures is the absence of communication mismatches, i.e., the type of a transmitted value has to be the type that the receiver expects. The global type $1 \rightarrow_u 2:l_1\langle\mathbb{N}\rangle.1 \rightarrow_u 2:l_2\langle\mathbb{B}\rangle.\text{end}$ specifies two subsequent unreliable communications in that values of different sorts are transmitted. If the first message with its natural number is lost but the second message containing a Boolean value is transmitted, 2 could wrongly receive a Boolean value although it still waits for a natural number. To avoid this mismatch, we add a label to unreliable communication and ensure (by the typing rules) that the same label is never associated with different types. Similarly, labels are used in [4] to avoid communication errors. Accordingly, we require $l_1 \neq l_2$ such that the reduction rules in Fig. 2.3 will not allow to consume the Boolean message before 2 has reduced its first prefix.

We do that, because we think of labels not only as identifiers for branching but also as some kind of meta data of messages as they can be often found in communication media or as they are assumed by many distributed algorithms. Our unreliable communication mechanism exploits such meta data to guarantee strong properties about the communication structure including the described absence of communication mismatches.

Note that we keep the failure patterns abstract and do not model how to check them in producing runs. Indeed system requirements such as bounds on the number of processes that can crash usually cannot be checked, but result from observations, i.e., system designers ensure that a violation of this bound is very unlikely and algorithm designers are willing to ignore these unlikely events. In particular, FP_{ml} and FP_{crash} are thus often implemented as oracles for verification, whereas e.g. FP_{uskip} and FP_{wskip} are often implemented by system specific time-outs. Note that we are talking about implementing these failure patterns and not formalising them. Failure patterns are abstractions of real world system requirements or software. We implement them by conditions providing the necessary guarantees that we need in general (i.e., for subject reduction and progress) or for the verification of concrete algorithms. In practise, we expect that the systems on that the verified algorithms are running satisfy the respective conditions. Accordingly, the session channels, roles, labels, and processes mentioned in Fig. 2.3 are not parameters of the failure patterns, but just a vehicle to more formally specify the conditions on failure patterns in Section 2.3.4.

Similarly, strongly reliable and weakly reliable interactions in potentially faulty systems are abstractions. They are usually implemented by handshakes and redundancy; replicated servers against crash failures and retransmission of late messages against message loss. Algorithm designers have to be aware of the additional costs of these interactions.

Consider the implementation of $G_{dice,u}$ in (2), i.e., an infinite variant of the dice game, where the players 1 and 2 use their respective last known sum x_i of former dice rolls as default values:

$$\begin{aligned} P_{dice,u} &= P_{3,u} \mid P_{i,u} \mid P_{2,u} \\ P_{3,u} &= \bar{a}[3](s).(\mu X)s[3,1]!_u \text{roll}\langle \text{roll}(x_1) \rangle.s[3,2]!_u \text{roll}\langle \text{roll}(x_2) \rangle.X \\ P_{i,u} &= a[i](s).(\mu X)s[i,3]?_u \text{roll}\langle x_i \rangle(x_i).X \end{aligned}$$

An unreliable communication in a global type specifies a communication that, due to system failures, may or may not happen. Moreover, regardless of the successful completion of this unreliable communication, the future behaviour of a well-typed system will follow its specification in the global type. Since the players 1 and 2 repeat the same kind of unreliable action, they may lose track of the current round. If they successfully receive a new sum of dice rolls from 3 they cannot be sure on how often 3 actually did roll the dice. Because of lost messages, they may have missed some former announcements of 3 and, because of their ability to skip the reception of messages, they may have proceeded to the next round before 3 rolled a dice. Because the information about the current round is irrelevant for the communication structure in this case, there is no need to enforce round information.

2.3.4 Typing Fault-Tolerant Processes

The type of a process P is checked in a *typed judgment*, i.e., triples $\Gamma \vdash P \triangleright \Delta$, where

$$\begin{aligned} \Gamma &::= \emptyset \mid \Gamma \cdot x:S \mid \Gamma \cdot a:G \mid \Gamma \cdot l:S \\ \Delta &::= \emptyset \mid \Delta \cdot s[r]:T \mid \Delta \cdot s_{r_1 \rightarrow r_2}:MT^* \end{aligned}$$

Assignments in Γ relate variables to their sort, shared channels to the type of the session they introduce, and connect labels with a sort. Session environments collect the local types of actors and the list of message types of queues, i.e., MT^* denotes a list of message types.

We write $x \nmid \Gamma$ and $x \nmid \Delta$ if the name x does not occur in Γ and Δ , respectively. We use \cdot to add an assignment provided that the new assignment is not in conflict with the type environment, i.e., $\Gamma \cdot A$ implies that the respective name/variable/label in A is not contained in Γ and $\Delta \cdot A$ implies that the respective actor/queue in A is not contained in Δ . These conditions on \cdot for global and session environments are referred to as *linearity*. We restrict in the following our attention to linear environments.

We write $\text{nsr}(\Delta)$ if $\text{nsr}(T)$ for all local types T in Δ and if Δ does not contain message queues. With $\Gamma \Vdash y:S$ we check that y is an expression of the sort S if all names x in y are replaced by arbitrary values of sort S_x for $x:S_x \in \Gamma$.

A process P is *well-typed* w.r.t. Γ and Δ if $\Gamma \vdash P \triangleright \Delta$ can be derived from the rules in the Fig. 2.5 and 2.6. We concentrate on the interaction cases, where we observe that all new cases are quite similar to their strongly reliable counterparts.

Rule (RSend) checks strongly reliable senders, i.e., requires a matching strongly reliable sending in the local type of the actor and compares the actor with this type. With $\Gamma \Vdash y:S$ we check that y is an expression of the sort S if all names x in y are replaced by arbitrary values of sort S_x for $x:S_x \in \Gamma$. Then the continuation of the process is checked against the continuation of the type. The unreliable case is very similar, but additionally checks that the label is assigned to the sort of the expression in Γ . Rule (RGet) type strongly reliable receivers, where again the prefix is checked against a corresponding type prefix and the assumption $x:S$ is added to the type check of the continuation. Again the unreliable case is very similar, but apart from the label also checks the sort of the default value.

Rule (RSeI) checks the strongly reliable selection prefix, that the selected label matches one of the specified labels, and that the process continuation is well-typed w.r.t. the type continuation following the selected label. The only difference in the weakly reliable case is the set of roles for the receivers. For strongly reliable branching we have to check the prefix and that for each branch in the type there is a matching branch in

$$\begin{array}{c}
\text{(Req)} \frac{a:G \in \Gamma \quad |\mathbf{R}(G)| = n \quad \Gamma \vdash P \triangleright \Delta \cdot s[n]:G|_n}{\Gamma \vdash \bar{a}[n](s).P \triangleright \Delta} \quad \text{(Acc)} \frac{a:G \in \Gamma \quad 0 < r < |\mathbf{R}(G)| \quad \Gamma \vdash P \triangleright \Delta \cdot s[r]:G|_r}{\Gamma \vdash a[r](s).P \triangleright \Delta} \\
\\
\text{(RSend)} \frac{\Gamma \Vdash y:S \quad \Gamma \vdash P \triangleright \Delta \cdot s[r_1]:T}{\Gamma \vdash s[r_1, r_2]!_r \langle y \rangle . P \triangleright \Delta \cdot s[r_1]:[r_2]!_r \langle S \rangle . T} \\
\text{(RGet)} \frac{x^\sharp(\Gamma, \Delta, s) \quad \Gamma \cdot x:S \vdash P \triangleright \Delta \cdot s[r_1]:T}{\Gamma \vdash s[r_1, r_2]?_r \langle x \rangle . P \triangleright \Delta \cdot s[r_1]:[r_2]?_r \langle S \rangle . T} \\
\text{(USend)} \frac{\Gamma \Vdash y:S \quad l \doteq l' \quad l':S \in \Gamma \quad \Gamma \vdash P \triangleright \Delta \cdot s[r_1]:T}{\Gamma \vdash s[r_1, r_2]!_u \langle y \rangle . P \triangleright \Delta \cdot s[r_1]:[r_2]!_u \langle l' \rangle \langle S \rangle . T} \\
\text{(UGet)} \frac{x^\sharp(\Gamma, \Delta, s) \quad \Gamma \Vdash v:S \quad l \doteq l' \quad l':S \in \Gamma \quad \Gamma \cdot x:S \vdash P \triangleright \Delta \cdot s[r_1]:T}{\Gamma \vdash s[r_1, r_2]?_u \langle v \rangle \langle x \rangle . P \triangleright \Delta \cdot s[r_1]:[r_2]?_u \langle l' \rangle \langle S \rangle . T} \\
\\
\text{(RSel)} \frac{j \in I \quad l \doteq l_j \quad \Gamma \vdash P \triangleright \Delta \cdot s[r_1]:T_j}{\Gamma \vdash s[r_1, r_2]!_r l . P \triangleright \Delta \cdot s[r_1]:[r_2]!_r \{l_i.T_i\}_{i \in I}} \quad \text{(Var)} \frac{}{\Gamma \cdot X:t \vdash X \triangleright s[r]:t} \\
\text{(RBran)} \frac{\forall j \in I_2. \exists i \in I_1. l_i \doteq l_j \wedge \Gamma \vdash P_i \triangleright \Delta \cdot s[r_1]:T_j}{\Gamma \vdash s[r_1, r_2]?_r \{l_i.P_i\}_{i \in I_1} \triangleright \Delta \cdot s[r_1]:[r_2]?_r \{l_i.T_i\}_{i \in I_2}} \\
\text{(WSel)} \frac{j \in I \quad l \doteq l_j \quad \Gamma \vdash P \triangleright \Delta \cdot s[r]:T_j}{\Gamma \vdash s[r, R]!_w l . P \triangleright \Delta \cdot s[r]:[R]!_w \{l_i.T_i\}_{i \in I}} \quad \text{(Crash)} \frac{\text{nsr}(\Delta)}{\Gamma \vdash \perp \triangleright \Delta} \\
\text{(WBran)} \frac{l_d \doteq l'_d \quad \forall j \in I_2. \exists i \in I_1. l_i \doteq l_j \wedge \Gamma \vdash P_i \triangleright \Delta \cdot s[r_1]:T_j}{\Gamma \vdash s[r_1, r_2]?_w \{l_i.P_i\}_{i \in I_1, l_d} \triangleright \Delta \cdot s[r_1]:[r_2]?_w \{l_i.T_i\}_{i \in I_2, l'_d}} \\
\\
\text{(Deleg)} \frac{\Gamma \vdash P \triangleright \Delta \cdot s[r_1]:T}{\Gamma \vdash s[r_1, r_2]! \langle s'[r] \rangle . P \triangleright \Delta \cdot s[r_1]:[r_2]! \langle s'[r]:T' \rangle . T \cdot s'[r]:T'} \\
\text{(SRecv)} \frac{\Gamma \vdash P \triangleright \Delta \cdot s[r_1]:T \cdot s'[r]:T'}{\Gamma \vdash s[r_1, r_2]? \langle s'[r] \rangle . P \triangleright \Delta \cdot s[r_1]:[r_2]? \langle s'[r]:T' \rangle . T} \quad \text{(End)} \frac{}{\Gamma \vdash \mathbf{0} \triangleright \emptyset} \\
\\
\text{(If)} \frac{\Gamma \Vdash e:\mathbb{B} \quad \Gamma \vdash P \triangleright \Delta \quad \Gamma \vdash P' \triangleright \Delta}{\Gamma \vdash \text{if } e \text{ then } P \text{ else } P' \triangleright \Delta} \quad \text{(Par)} \frac{\Gamma \vdash P \triangleright \Delta \quad \Gamma \vdash P' \triangleright \Delta'}{\Gamma \vdash P \mid P' \triangleright \Delta \cdot \Delta'} \\
\\
\text{(Res1)} \frac{x^\sharp(\Gamma, \Delta) \quad \Gamma \cdot x:S \vdash P \triangleright \Delta}{\Gamma \vdash (\nu x)P \triangleright \Delta} \quad \text{(Rec)} \frac{\Gamma \cdot X:t \vdash P \triangleright \Delta \cdot s[r]:T}{\Gamma \vdash (\mu X)P \triangleright \Delta \cdot s[r]:(\mu t)T}
\end{array}$$

Figure 2.5: Typing Rules for Fault-Tolerant Systems.

$$\begin{array}{c}
\frac{\{s[r]:G\}_r \mid r \in R(G)\} \cdot \{s_{r \rightarrow r'}:[] \mid r, r' \in R(G') \wedge r \neq r'\} \stackrel{s}{\Rightarrow} \Delta' \quad s^\sharp(\Gamma, \Delta) \quad a:G \in \Gamma \quad \Gamma \vdash P \triangleright \Delta \cdot \Delta'}{(\text{Res2}) \quad \Gamma \vdash (\nu s)P \triangleright \Delta} \\
\\
(\text{MQComR}) \quad \frac{\Gamma \Vdash v:S \quad \Gamma \vdash s_{r_1 \rightarrow r_2}:M \triangleright s_{r_1 \rightarrow r_2}:MT}{\Gamma \vdash s_{r_1 \rightarrow r_2}:\langle v \rangle^r \# M \triangleright s_{r_1 \rightarrow r_2}:\langle S \rangle^r \# MT} \\
(\text{MQComU}) \quad \frac{\Gamma \Vdash v:S \quad l \doteq l' \quad l':S \in \Gamma \quad \Gamma \vdash s_{r_1 \rightarrow r_2}:M \triangleright s_{r_1 \rightarrow r_2}:MT}{\Gamma \vdash s_{r_1 \rightarrow r_2}:l\langle v \rangle^u \# M \triangleright s_{r_1 \rightarrow r_2}:l'\langle S \rangle^u \# MT} \\
(\text{MQBranR}) \quad \frac{l \doteq l' \quad \Gamma \vdash s_{r_1 \rightarrow r_2}:M \triangleright s_{r_1 \rightarrow r_2}:MT}{\Gamma \vdash s_{r_1 \rightarrow r_2}:l^r \# M \triangleright s_{r_1 \rightarrow r_2}:l'^r \# MT} \\
(\text{MQBranW}) \quad \frac{l \doteq l' \quad \Gamma \vdash s_{r_1 \rightarrow r_2}:M \triangleright s_{r_1 \rightarrow r_2}:MT}{\Gamma \vdash s_{r_1 \rightarrow r_2}:l^w \# M \triangleright s_{r_1 \rightarrow r_2}:l'^w \# MT} \\
(\text{MQDeleg}) \quad \frac{\Gamma \vdash s_{r_1 \rightarrow r_2}:M \triangleright s_{r_1 \rightarrow r_2}:MT}{\Gamma \vdash s_{r_1 \rightarrow r_2}:s'[r] \# M \triangleright s_{r_1 \rightarrow r_2}:s'[r] \# MT} \quad (\text{MQNil}) \quad \frac{}{\Gamma \vdash s_{r_1 \rightarrow r_2}:[] \triangleright s_{r_1 \rightarrow r_2}:[]}
\end{array}$$

Figure 2.6: Runtime Typing Rules for Fault-Tolerant Systems.

the process that is well-typed w.r.t. the respective branch in the type. For the weakly reliable case we have to additionally check that the default labels of the process and the type coincide.

Rule (Crash) for crashed processes checks that $\text{nsr}(\Delta)$.

Figure 2.6 presents the runtime typing rules, i.e., the typing rules for processes that may result from steps of a system that implements a global type. Since it covers only operators that are not part of initial systems, a type checking tool might ignore them. We need these rules however for the proofs of progress and subject reduction. Under the assumption that initial systems cannot contain crashed processes, Rule (Crash) may be moved to the set of runtime typing rules.

We have to prove that our extended type system satisfies the standard properties of MPST, i.e., subject reduction and progress. *Subject reduction* tells us that derivatives of well-typed systems are again well-typed. This is fundamental, since it ensures that our formalism can be used to analyse processes by static type checking. We extend subject reduction such that it provides some information on how the session environment evolves alongside reductions of the system. Therefore we introduce a reduction relation $\stackrel{s}{\rightarrow}$ on session environments, that emulates the reduction steps of processes. As an example consider the rule $\Delta \cdot s[r_1]:[r_2]!_r \langle S \rangle.T \cdot s_{r_1 \rightarrow r_2}:MT \stackrel{s}{\rightarrow} \Delta \cdot s[r_1]:T \cdot s_{r_1 \rightarrow r_2}:MT \# \langle S \rangle^r$ that emulates the transfer of a value in (RSend). Let $\stackrel{s}{\Rightarrow}$ denote the reflexive and transitive closure of $\stackrel{s}{\rightarrow}$.

Coherence intuitively describes that a session environment captures all local endpoints of a collection of global types. Since we capture all relevant global types in the global environment, we define coherence on pairs of global and session environments.

Definition 1 (Coherence) *The type environments Γ, Δ are coherent if, for all session channels s in Δ , there exists a global type G in Γ such that the restriction of Δ on assignments with s is the set Δ' such that $\{s[r]:G\}_r \mid r \in R(G)\} \cdot \{s_{r \rightarrow r'}:[] \mid r, r' \in R(G)\} \stackrel{s}{\Rightarrow} \Delta'$.*

Because of the failure pattern in the reduction semantics in Fig. 2.3, subject reduction and progress do not hold in general. Instead we have to fix conditions on failure patterns that ensure these properties. Subjection

reduction needs one condition on crashed processes and progress requires that the instantiation of the failure patterns is such that they do not block parts of the system. In fact, different instantiations of these failure patterns may allow for progress. We leave it to further work to determine what kind of conditions on failure patterns or requirements on their interactions are necessary to prove these properties. Here, we consider only one such set of conditions.

Condition 1 (Failure Pattern) 1. If $\text{FP}_{\text{crash}}(P)$, then $\text{nsr}(P)$.

2. The failure pattern $\text{FP}_{\text{uget}}(s, r_1, r_2, l)$ is always valid.

3. The pattern $\text{FP}_{\text{ml}}(s, r_1, r_2, l)$ is valid iff $\text{FP}_{\text{uskip}}(s, r_2, r_1, l)$ is valid.

4. If $\text{FP}_{\text{crash}}(P)$ and $s[r] \in A(P)$ then eventually $\text{FP}_{\text{uskip}}(s, r_2, r, l)$ and also $\text{FP}_{\text{wskip}}(s, r_2, r, l)$ for all r_2, l .

5. If $\text{FP}_{\text{crash}}(P)$ and $s[r] \in A(P)$ then eventually $\text{FP}_{\text{ml}}(s, r_1, r, l)$ for all r_1, l .

6. If $\text{FP}_{\text{wskip}}(s, r_1, r_2)$ then $s[r_2]$ is crashed, i.e., the system does no longer contain an actor $s[r_2]$ and the message queue $s_{r_2 \rightarrow r_1}$ is empty.

The crash of a process should not block strongly reliable actions, i.e., only processes with $\text{nsr}(P)$ can crash (Condition 1.1). Condition 1.2 requires that no process can refuse to consume a message on its queue. This condition prevents deadlocks that may arise from refusing a message m that is never dropped from the message queue. Condition 1.3 requires that if a message can be dropped from a message queue then the corresponding receiver has to be able to skip this message and vice versa. Similarly, processes that wait for messages from a crashed process have to be able to skip (Condition 1.4) and all messages of a queue towards a crashed receiver can be dropped (Condition 1.5). Finally, weakly reliable branching requests should not be lost. To ensure that the receiver of such a branching request can proceed if the sender is crashed but is not allowed to skip the reception of the branching request before the sender crashed, we require that $\text{FP}_{\text{wskip}}(s, r_1, r_2)$ is false as long as $s[r_2]$ is alive or messages on the respective queue are still in transit (Condition 1.6).

The combination of these 6 conditions might appear quite restrictive on a first glance. For example the combination of the Conditions 1.4 and 1.6 ensures the correct behaviour of weakly reliable branching such that branching messages can be skipped if and only if the respective sender has crashed. An implementation of such a weakly reliable interaction in an asynchronous system that is subject to message losses and process crashes, might require something like a perfect failure detector or actually solving consensus¹. Here it is important to remember that these conditions are minimal assumptions on the system requirements and that system requirements are just abstractions. Parts of them may be realised by actual software-code (which then allows to check them), whereas other parts of the system requirements may not be realised at all but rather observed (which then does not allow to verify them). Weakly reliable branching is a good example of this case. The easiest way to obtain a weakly reliable interaction, is by using a handshake communication and time-outs. If the sender time-outs while waiting for an acknowledgement, it resends the message. If the sender does not hear from its receiver for a long enough period of time, it assumes that the receiver has crashed and proceeds. With carefully chosen time-frames for the time-outs, this approach is a compromise between correctness and efficiency. In a theoretical sense, it is clearly not correct. There is no time-frame such that the sender can be really sure that the receiver has crashed. From a practical point of view, this is not so problematic, since in many systems failures are very unlikely. If failures that are so severe that they are not captured by the time-outs are extremely unlikely, then it is often much more efficient to just accept

¹Note that the example we present in Section ?? is a consensus algorithm. So, if the Conditions 1 require a solution of consensus, an example on top of that solving consensus would be pointless.

that the algorithm is not correct in these cases. Trying to obtain an algorithm that is always correct might be impossible or at least usually results into very inefficient algorithms. Moreover, verifying this requires to also verify the underlying communication infrastructure and the way in that failures may occur, which is impossible or at least impracticable. Because of that, it is an established method to verify the correctness of algorithms w.r.t. given system requirements (e.g. in [5, 11, 15]), even if these system requirements are not verified and often do not hold in all (but only nearly all) cases.

With Conditions 1, we can now analyse our fault-tolerant type system.

Theorem 1 (Subject Reduction) *If $\Gamma \vdash P \triangleright \Delta$, Γ, Δ are coherent, and $P \mapsto P'$, then there are some Δ', s such that $\Gamma \vdash P' \triangleright \Delta'$, Γ, Δ' are coherent, and $\Delta \xrightarrow{s} \Delta'$.*

The proof is by induction on the derivation of $P \mapsto P'$. In every case, we use the information about the structure of the processes to generate partial proof trees for the respective typing judgement. Additionally, we use Condition 1.1 to ensure that the type environment of a crashed process cannot contain the types of reliable communication prefixes.

Progress states that no part of a well-typed and coherent system can block other parts, that eventually all matching communication partners of strongly reliable and weakly reliable communications (that are not crashed) are unguarded, and that there are no communication mismatches. Subject reduction and progress together then imply *session fidelity*, i.e., that processes behave as specified in their global types.

To ensure that the interleaving of sessions and session delegation cannot introduce deadlocks, we assume an interaction type system as introduced in [2, 9]. For this type system it does not matter whether the considered actions are strongly reliable, weakly reliable, or unreliable. More precisely, we can adapt the interaction type system of [2] in a straightforward way to the above session calculus, where unreliable communication and weakly reliable branching is treated in exactly the same way as strongly reliable communication/branching. We say that P is *free of cyclic dependencies between sessions* if this interaction type system does not detect any cyclic dependencies. In this sense fault-tolerance is more flexible than explicit failure handling as e.g. [16] has to exclude interleaved sessions.

In the literature there are different formulations of progress. We are interested in a rather strict definition of progress that ensures that well-typed systems cannot block. Therefore, we need an additional assumption on session requests and acceptances. Coherence ensures the existence of communication partners within sessions only. If we want to avoid blocking, we need to be sure, that no participant of a session is missing during its initialisation. Note that without action prefixes all participants either terminated or crashed.

Theorem 2 (Progress/Session Fidelity) *Let $\Gamma \vdash P \triangleright \Delta$, Γ, Δ be coherent, and let P be free of cyclic dependencies between sessions. Assume that in the derivation of $\Gamma \vdash P \triangleright \Delta$, whenever $\bar{a}[n](s).Q$ or $a[r](s).Q$ with $a:G$, then there are $\bar{a}[n](s).Q_n$ or $a[r_i](s).Q_i$ for all $1 \leq r_i < n$.*

1. *Then either P does not contain any action prefixes or $P \mapsto P'$.*
2. *If P does not contain recursion, then there exists P' such that $P \mapsto^* P'$ and P' does not contain any action prefixes.*

2.4 Additional Notation

In [13] the authors introduced notation for the construction of global types and processes.

Let $(\odot_{1 \leq i \leq n} \pi_i).G$ abbreviate the sequence $\pi_1 \dots \pi_n.G$ to simplify the presentation, where $G \in \mathcal{G}$ is a global type and π_1, \dots, π_n are sequences of prefixes. More precisely, each π_i is of the form $\pi_{i,1} \dots \pi_{i,m}$ and each $\pi_{i,j}$ is a type prefix of the form $r_1 \rightarrow_u r_2 : l \langle S \rangle$ or $r \rightarrow_w R : l_1.T_1 \oplus \dots \oplus l_n.T_n \oplus l_d$, where the latter case represents a weakly reliable branching prefix with the branches l_1, \dots, l_n, l_d , the default branch l_d , and where the next global type provides the missing specification for the default case.

Let $(\odot_{1 \leq i \leq n} \pi_i).P$ abbreviate the sequence $\pi_1 \dots \pi_n.P$, where $P \in \mathcal{P}$ is a process and π_1, \dots, π_n are sequences of prefixes.

3 Model

First, we specify some sorts with which we can then define the global type. Afterwards, we define the processes for the proposer and the acceptor. Finally, we will study an example run of the model.

3.1 Sorts

Sorts are basic data types. We assume the following sorts.

First, we have `Bool` which we define as a set.

$$\text{Bool} = \{\text{true}, \text{false}\}$$

Second, we assume a set of values `Value`.

Then, we have some sorts which we define using a grammar. Each of these definitions contains a type variable, which is a variable ranging over types. In this case the type variable in each definition is called a .

$$\text{Maybe } a = \text{Just } a \mid \text{Nothing}$$

A value of type `Maybe a` can have the form `Just a` or `Nothing`. Some examples include `Just 4` of type `Maybe \mathbb{N}` , `Just false` of type `Maybe Bool`, and `Nothing`. `Nothing` itself does not dictate an exact type because its definition does not include the type variable a . The type is underspecified and is specified manually or through the context in which `Nothing` is used. It can be of type `Maybe \mathbb{N}` , `Maybe Bool`, or any other type b in `Maybe b` . We use `Maybe a` where optional values are needed.

$$\text{Proposal } a = \text{Proposal } \mathbb{N} \ a$$

`Proposal a` only has one possible form, which is `Proposal $\mathbb{N} \ a$` . A proposal contains its proposal number of type \mathbb{N} and its value of type a . Again, a is a variable ranging over types. An example for a value of type `Proposal Bool` could be `Proposal 1 true` and an example for a value of type `Proposal Maybe \mathbb{N}` could be `Proposal 1 Just 1`. Note that `Proposal 1 Just 1` is of type `Proposal a` where $a = \text{Maybe } b$ and $b = \mathbb{N}$. This sort models the proposals issued by the proposers in phase $2a$.

$$\text{Promise } a = \text{Promise Maybe Proposal } a \mid \text{Nack } \mathbb{N}$$

Promise a has two possible forms. Promise Maybe Proposal a and Nack \mathbb{N} . Promise Maybe Proposal a is the same as Promise c where $c = \text{Maybe } b$ and $b = \text{Proposal } a$. Possible values include Nack 1 and Promise Just Proposal 1 – 1 of type Promise \mathbb{Z} . The actual type of Nack 1, much like that of Nothing, is underspecified. Again, we have to specify the exact type manually or through context.

In phase 1b the acceptors respond to the proposers prepare request with a value of type Promise Value. The prepare request contains a number n . The acceptors may respond to the prepare request with a promise to not accept any proposal numbered less than n or with a rejection. In the first case the acceptor's response optionally includes the last proposal it accepted, if available, and is of the form Promise Maybe Proposal a . In the second case it includes the highest n that acceptor promised and is of the form Nack \mathbb{N} .

3.2 Global Type

Since each proposer initiates its own session the global type can be defined for one proposer p and a quorum of acceptors A_Q .

The last phase of Paxos contains no inter-process communication, so it is not modeled in the global type.

$$G_{P,A_Q} = (\mu X) \left(\bigodot_{a \in A_Q} p \rightarrow_u a : l1a \langle \mathbb{N} \rangle \right) . \left(\bigodot_{a \in A_Q} a \rightarrow_u p : l1b \langle \text{Promise Value} \rangle \right) . \\ \left(p \rightarrow_w A_Q : \text{Accept} . \left(\bigodot_{a \in A_Q} p \rightarrow_u a : l2a \langle \text{Proposal Value} \rangle \right) . 0 \oplus \text{Restart} . X \oplus \text{Abort} . 0 \right)$$

We can distinguish the individual phases of the Paxos algorithm by the labels $l1a$, $l1b$, and $l2a$.

In the first two steps, 1a and 1b, the proposer sends its proposal number to each acceptor in A_Q and listens for their responses. In step 2a the proposer decides whether to send an *Accept* or *Restart* message to restart the algorithm. This decision is broadcast to all acceptors in A_Q . Should the proposer crash the algorithm ends for this particular proposer and quorum of acceptors.

3.3 Functions

We define some functions which we use in the next section to define the processes.

$$\text{proposalNumber} : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$$

$\text{proposalNumber}_p(n)$ returns a proposal number for proposer p when given a natural number n . It is used to pick a number for the prepare request in phase 1a, which is also used in phase 2a in the actual proposal. We have two requirements for this function.

Let \mathbb{P} be the set of proposers.

$$\forall p, q \in \mathbb{P}. \forall n, m \in \mathbb{N} : p \neq q \rightarrow \text{proposalNumber}_p(n) \neq \text{proposalNumber}_q(m)$$

Different proposers pick their proposal numbers from disjoint sets of numbers. This way different proposers never issue a proposal with the same proposal number.

$$\forall p \in \mathbb{P}. \forall n, m \in \mathbb{N} : n > m \rightarrow \text{proposalNumber}_p(n) > \text{proposalNumber}_p(m)$$

We require $\text{proposalNumber}_p(n)$ to be strictly increasing for each proposer p so every proposer uses a higher proposal number than any it has already used.

$$\text{promiseValue} : \text{list of Promise } a \rightarrow a$$

$\text{promiseValue}(ps)$ returns a fresh value if none of the promises in ps contain a value. Otherwise, the best value is returned. Usually, that means the value with the highest associated proposal number. A promise contains a value v if it is of the form $\text{Promise Just } v$. With this function we can model the picking of a value for a proposal in phase 2a.

$$\begin{aligned} \text{anyNack} &: \text{list of Promise } a \rightarrow \text{Bool} \\ \text{anyNack}([]) &= \text{false} \\ \text{anyNack}((\text{Nack } _ \# _)) &= \text{true} \\ \text{anyNack}((_ \# xs)) &= \text{anyNack}(xs) \end{aligned}$$

$\text{anyNack}(ps)$ returns true if the list contains at least one promise of the form $\text{Nack } n$. Otherwise, it returns false.

$$\begin{aligned} \text{promiseCount} &: \text{list of Promise } a \rightarrow \mathbb{N} \\ \text{promiseCount}([]) &= 0 \\ \text{promiseCount}((\text{Promise } _ \# xs)) &= 1 + \text{promiseCount}(xs) \\ \text{promiseCount}((_ \# xs)) &= \text{promiseCount}(xs) \end{aligned}$$

$\text{promiseCount}(ps)$ takes a list of promises ps and calculates the number of promises in that list of that have the form $\text{Promise } m$.

$\text{anyNack}(ps)$ and $\text{promiseCount}(ps)$ are used in the proposer to decide which branch to take in phase 2a.

$$\begin{aligned} \text{gt} &: a \rightarrow \text{Maybe } a \rightarrow \text{Bool} \\ \text{gt}(_, \text{Nothing}) &= \text{true} \\ \text{gt}(a, \text{Just } b) &= a > b \end{aligned}$$

$$\begin{aligned} \text{ge} &: a \rightarrow \text{Maybe } a \rightarrow \text{Bool} \\ \text{ge}(_, \text{Nothing}) &= \text{true} \\ \text{ge}(a, \text{Just } b) &= a \geq b \end{aligned}$$

$$\begin{aligned} \text{nFromProposal} &: \text{Proposal } a \rightarrow \mathbb{N} \\ \text{nFromProposal}(\text{Proposal } n _) &= n \end{aligned}$$

$\text{nFromProposal}(p)$ retrieves the proposal number n inside proposal p , which has the form $\text{Proposal } n \text{ } pr$. $\text{nFromProposal}(p)$, $\text{gt}(a, ma)$, and $\text{ge}(a, ma)$ are used to extract and compare proposal numbers in phase 2b of the acceptor.

$$\text{genA}_Q : \mathbb{N} \times \mathbb{N} \times \mathbb{N} \rightarrow \text{list of } \mathbb{N}$$

$\text{genA}_Q(p, c_A, c_P)$ returns a randomly selected set A_Q with $A_Q \subseteq A = \{1, \dots, c_A\}$ and $|A_Q| > \frac{|A|}{2}$. A_Q consists of any majority of acceptors. In Paxos a majority of acceptors forms a quorum, i.e. an accepting set with which a value can be chosen [10]. We use this function when initiating the proposers to give them a quorum of acceptors with which they communicate.

3.4 Processes

3.4.1 System Initialization

$$\begin{aligned} \text{Sys}(c_A, c_P) &= \bar{o}[2](t) \cdot \text{P}_{\text{init}}^P(c_A + 1, \text{genA}_Q(c_A + c_P, c_A, c_P), c_A + c_P, c_A + c_P, []) \\ &\quad | \bar{o}[1](t) \cdot \Pi_{c_A < k \leq c_A + c_P} \text{P}_{\text{init}}^P(c_A + 1, \text{genA}_Q(k, c_A, c_P), k, k, []) \\ &\quad | \Pi_{1 \leq j \leq c_A} \text{P}_{\text{init}}^A(j, c_A + 1, c_A, c_P, n_a, pr_a) \\ \\ \text{P}_{\text{init}}^P(p, A_Q, n, m, \vec{V}) &= \bar{b}_n[i](s) \cdot \text{P}^P \\ \text{P}_{\text{init}}^A(a, p, c_A, c_P, n, pr) &= \Pi_{c_A < k \leq c_A + c_P} b_k[a](s) \cdot \text{P}_1^A \end{aligned}$$

$\text{Sys}(c_A, c_P)$, $\text{P}_{\text{init}}^P(p, A_Q, n, m, \vec{V})$, and $\text{P}_{\text{init}}^A(a, p, c_A, c_P, n, pr)$ describe the system initialization. c_A and c_P are the number of acceptors and proposers respectively.

An outer session is created through shared-point \bar{o} . This outer session is not strictly necessary but was left in to allow for easier extension of the model. The acceptors are initialized using indices from 1 to c_A and the proposers are initialized using indices from $c_A + 1$ to $c_A + c_P$.

$\text{P}_{\text{init}}^P(p, A_Q, n, m, \vec{V})$ is initialized with the proposer's role in its own session p , which is always $c_A + 1$, a quorum of acceptors A_Q , an index n , and a vector \vec{V} . Each proposer has the same role $p = c_A + 1$ but uses a different shared-point b_n according to its index n . m is initialized to the same value as n but is never updated. \vec{V} is used in the proposer to collect and evaluate the responses from the acceptors. It is always initialized with an empty list $[]$. Shared-point b_n is used to initiate a session. Afterwards, the process behaves like P^P . We assume a mechanism for electing a distinguished proposer, which acts as the leader [11]. The leader is

the only proposer that can try issuing proposers. A new leader is elected via the same mechanism when the previous leader terminates or crashes.

$P_{\text{init}}^A(a, p, c_A, c_P, n, pr)$ is initialized with the acceptor's index a , the proposer index p , which is always $c_A + 1$, c_A , c_P , initial knowledge for the highest promised proposal number n , if available, and initial knowledge for the most recently accepted proposal pr , if available. n is of type Maybe \mathbb{N} and pr is of type Maybe (Proposal Value) thus both can be Nothing. Each of the proposers' session requests are accepted in a separate subprocess. These subprocesses run parallel to each other but still access the same values for n and pr . We observe that each subprocess in an acceptor accesses a different channel s , since it is generated by the proposer and passed through when the proposers' session request is accepted. Afterwards, each subprocess behaves like P_1^A .

3.4.2 Proposer

To define the proposer and the acceptor we introduce a function $\text{update}(n, m)$ which replaces the value inside n with the value of m . We use this function to update the local variables of the processes.

$$\begin{aligned}
P^P &= (\mu X) \text{ update}(n, n + 1) . \\
&\left(\bigodot_{a \in A_Q} s[p, a]!_u l1a \langle \text{proposalNumber}_m(n) \rangle \right) . \\
&\left(\bigodot_{a \in A_Q} s[a, p]?_u l1b \langle \perp \rangle (v_a) \right) . \\
&\text{if anyNack}(\vec{V}) \text{ or } \text{promiseCount}(\vec{V}) < \left\lceil \frac{p}{2} \right\rceil \\
&\text{then } s[p, A_Q]!_w \text{Restart}.X \\
&\text{else} \\
&\quad s[p, A_Q]!_w \text{Accept}. \\
&\left(\bigodot_{a \in A_Q} s[p, a]!_u l2a \langle \text{Proposal } \text{proposalNumber}_m(n) \text{ promiseValue}(\vec{V}) \rangle \right) . 0
\end{aligned}$$

At the start of the recursion n is incremented to make sure every run of the recursion uses a different n and thus a different proposal number. The proposal number is sent to every acceptor in A_Q and their replies are gathered in \vec{V} through v_a . Because $p = c_A + 1$, the minimum number of acceptors needed to form a majority is $\left\lceil \frac{p}{2} \right\rceil = \left\lceil \frac{c_A + 1}{2} \right\rceil$. If any Nack x was received or the number of Promise y received is less than that needed for the smallest majority the proposer restarts the algorithm. Otherwise, the proposer sends its proposal to the acceptors and terminates.

3.4.3 Acceptor

$$\begin{aligned}
P_1^A &= (\mu X) s[p, a]?_u l1a \langle \perp \rangle (n') . \\
&\quad \text{if } n' = \perp \\
&\quad \text{then } s[a, p]!_u l1b \langle \perp \rangle . P_2^A \\
&\quad \text{else} \\
&\quad \quad \text{if } \text{gt}(n', n) \\
&\quad \quad \text{then } \text{update}(n, n') . s[a, p]!_u l1b \langle \text{Promise } pr \rangle . P_2^A \\
&\quad \quad \text{else } s[a, p]!_u l1b \langle \text{Nack } n \rangle . P_2^A \\
\\
P_2^A &= s[p, a]?_w \text{Accept} . s[p, a]?_u l2a \langle \perp \rangle (pr') . \\
&\quad \text{if } pr' = \perp \\
&\quad \text{then } 0 \\
&\quad \text{else} \\
&\quad \quad \text{if } \text{ge}(\text{nFromProposal}(pr'), n) \\
&\quad \quad \text{then } \text{update}(pr, pr') . \text{update}(n, \text{Just } \text{nFromProposal}(pr')) . 0 \\
&\quad \quad \text{else } 0 \\
&\quad \oplus \text{Restart} . X \\
&\quad \oplus \text{Abort} . 0
\end{aligned}$$

For each proposer an acceptor has a corresponding subprocess, which behaves like P_1^A . These subprocesses access the same values for n and pr . This means that updating these values with $\text{update}(n, m)$ updates them for all subprocesses of an acceptor.

Each subprocess can communicate with one proposer. Thus, if that proposer does not or can not communicate with a particular subprocess of an acceptor then there is no need for that subprocess. It is possible that an acceptor participates in a proposer's session but is not contained in the proposer's quorum of acceptors A_Q , in which case the proposer does not communicate with that acceptor. It is also possible for a proposer to crash or otherwise terminate, in which case the proposer can not communicate with that acceptor.

Each subprocess starts out by potentially receiving a proposal number n' from the corresponding proposer. If the acceptor does receive a proposal number n' it responds with either $\text{Promise } pr$ or $\text{Nack } n$, depending on the values of n' and n . If the acceptor does not receive a proposal number then it sends \perp to the proposer. Sending \perp to the proposer is only necessary to maintain the global type. In either case the subprocess moves on to receive the proposers' decision in phase 2a.

Since the proposers' decision broadcast is weakly reliable, there are two cases in which the acceptor receives no decision. The proposer might have terminated or this particular acceptor is not in the proposers' quorum of acceptors A_Q . In either case this particular subprocess of the acceptor is no longer needed, because each subprocess of the acceptor exclusively communicates with one proposer. Thus, the subprocess terminates in the default branch *Abort*.

In the *Restart* branch this particular subprocess of the acceptor restarts the algorithm to match the corresponding proposer.

In the *Accept* branch the acceptor potentially receives a proposal pr' from the corresponding proposer. The acceptor updates n and pr if the proposal number in pr' is greater or equal to n . Then the subprocess terminates. If the acceptor does not receive a proposal or the proposal number of pr' is less than n the subprocess terminates without updating n or pr .

3.5 Failure Patterns

Chandra and Toueg introduce a class of failure detectors $\diamond\mathcal{S}$, which is called *eventually strong* in [5]. Failure detectors in $\diamond\mathcal{S}$ satisfy the following properties: (1) eventually every process that crashes is permanently suspected by every correct process and (2) eventually some correct process is never suspected by any correct process.

In all three phases modeled in the global type it is possible to suspect senders. In phases 1a and 2a, with labels $l1a$ and $l2a$ respectively, the acceptors may suspect some proposers. The proposers may suspect some acceptors in phase 1b with label $l1b$. Accordingly, FP_{uskip} is implemented with a failure detector in $\diamond\mathcal{S}$ for phases 1a, 1b, and 2a.

Similarly, message loss is possible in all phases modeled in the global type. Thus, FP_{ml} is also implemented with a failure detector in $\diamond\mathcal{S}$ with one exception. $FP_{\text{ml}}(s, p, a, l)$ returns true if p is a proposer, a is an acceptor that is not contained in the proposers' quorum of acceptors, and $l = l1b$. Since any proposer only communicates with the acceptors in its quorum, we can discard any messages from acceptors outside it.

For the weakly reliable broadcast in phase 2a, the failure pattern FP_{wskip} returns true for sub-processes of acceptors if, and only if, the corresponding proposer crashed, otherwise terminated, or the corresponding proposer's quorum does not include that particular acceptor.

For Paxos to work a majority of acceptors needs to be alive. That means that the number of failed acceptors f needs to satisfy $n > 2f$ where n is the total number of acceptors, except in one case where there are 2 acceptors. Then, at most one acceptor may crash [10]. For acceptors FP_{crash} returns true if, and only if, at least one more acceptor may crash, i.e. $n > 2(f + 1)$ is satisfied. Let \mathbb{F} be the set of processes permanently suspected by a failure detector in $\diamond\mathcal{S}$. For proposers FP_{crash} returns true if $A_Q \setminus \mathbb{F}$ is not a quorum, i.e. if the set of acceptors in A_Q that are not permanently suspected is not enough to form a majority of acceptors.

In Paxos there is no need to reject outdated messages so FP_{uget} is implemented with a constant true.

3.6 Example

In this section we will study an example run of the model with 3 acceptors and 2 proposers. First, we will take a look at the example scenario. Then we will examine the scenario using reduction rules starting at system initialization.

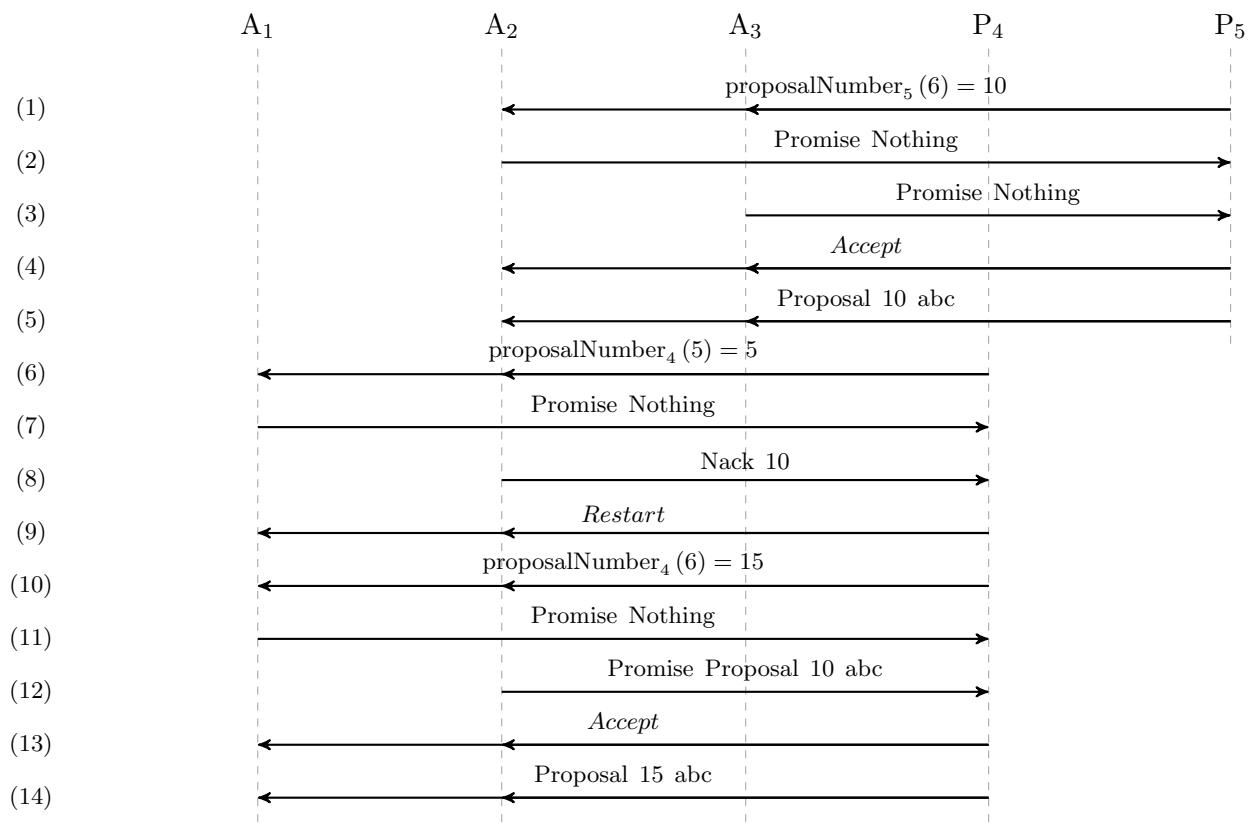


Figure 3.1: Example scenario with 3 acceptors and 2 proposers.

3.6.1 Scenario

Figure 3.1 provides an overview where A_1 , A_2 , and A_3 are the acceptors and P_4 and P_5 are the proposers. P_5 is elected to be the leader. In steps (1) to (5), P_5 completes the Paxos algorithm with A_2 and A_3 and terminates.

At this point A_2 has promised not to accept any proposal numbered less than 10 and has accepted the value abc. So, when P_4 tries to use 5 as its proposal number (6), it receives Nack 10 from A_2 (8) and has to restart the algorithm (9).

P_4 then runs through the Paxos algorithm with A_1 and A_2 starting with a new prepare request (10) with a higher proposal number. In step (12) P_4 learns that value abc with proposal number 10 has already been accepted by A_2 . Later, in step (14), P_4 issues a proposal with the value of the highest-numbered proposal that it receives as a response to its prepare request. In this case there is only one such proposal, which is Proposal 10 abc.

In the end all 3 acceptors have accepted the value abc. A_1 and A_2 have accepted Proposal 15 abc and A_3 has accepted Proposal 10 abc.

3.6.2 Formulae

We set $c_A = 3$, $c_P = 2$, and $V = \{\text{abc}, \text{def}, \dots, \text{vwxyz}\}$.

System Initialization

After inserting c_A and c_P and applying (Init) once for shared-point a we have:

$$\begin{aligned}
 \text{Sys}(c_A, c_P) &= \text{Sys}(3, 2) = \\
 &o[1](t) \cdot \Pi_{3 < k < 5} P_{\text{init}}^P(4, \text{genA}_Q(k, 3, 2), k, k, []) \\
 &\quad | \bar{o}[2](t) \cdot P_{\text{init}}^P(4, \text{genA}_Q(5, 3, 2), 5, 5, []) \\
 &\quad | \Pi_{1 \leq a \leq 3} P_{\text{init}}^A(a, 4, 3, 2, n_a, pr_a) \\
 &\xrightarrow{*} (\nu t) (\bar{b}_4[4](s) \cdot P^P = P_4 \\
 &\quad | \bar{b}_5[4](r) \cdot P^P = P_5 \\
 &\quad | (b_4[1](s) \cdot P_1^A \mid b_5[1](r) \cdot P_1^A) = A_1 \\
 &\quad | (b_4[2](s) \cdot P_1^A \mid b_5[2](r) \cdot P_1^A) = A_2 \\
 &\quad | (b_4[3](s) \cdot P_1^A \mid b_5[3](r) \cdot P_1^A) = A_3 \\
 &\quad | \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
 \end{aligned}$$

Note that the outer session created via shared-point o isn't strictly necessary in the model. We apply (Init) once for shared-point b_4 and once again for shared-point b_5 to obtain:

$$\begin{aligned}
& \mapsto^* (\nu t) (\nu s) (\nu r) (\\
& (\mu X) \text{ update } (n, 5) . \left(\bigodot_{a \in \{1,2\}} s[4, a]!_u l1a \langle \text{proposalNumber}_4(n) \rangle \right) \dots = P_4 \\
& | (\mu X) \text{ update } (n, 6) . \left(\bigodot_{a \in \{2,3\}} r[4, a]!_u l1a \langle \text{proposalNumber}_5(n) \rangle \right) \dots = P_5 \\
& | ((\mu X) s[4, 1]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) r[4, 1]?_u l1a \langle \perp \rangle (n') . \text{if } \dots) = A_1 \\
& | ((\mu X) s[4, 2]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) r[4, 2]?_u l1a \langle \perp \rangle (n') . \text{if } \dots) = A_2 \\
& | ((\mu X) s[4, 3]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) r[4, 3]?_u l1a \langle \perp \rangle (n') . \text{if } \dots) = A_3 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [] \rangle
\end{aligned}$$

Note that each process is shortened to only show the next few steps instead of the entire process.

The Happy Path

After applying the updates in P_4 and P_5 the first inter-process communication can take place. In this case P_5 communicates with A_2 and A_3 . We apply (USend) and (UGet) twice to send $\text{proposalNumber}_5(6) = 10$ to A_2 and A_3 .

$$\begin{aligned}
& \mapsto^* (\nu t) (\nu s) (\nu r) (\\
& (\mu X) s[4, 1]!_u l1a \langle \text{proposalNumber}_4(5) \rangle . s[4, 2]!_u l1a \langle \text{proposalNumber}_4(5) \rangle \dots = P_4 \\
& | (\mu X) r[2, 4]?_u l1b \langle \perp \rangle (v_2) . r[3, 4]?_u l1b \langle \perp \rangle (v_3) \dots = P_5 \\
& | ((\mu X) s[4, 1]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) r[4, 1]?_u l1a \langle \perp \rangle (n') . \text{if } \dots) = A_1 \\
& | ((\mu X) s[4, 2]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) \text{if } 10 = \perp \text{ then } s[2, 4]!_u l1b \langle \perp \rangle . P_2^A \text{ else } \dots) = A_2 \\
& | ((\mu X) s[4, 3]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) \text{if } 10 = \perp \text{ then } s[3, 4]!_u l1b \langle \perp \rangle . P_2^A \text{ else } \dots) = A_3 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [] \rangle
\end{aligned}$$

Since $10 \neq \perp$ both A_2 and A_3 move into their respective else branches.

$$\begin{aligned}
& = (\nu t) (\nu s) (\nu r) (\\
& (\mu X) s[4, 1]!_u l1a \langle \text{proposalNumber}_4(5) \rangle . s[4, 2]!_u l1a \langle \text{proposalNumber}_4(5) \rangle \dots = P_4 \\
& | (\mu X) r[2, 4]?_u l1b \langle \perp \rangle (v_2) . r[3, 4]?_u l1b \langle \perp \rangle (v_3) \dots = P_5 \\
& | ((\mu X) s[4, 1]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) r[4, 1]?_u l1a \langle \perp \rangle (n') . \text{if } \dots) = A_1 \\
& | ((\mu X) s[4, 2]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) \text{if } \text{gt}(10, \text{Nothing}) \text{ then } \dots \text{ else } \dots) = A_2 \\
& | ((\mu X) s[4, 3]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) \text{if } \text{gt}(10, \text{Nothing}) \text{ then } \dots \text{ else } \dots) = A_3 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [] \rangle
\end{aligned}$$

Because $\text{gt}(10, \text{Nothing})$ returns true, A_2 and A_3 move into their respective then branches. After executing $\text{update}(n, 10)$, A_2 and A_3 are ready to send their responses to P_5 .

$$\begin{aligned}
&= (\nu t) (\nu s) (\nu r) (\\
&(\mu X) s [4, 1]!_u l1a \langle \text{proposalNumber}_4 (5) \rangle . s [4, 2]!_u l1a \langle \text{proposalNumber}_4 (5) \rangle \dots = P_4 \\
&| (\mu X) r [2, 4]?_u l1b \langle \perp \rangle (v_2) . r [3, 4]?_u l1b \langle \perp \rangle (v_3) \dots = P_5 \\
&| ((\mu X) s [4, 1]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) r [4, 1]?_u l1a \langle \perp \rangle (n') . \text{if } \dots) = A_1 \\
&| ((\mu X) s [4, 2]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) r [2, 4]?_u l1b \langle \text{Promise Nothing} \rangle . P_2^A) = A_2 \\
&| ((\mu X) s [4, 3]?_u l1a \langle \perp \rangle (n') . \text{if } \dots | (\mu X) r [3, 4]?_u l1b \langle \text{Promise Nothing} \rangle . P_2^A) = A_3 \\
&| \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

We apply (USend) and (UGet) twice to do just that. Note that we also apply (USkip) to A_1 , evaluate its branches, and apply (USend) to A_1 and then (ML) to P_5 to discard the dummy message. All three acceptors move into P_2^A .

$$\begin{aligned}
&\mapsto^* (\nu t) (\nu s) (\nu r) (\\
&(\mu X) s [4, 1]!_u l1a \langle \text{proposalNumber}_4 (5) \rangle . s [4, 2]!_u l1a \langle \text{proposalNumber}_4 (5) \rangle \dots = P_4 \\
&| (\mu X) r [4, \{2, 3\}]!_w \text{Accept} . r [4, 2]!_u l2a \langle \text{Proposal 10 abc} \rangle \dots = P_5 \\
&| ((\mu X) \dots | (\mu X) r [4, 1]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0) = A_1 \\
&| ((\mu X) \dots | (\mu X) r [4, 2]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0) = A_2 \\
&| ((\mu X) \dots | (\mu X) r [4, 3]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0) = A_3 \\
&| \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

P_5 broadcasts its decision *Accept* to A_2 and A_3 . By applying (WSel) once, (WBran) twice we obtain:

$$\begin{aligned}
&\mapsto^* (\nu t) (\nu s) (\nu r) (\\
&(\mu X) s [4, 1]!_u l1a \langle \text{proposalNumber}_4 (5) \rangle . s [4, 2]!_u l1a \langle \text{proposalNumber}_4 (5) \rangle \dots = P_4 \\
&| (\mu X) r [4, 2]!_u l2a \langle \text{Proposal 10 abc} \rangle . r [4, 3]!_u l2a \langle \text{Proposal 10 abc} \rangle \dots = P_5 \\
&| ((\mu X) \dots | (\mu X) r [4, 1]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0) = A_1 \\
&| ((\mu X) \dots | (\mu X) r [4, 2]?_u l2a \langle \perp \rangle (pr') . \text{if } \dots) = A_2 \\
&| ((\mu X) \dots | (\mu X) r [4, 3]?_u l2a \langle \perp \rangle (pr') . \text{if } \dots) = A_3 \\
&| \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

Now P_5 can send its proposal to A_2 and A_3 and terminate. P_4 will be the new leader. To do so we apply (USend) and (UGet) twice. A_2 and A_3 accept the proposal and the respective subprocesses terminate. Note that we apply (WSkip) in A_1 and terminate that subprocess as well.

$$\begin{aligned}
& \mapsto^* (\nu t) (\nu s) (\nu r) (\\
& (\mu X) s [4, 1]!_u l1a \langle \text{proposalNumber}_4 (5) \rangle . s [4, 2]!_u l1a \langle \text{proposalNumber}_4 (5) \rangle \dots = P_4 \\
& | (\mu X) s [4, 1]?_u l1a \langle \perp \rangle (n') . \text{if } \dots = A_1 \\
& | (\mu X) s [4, 2]?_u l1a \langle \perp \rangle (n') . \text{if } \dots = A_2 \\
& | (\mu X) s [4, 3]?_u l1a \langle \perp \rangle (n') . \text{if } \dots = A_3 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

At this point the local variables of A_2 and A_3 are $n = 10$ and $pr = \text{Proposal } 10 \text{ abc}$. A_1 has not updated its local variables $n = \text{Nothing}$ and $pr = \text{Nothing}$.

Restarting the Algorithm

Next, P_4 sends prepare requests with a proposal number less than 10, which is rejected by A_2 . P_4 then decides to restart the algorithm. We apply (USend) and (UGet) twice. We also apply (USkip) once in A_3 .

$$\begin{aligned}
& \mapsto^* (\nu t) (\nu s) (\nu r) (\\
& (\mu X) s [1, 4]?_u l1b \langle \perp \rangle (v_1) . s [2, 4]?_u l1b \langle \perp \rangle (v_2) \dots = P_4 \\
& | (\mu X) \text{if } 5 = \perp \text{ then } P_2^A \text{ else } \dots = A_1 \\
& | (\mu X) \text{if } 5 = \perp \text{ then } P_2^A \text{ else } \dots = A_2 \\
& | (\mu X) \text{if } \perp = \perp \text{ then } P_2^A \text{ else } \dots = A_3 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

A_3 moves directly to P_2^A whereas A_1 and A_2 send their responses to P_4 before moving to P_2^A . A_1 also updates its local variable $n = 5$.

$$\begin{aligned}
& = (\nu t) (\nu s) (\nu r) (\\
& (\mu X) s [1, 4]?_u l1b \langle \perp \rangle (v_1) . s [2, 4]?_u l1b \langle \perp \rangle (v_2) \dots = P_4 \\
& | (\mu X) s [1, 4]!_u l1b \langle \text{Promise Nothing} \rangle . P_2^A = A_1 \\
& | (\mu X) s [2, 4]!_u l1b \langle \text{Nack } 10 \rangle . P_2^A = A_2 \\
& | (\mu X) r [4, 3]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0 = A_3 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

Applying (USend) and (UGet) twice and evaluating the branching in P_4 yields:

$$\begin{aligned}
& \mapsto^* (\nu t) (\nu s) (\nu r) (\\
& (\mu X) s [4, \{1, 2\}]!_w \text{Restart}.X = P_4 \\
& | (\mu X) s [4, 1]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0 = A_1 \\
& | (\mu X) s [4, 1]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0 = A_2 \\
& | (\mu X) r [4, 3]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0 = A_3 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

P_4 sends its decision to restart the algorithm to A_1 and A_2 by applying (WSel) once and (WBran) twice. A_3 terminates after applying (WSkip).

$$\begin{aligned}
& \mapsto^* (\nu t) (\nu s) (\nu r) (\\
& (\mu X) s [4, 1]!_u l1a \langle 15 \rangle . s [4, 2]!_u l1a \langle 15 \rangle \dots = P_4 \\
& | (\mu X) s [4, 1]?_u l1a \langle \perp \rangle (n') . \text{if} \dots = A_1 \\
& | (\mu X) s [4, 2]?_u l1a \langle \perp \rangle (n') . \text{if} \dots = A_2 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

The Happy Path, Again

This time P_4 uses a high enough proposal number so that A_1 and A_2 both promise not to accept any proposal numbered less than that. By applying (USend) and (UGet) and evaluating the branches in the remaining acceptors we arrive at:

$$\begin{aligned}
& \mapsto^* (\nu t) (\nu s) (\nu r) (\\
& (\mu X) s [1, 4]?_u l1b \langle \perp \rangle (v_1) . s [2, 4]?_u l1b \langle \perp \rangle (v_2) . \text{if} \dots = P_4 \\
& | (\mu X) s [1, 4]!_u l1b \langle \text{Promise Nothing} \rangle . P_2^A = A_1 \\
& | (\mu X) s [2, 4]!_u l1b \langle \text{Promise Proposal 10 } abc \rangle . P_2^A = A_2 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

Note that, at this point, A_1 and A_2 have updated their respective n to 15.

Because A_2 has already accepted a proposal, it responds to P_4 's prepare request with that proposal. Twice more we apply (USend) and (UGet) and evaluate the branch in P_4 to obtain:

$$\begin{aligned}
& \mapsto^* (\nu t) (\nu s) (\nu r) (\\
& (\mu X) s [4, \{1, 2\}]!_w \text{Accept} \dots = P_4 \\
& | (\mu X) s [4, 1]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0 = A_1 \\
& | (\mu X) s [4, 1]?_w \text{Accept} \dots \oplus \text{Restart}.X \oplus \text{Abort}.0 = A_2 \\
& | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])
\end{aligned}$$

P_4 has received enough promises to send its own proposal. The value for that proposal is `abc` because that is the value of the highest-numbered proposal P_4 received as a response to its prepare request. First, we apply (WSel) and (WBran).

$$\begin{aligned} & \mapsto^* (\nu t) (\nu s) (\nu r) (\\ & (\mu X) s [4, 1]!_u l2a \langle \text{Proposal } 15 \text{ abc} \rangle . s [4, 2]!_u l2a \langle \text{Proposal } 15 \text{ abc} \rangle . 0 = P_4 \\ & | (\mu X) s [4, 1]?_u l2a \langle \perp \rangle (pr') . \text{if } \dots = A_1 \\ & | (\mu X) s [4, 2]?_u l2a \langle \perp \rangle (pr') . \text{if } \dots = A_2 \\ & | \Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : []) \end{aligned}$$

Then we apply (USend) and (UGet) to send the proposal from P_4 to the acceptors. P_4 terminates and the acceptors accept the received proposal and then terminate as well.

$$\mapsto^* (\nu t) (\nu s) (\nu r) (\Pi_{1 \leq k, l \leq 4, k \neq l} s_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 4, k \neq l} r_{k \rightarrow l} : [] \mid \Pi_{1 \leq k, l \leq 2, k \neq l} t_{k \rightarrow l} : [])$$

Afterwards A_1 and A_2 have $n = 15$ and $pr = \text{Proposal } 15 \text{ abc}$ and A_3 has $n = 10$ and $pr = \text{Proposal } 10 \text{ abc}$. All acceptors have accepted the value `abc`.

4 Analysis

We take the model from the previous chapter, type-check it, and discuss what the type check means for agreement, validity, and termination of the Paxos algorithm. To execute the type check we project the global type to local types and use the typing rules given in [13] to prove that our model is well-typed.

4.1 Local Types

Because no communication takes place in the outer session, the outer session's type is $G = 0$. Every projection of G to a local type is $G \upharpoonright_k = 0$ for every k .

For $1 \leq a \leq c_A$ and $c_A + 1 \leq p \leq c_A + c_P$ we define the projections of the global type G_{p,A_Q} .

$$G_{p,A_Q} \upharpoonright_p = (\mu x) \left(\bigodot_{a \in A_Q} [a]!_u l1a \langle \mathbb{N} \rangle \right) \cdot \left(\bigodot_{a \in A_Q} [a]?_u l1b \langle \text{Promise Value} \rangle \right) \cdot \\ \left([A_Q]!_w \text{Accept}. \left(\bigodot_{a \in A_Q} [a]!_u l2a \langle \text{Proposal Value} \rangle \right) . 0 \oplus \text{Restart}.x \oplus \text{Abort}.0 \right)$$

$G_{p,A_Q} \upharpoonright_p$ defines the local type for proposers. First, the proposer sends a proposal number to all acceptors in its quorum in phase 1a. It receives their responses in phase 1b and then branches in phase 2a. We can see that the proposer communicates with all acceptors in its quorum in every phase.

$$G_{p,A_Q} \upharpoonright_a = (\mu x) [p]?_u l1a \langle \mathbb{N} \rangle \cdot [p]!_u l1b \langle \text{Promise Value} \rangle \cdot \\ ([p]?_w \text{Accept}. [p]?_u l2a \langle \text{Proposal Value} \rangle . 0 \oplus \text{Restart}.x \oplus \text{Abort}.0)$$

$G_{p,A_Q} \upharpoonright_a$ defines the local type for acceptors, assuming a proposer p . Since Paxos defines two roles that communicate with each other, their local types complement each other. An acceptor first receives a proposal number, then it responds with a Promise Value. Finally, it receives the proposer's branching choice.

4.2 Type Check

$$\Gamma = o : G \cdot b_{c_A+1} : G_{p,A_Q} \cdot b_{c_A+2} : G_{p,A_Q} \cdot \dots \cdot b_{c_A+c_P} : G_{p,A_Q} \cdot c_A : \mathbb{N} \cdot c_P : \mathbb{N}$$

Γ contains the type for our shared-points o and b_n where $c_A + 1 \leq n \leq c_A + c_P$.

We start the type-check with the global environment Γ and the entry-point of the model $\text{Sys}(c_A, c_P)$. Then, we apply the typing rules in [13] in a proof tree and show that the model can be derived from the axioms (Var) and (End).

4.2.1 System Initialization

$$\frac{\frac{(S_1)}{\Gamma \vdash \bar{o}[2](t) \dots \triangleright \emptyset} \quad \frac{\frac{(S_2)}{\Gamma \vdash o[1](t) \dots \triangleright \emptyset} \quad \frac{(S_3)}{\Gamma \vdash \Pi_{1 \leq a \leq c_A} P_{\text{init}}^A(\dots) \triangleright \emptyset}}{\Gamma \vdash o[1](t) \dots \mid \Pi_{1 \leq a \leq c_A} P_{\text{init}}^A(\dots) \triangleright \emptyset} (\text{Par})}{\Gamma \vdash \bar{o}[2](t) \cdot P_{\text{init}}^P(\dots) \mid o[1](t) \dots \mid \Pi_{1 \leq a \leq c_A} P_{\text{init}}^A(\dots) \triangleright \emptyset} (\text{Par})$$

We apply (Par) twice and split off into three sub-proofs (S_1) , (S_2) , and (S_3) .

$$(S_1) = \frac{\frac{(P)}{\Gamma \vdash \bar{b}_{c_A+c_P}[c_A+1](s) \cdot P^P \triangleright t[2] : G \mid_2}}{\Gamma \vdash \bar{o}[2](t) \cdot P_{\text{init}}^P(c_A+1, \text{genA}_Q(c_A+c_P, c_A, c_P), c_A+c_P, c_A+c_P, []) \triangleright \emptyset} (\text{Rec})$$

In (S_1) we apply (Rec) once and defer the rest of the proof-tree. Note that, since $G \mid_2 = 0$, $t[2] : G \mid_2 = \emptyset$. This is relevant later when continuing (P) .

$$(S_2) = \frac{\frac{\frac{(P)}{\Gamma \vdash \bar{b}_{c_A+1}[c_A+1](s) \cdot P^P \triangleright \emptyset} \quad \dots \quad \frac{(P)}{\Gamma \vdash \bar{b}_{c_A+c_P-1}[c_A+1](s) \cdot P^P \triangleright \emptyset}}{\Gamma \vdash \Pi_{c_A < k < c_A+c_P} P_{\text{init}}^P(c_A+1, \text{genA}_Q(k, c_A, c_P), k, k, []) \triangleright t[1] : G \mid_1} (\text{Par})^{c_P-1}}{\Gamma \vdash o[1](t) \cdot \Pi_{c_A < k < c_A+c_P} P_{\text{init}}^P(\dots) \triangleright \emptyset} (\text{Acc})$$

Applying (Acc) in (S_2) requires that $o : G \in \Gamma$. (Par) is applied $c_P - 1$ times to separate all the proposer processes. Each individual proposer can be type-checked with the same proof-tree (P) . Because $G \mid_1 = 0$, $t[1] : G \mid_1 = \emptyset$. The session environment Δ in (P) is empty for every proposer.

$$(S_3) = \frac{\frac{\frac{(A_1)}{\Gamma \vdash P_1^A \triangleright s[a] : G_{p,A_Q} \vdash_a} \quad (\text{Acc})}{\Gamma \vdash b_k[a](s) \cdot P_1^A \triangleright \emptyset} \quad \dots \quad (\text{Par})^{c_P}}{\Gamma \vdash \Pi_{c_A < k \leq c_A + c_P} b_k[c_A](s) \cdot P_1^A \triangleright \emptyset} \quad \dots \quad (\text{Par})^{c_A}$$

(Par) is applied c_A times to separate the individual acceptors and then c_P times for each acceptor to separate the individual subprocesses. Since every subprocess of every acceptor behaves like P_1^A and has the same local type, the same proof-tree (A_1) can be applied. Applying (Acc) to every subprocess of every acceptor requires $\forall k \in \mathbb{N} : (c_A + 1 \leq k \wedge k \leq c_A + c_P) \rightarrow b_k : G_{p,A_Q} \in \Gamma$.

Note that only one acceptor and one of its subprocesses is shown in (S_3) . The rest has been left out to improve readability.

4.2.2 Proposer

Let $p = c_A + 1, A_Q = \text{genA}_Q(k, c_A, c_P), n = k, m = k, \vec{V} = []$ where $c_A < k \leq c_A + c_P$. This gives us the values for the arguments of P_{init}^P . We observe that $\Gamma \Vdash p : \mathbb{N}, \Gamma \Vdash k : \mathbb{N}, \Gamma \Vdash n : \mathbb{N}, \Gamma \Vdash m : \mathbb{N}$, and $\Gamma \Vdash A_Q : \text{list of } \mathbb{N}$. p, k, n , and m are natural numbers and A_Q is a list of natural numbers under global environment Γ .

To abbreviate the proposer's local type in the following proof-trees we define the following sub-formulae.

$$\begin{aligned} T_{\text{acc}}^P &= \left(\odot_{a \in A_Q} [a]!_u l2a \langle \text{Proposal Value} \rangle \right) . 0 \\ T_{\text{branch}}^P &= ([A_Q]!_w \text{Accept} . T_{\text{acc}}^P \oplus \text{Restart} . x \oplus \text{Abort} . 0) \end{aligned}$$

Note that $G_{p,A_Q} \vdash_p = (\mu x) \left(\odot_{a \in A_Q} [a]!_u l1a \langle \mathbb{N} \rangle \right) . \left(\odot_{a \in A_Q} [a]?_u l1b \langle \text{Promise Value} \rangle \right) . T_{\text{branch}}^P$.

In order to shorten the proposer's process we define some variables.

$$\begin{aligned} e &= \text{anyNack} \left(\vec{V} \right) \text{ or } \text{promiseCount} \left(\vec{V} \right) < \left\lceil \frac{p}{2} \right\rceil \\ pn &= \text{proposalNumber}_m(n) \\ prop &= \text{Proposal } \text{proposalNumber}_m(n) \text{ promiseValue} \left(\vec{V} \right) \end{aligned}$$

The actual values of e , pn , and $prop$ are not relevant for the type check. We observe that $\Gamma \Vdash e : \text{Bool}$, $\Gamma \Vdash pn : \mathbb{N}$, and $\Gamma \Vdash prop : \text{Proposal Value}$.

To further abbreviate the terms in the proof-trees we define two global environments Γ' and Γ'' .

$$\Gamma' = \Gamma \cdot X : x$$

Γ' contains Γ and a type for the recursion variable X .

$$\Gamma'' = \Gamma' \cdot v_a : \text{Promise Value}, \forall a \in A_Q$$

Γ'' contains Γ' and types for the entries of \vec{V} . These are added to the global environment when applying (UGet) in phase 1b.

$$\begin{array}{c}
\frac{(P_t)}{\Gamma'' \vdash s[p, A_Q]!_w \text{Restart}.X \triangleright s[p] : T_{\text{branch}}^P} \quad \frac{(P_f)}{\Gamma'' \vdash s[p, A_Q]!_w \text{Accept} \dots \triangleright s[p] : T_{\text{branch}}^P} \quad (\text{If}) \\
\frac{\Gamma'' \vdash \text{if } e \text{ then } s[p, A_Q]!_w \text{Restart}.X \text{ else } s[p, A_Q]!_w \text{Accept} \dots \triangleright s[p] : T_{\text{branch}}^P}{\Gamma' \vdash \left(\bigodot_{a \in A_Q} s[a, p]?_u l1b \langle \perp \rangle (v_a) \right) \dots \triangleright s[p] : \left(\bigodot_{a \in A_Q} [a]?_u l1b \langle \text{Promise Value} \rangle \right)} \quad (\text{UGet})^{|A_Q|} \\
\frac{\Gamma' \vdash \left(\bigodot_{a \in A_Q} s[p, a]?_u l1a \langle pn \rangle \right) \dots \triangleright s[p] : \left(\bigodot_{a \in A_Q} [a]?_u l1a \langle \mathbb{N} \rangle \right) \dots}{\Gamma' \vdash \text{update}(n, n+1) \dots \triangleright s[p] : \left(\bigodot_{a \in A_Q} [a]?_u l1a \langle \mathbb{N} \rangle \right) \dots} \quad (\text{USend})^{|A_Q|} \\
\frac{\Gamma' \vdash \text{update}(n, n+1) \dots \triangleright s[p] : \left(\bigodot_{a \in A_Q} [a]?_u l1a \langle \mathbb{N} \rangle \right) \dots}{\Gamma' \vdash (\mu X) \text{update}(n, n+1) \dots \triangleright s[p] : (\mu x) \left(\bigodot_{a \in A_Q} [a]?_u l1a \langle \mathbb{N} \rangle \right) \dots} \quad (??) \\
\frac{\Gamma' \vdash (\mu X) \text{update}(n, n+1) \dots \triangleright s[p] : (\mu x) \left(\bigodot_{a \in A_Q} [a]?_u l1a \langle \mathbb{N} \rangle \right) \dots}{\Gamma \vdash (\mu X) \text{update}(n, n+1) \dots \triangleright s[p] : (\mu x) \left(\bigodot_{a \in A_Q} [a]?_u l1a \langle \mathbb{N} \rangle \right) \dots} \quad (\text{Rec}) \\
(P) = \frac{\Gamma \vdash (\mu X) \text{update}(n, n+1) \dots \triangleright s[p] : (\mu x) \left(\bigodot_{a \in A_Q} [a]?_u l1a \langle \mathbb{N} \rangle \right) \dots}{\Gamma \vdash \bar{b}_n[p](s) \cdot P^P \triangleright \emptyset} \quad (\text{Req})
\end{array}$$

(P) is the continuation of (S_1) and (S_2). In both proof-trees the session environment Δ was empty. Here, we apply (Req) and add $s[p] : G_{p, A_Q} \vdash_p$ to the session environment. Applying (Rec) changes the global environment from Γ to Γ' . (??) only changes the process and lets us continue. First (USend) and then (UGet) is applied for every acceptor in A_Q . (UGet) expands the session environment to Γ'' . (If) splits the proof-tree into (P_t) and (P_f).

$$(P_t) = \frac{\frac{\Gamma'' \vdash X \triangleright s[p] : x}{\Gamma'' \vdash s[p, A_Q]!_w \text{Restart}.X \triangleright s[p] : ([A_Q]!_w \text{Accept}. T_{\text{acc}}^P \oplus \text{Restart}.x \oplus \text{Abort}.0)} \quad (\text{Var})}{\Gamma'' \vdash s[p, A_Q]!_w \text{Restart}.X \triangleright s[p] : ([A_Q]!_w \text{Accept}. T_{\text{acc}}^P \oplus \text{Restart}.x \oplus \text{Abort}.0)} \quad (\text{WSel})$$

We apply (WSel) and then (Var) to finish (P_t).

$$(P_f) = \frac{\frac{\frac{\Gamma'' \vdash 0 \triangleright s[p] : 0}{\Gamma'' \vdash \left(\bigodot_{a \in A_Q} s[p, a]?_u l2a \langle prop \rangle \right) . 0 \triangleright s[p] : \left(\bigodot_{a \in A_Q} [a]?_u l2a \langle \text{Proposal Value} \rangle \right) . 0} \quad (\text{USend})^{|A_Q|}}{\Gamma'' \vdash s[p, A_Q]!_w \text{Accept} \dots \triangleright s[p] : ([A_Q]!_w \text{Accept}. T_{\text{acc}}^P \oplus \text{Restart}.x \oplus \text{Abort}.0)} \quad (\text{WSel})$$

After applying (USend) we can apply (USend) once for every acceptor in A_Q . Finally, we can finish (P_f) — and with it (P) — by applying (End).

4.2.3 Acceptor

First, we define the arguments of P_{init}^A and P_1^A . Let $a = j$ and $p = c_A + 1$ where $1 \leq j \leq c_A$. With session environment Γ we have $\Gamma \vdash a : \mathbb{N}$ and $\Gamma \vdash p : \mathbb{N}$.

To improve readability of the proof-trees we break down the acceptor's process and local type.

$$\begin{aligned} P_{\text{acc}}^A &= s[p, a]?_u l2a \langle \perp \rangle (pr') . \text{if } pr' = \perp \\ &\quad \text{then } 0 \\ &\quad \text{else if } \text{ge}(\text{nFromProposal}(pr'), n) \\ &\quad \quad \text{then } \text{update}(pr, pr') . \text{update}(n, \text{Just nFromProposal}(pr')) . 0 \\ &\quad \text{else } 0 \end{aligned}$$

We can see that P_2^A contains P_{acc}^A as $P_2^A = s[p, a]?_w \text{Accept}. P_{\text{acc}}^A \oplus \text{Restart}.X \oplus \text{Abort}.0$.

$$P_t^A = \text{update}(n, n') . s[a, p]!_u l1b \langle \text{Promise } pr \rangle . P_2^A$$

$$P_f^A = s[a, p]!_u l1b \langle \text{Nack } n \rangle . P_2^A$$

$$P_{\text{gt}}^A = \text{if } \text{gt}(n', n) \text{ then } P_t^A \text{ else } P_f^A$$

With P_{gt}^A , P_1^A can be written as $P_1^A = (\mu X) s[p, a]?_u l1a \langle \perp \rangle (n') . \text{if } n' = \perp \text{ then } s[a, p]!_u l1b \langle \perp \rangle . P_2^A \text{ else } P_{\text{gt}}^A$.

$$T_{\text{acc}}^A = [p]?_u l2a \langle \text{Proposal Value} \rangle . 0$$

$$T_{\text{branch}}^A = ([p]?_w \text{Accept}. T_{\text{acc}}^A \oplus \text{Restart}.x \oplus \text{Abort}.0)$$

$$T_{1b}^A = [p]!_u l1b \langle \text{Promise Value} \rangle . T_{\text{branch}}^A$$

The acceptor's local type $G_{p, A_Q} \vdash_a$ can be written as $G_{p, A_Q} \vdash_a = (\mu x) [p]?_u l1a \langle \mathbb{N} \rangle . T_{1b}^A$.

Finally, we define the global environments Γ' , Γ'' , and Γ''' .

$$\begin{aligned} \Gamma' &= \Gamma \cdot X : x \\ \Gamma'' &= \Gamma' \cdot n' : \mathbb{N} \\ \Gamma''' &= \Gamma'' \cdot pr' : \text{Proposal Value} \end{aligned}$$

Γ' contains Γ and assigns the type x to X . Γ'' additionally maps n' to type \mathbb{N} . Γ''' adds type Proposal Value for pr' .

$$\begin{aligned}
& \frac{(A_2)}{\frac{\Gamma'' \vdash P_2^A \triangleright s[a] : T_{\text{branch}}^A}{\Gamma'' \vdash s[a, p]!_u l1b \langle \perp \rangle . P_2^A \triangleright s[a] : T_{1b}^A} \text{ (USend)} \quad \frac{(A_{gt})}{\Gamma'' \vdash P_{gt}^A \triangleright s[a] : T_{1b}^A} \\
& \frac{\Gamma'' \vdash \text{if } n' = \perp \text{ then } s[a, p]!_u l1b \langle \perp \rangle . P_2^A \text{ else if } \text{gt}(n', n) \dots \triangleright s[a] : T_{1b}^A}{\Gamma'' \vdash \text{if } n' = \perp \text{ then } s[a, p]!_u l1a \langle \perp \rangle (n') . \text{if } n' = \perp \dots \triangleright s[a] : [p]?_u l1a \langle \mathbb{N} \rangle . T_{1b}^A} \text{ (UGet)} \\
(A_1) = & \frac{\Gamma' \vdash s[p, a]?_u l1a \langle \perp \rangle (n') . \text{if } n' = \perp \dots \triangleright s[a] : [p]?_u l1a \langle \mathbb{N} \rangle . T_{1b}^A}{\Gamma \vdash (\mu X) s[p, a]?_u l1a \langle \perp \rangle (n') \dots \triangleright s[a] : (\mu x) [p]?_u l1a \langle \mathbb{N} \rangle \dots} \text{ (Rec)}
\end{aligned}$$

After applying (Acc) in (S_3) the session environment contains the acceptor's local type $G_{p, A_Q} \upharpoonright_a$. We apply (Rec) and (UGet) and then split the proof-tree with (If). By applying (Rec) and (UGet) the global environment expands from Γ to Γ' to Γ'' . On the left branch we apply (USend) and defer to (A_2) . The right branch is deferred to (A_{gt}) .

Since the process of the right branch contains an if-then-else and unreliable-send statements before continuing to P_2^A , we will examine this branch first. Much like the left branch, the proof-tree of the right branch can later be deferred to (A_2) .

$$(A_{gt}) = \frac{\frac{(A_t)}{\Gamma'' \vdash P_t^A \triangleright s[a] : T_{1b}^A} \quad \frac{(A_f)}{\Gamma'' \vdash P_f^A \triangleright s[a] : T_{1b}^A}}{\Gamma'' \vdash \text{if } \text{gt}(n', n) \text{ then } P_t^A \text{ else } P_f^A \triangleright s[a] : T_{1b}^A} \text{ (If)}$$

First, we split the proof-tree with (If). We defer the resulting branches to separate proof-trees (A_t) and (A_f) .

$$(A_t) = \frac{\frac{(A_2)}{\Gamma'' \vdash P_2^A \triangleright s[a] : T_{\text{branch}}^A}}{\Gamma'' \vdash s[a, p]!_u l1b \langle \text{Nack } n \rangle . P_2^A \triangleright s[a] : [p]!_u l1b \langle \text{Promise Value} \rangle . T_{\text{branch}}^A} \text{ (USend)}$$

$$(A_f) = \frac{\frac{(A_2)}{\Gamma'' \vdash P_2^A \triangleright s[a] : T_{\text{branch}}^A}}{\Gamma'' \vdash s[a, p]!_u l1b \langle \text{Nack } n \rangle . P_2^A \triangleright s[a] : [p]!_u l1b \langle \text{Promise Value} \rangle . T_{\text{branch}}^A} \text{ (USend)}$$

In both, (A_t) and (A_f) , we apply (USend). Now we can defer to (A_2) , which is the proof-tree for P_2^A .

$$(A_2) = \frac{\frac{(A_{\text{Accept}})}{\Gamma'' \vdash P_{\text{acc}}^A \triangleright s[a] : T_{\text{acc}}^A} \quad \frac{\Gamma'' \vdash X \triangleright s[a] : x}{\Gamma'' \vdash X \triangleright s[a] : x} \text{ (Var)} \quad \frac{\Gamma'' \vdash 0 \triangleright s[a] : 0}{\Gamma'' \vdash 0 \triangleright s[a] : 0} \text{ (End)}}{\Gamma'' \vdash P_2^A \triangleright s[a] : T_{\text{branch}}^A} \text{ (WBranch)}$$

By applying (WBran) we separate the three branches. From left to right we get an *Accept*-, a *Restart*-, and an *Abort*-branch. We defer the *Accept*-branch to (A_{Accept}) . The *Restart*-branch can be finished by applying (Var) and the *Abort*-branch by applying (End).

$$(A_{Accept}) = \frac{\frac{\Gamma''' \vdash 0 \triangleright s[a] : 0}{\Gamma''' \vdash 0 \triangleright s[a] : 0} \text{ (End)} \quad \frac{\frac{\frac{\Gamma''' \vdash \text{update}(pr, pr') \dots \triangleright s[a] : 0}{\Gamma''' \vdash \text{if } \text{ge}(\text{nFromProposal}(pr'), n) \text{ then } \dots \text{else } 0 \triangleright s[a] : 0} \text{ (If)} \quad \frac{\Gamma''' \vdash 0 \triangleright s[a] : 0}{\Gamma''' \vdash 0 \triangleright s[a] : 0} \text{ (End)}}{\Gamma''' \vdash \text{if } pr' = \perp \text{ then } 0 \text{ else } \dots \triangleright s[a] : 0} \text{ (If)} \\ \frac{\Gamma'' \vdash s[p, a]?_u l2a \langle \perp \rangle (pr') \dots \triangleright s[a] : [p]?_u l2a \langle \text{Proposal Value} \rangle .0}{\Gamma'' \vdash s[p, a]?_u l2a \langle \perp \rangle (pr') \dots \triangleright s[a] : [p]?_u l2a \langle \text{Proposal Value} \rangle .0} \text{ (UGet)}$$

We apply (UGet) and expand the global session to Γ''' . The proof-tree is split twice by applying (If) twice. The right-most and left-most proof-trees are finished by applying (End). We defer the proof in the middle to keep (A_{Accept}) readable.

$$(A_{update}) = \frac{\frac{\Gamma''' \vdash \text{update}(n, \text{Just } \text{nFromProposal}(pr')) .0 \triangleright s[a] : 0}{\Gamma''' \vdash \text{update}(pr, pr') \dots \triangleright s[a] : 0} \text{ (??)} \quad \frac{\Gamma''' \vdash 0 \triangleright s[a] : 0}{\Gamma''' \vdash 0 \triangleright s[a] : 0} \text{ (End)}}{\Gamma''' \vdash \text{update}(pr, pr') \dots \triangleright s[a] : 0} \text{ (??)}$$

Finally, we apply (??) twice and (End) once. This concludes the type check and proves that the model is well-typed.

4.3 Termination, Agreement, Validity

4.3.1 Termination

The global type and well-typedness ensure the absence of deadlocks. This means that the processes either loop forever or terminate. Acceptors terminate if all of their sub-processes terminate. Each sub-process of an acceptor corresponds to one proposer. A sub-process can only terminate via the weakly reliable broadcast in P_2^A , which depends on the corresponding proposer. If that proposer crashes or its quorum does not include the acceptor, the sub-process terminates because FP_{wskip} returns true and the default branch is *Abort*, which terminates immediately. The termination of a sub-process with a correct proposer requires the termination of that proposer. Thus, we need to prove that correct proposers terminate to prove termination for our model.

If the set of acceptors in a proposer's quorum A_Q that are correct is not enough to form a majority of acceptors, that proposer repeatedly restarts the algorithm. In this case the proposer will be unable to issue a valid proposal. Because FP_{crash} returns true if $A_Q \setminus \mathbb{F}$, where \mathbb{F} is the set of processes permanently suspected by a failure detector in $\diamond \mathcal{S}$, is not a quorum, the proposer eventually crashes. Proposers either complete the Paxos algorithm after phase 2a or crash.

In [11] Lamport describes a scenario in which two proposers loop endlessly, never having their proposals accepted: Proposer p completes phase 1 for a proposal number n_1 . Another proposer q then completes phase 1 for a proposal number $n_2 > n_1$. Proposer p 's phase 2 accept requests for a proposal numbered n_1 are ignored

because the acceptors have all promised not to accept any new proposal numbered less than n_2 . So, proposer p then begins and completes phase 1 for a new proposal number $n_3 > n_2$, causing the second phase 2 accept requests of proposer q to be ignored. And so on.

From [11] we know that this problem is solved by electing a single distinguished proposer to be the leader. The leader eventually picks a proposal number high enough for its proposal to be accepted. The model assumes some sort of leader selection. A new leader is elected when the previous leader terminates.

4.3.2 Agreement

Any proposer p requires that the set of correct acceptors in its quorum of acceptors is itself a quorum, i.e. an accepting set with which a value can be chosen [10]. Should message loss occur in labels $l1a$ or $l1b$, p restarts the algorithm. This broadcast is weakly reliable and thus only fails when p crashes or terminates because FP_{wskip} disallows suspicion of correct live proposers in acceptors. Given the definition of `promiseValue`, p will only propose a fresh value if none of the acceptors have accepted a proposal yet. At least one acceptor in every other proposer's quorum is contained in p 's quorum. Thus, if a majority of acceptors accept p 's proposal it is sent to every other proposer when they reach phase 1b. These proposers then propagate the accepted value by proposing it again but with a higher proposal number. This way all correct acceptors accept the same value.

4.3.3 Validity

To prove validity for our model we examine the communication structure and the origins of the accepted values. Because the model is well-typed we know the communication structure is as specified in global type $G_{P,AQ}$. From [13] we know that validity then holds globally if it holds for each local process.

Labels $l1b$ and $l2a$ are used to send values that can be accepted.

Label $l1b$ is used to send messages of sort `Promise Value`. These messages are sent from the acceptors to a proposer and may contain the acceptors' accepted proposal. Should an acceptor previously have accepted a proposal, that proposal then contains the accepted value. The accepted proposal pr is either the acceptor's initial accepted proposal pr_a or a proposal that was previously proposed by a proposer. pr is sent over $l1b$ without alteration. The proposer receiving these messages stores them in \vec{V} without changing their values.

Proposers send a message of sort `Proposal Value` to their quorum of acceptors over label $l2a$. To do so, proposers pick the best value from a proposal in \vec{V} , if any is available, with `promiseValue`. An entry in \vec{V} contains a proposal $prop$ if it is of the form `Promise Just prop`. These proposals are either some acceptor's initial accepted proposal or a proposal proposed by a proposer. Not all entries in \vec{V} contain a proposal but if at least one does, `promiseValue` returns the value of one of them. If no entry in \vec{V} contains a proposal a fresh value is chosen and returned. In both cases the return value of `promiseValue` is not altered before being sent over label $l2a$, which constitutes proposing that value. Acceptors that receive and accept this proposal store it without alteration.

Since label $l1a$ is not used to transmit values that can be accepted, we conclude that validity holds for each local process and thus globally.

Bibliography

- [1] M.K. Aguilera, W. Chen, and S. Toueg. “Heartbeat: A timeout-free failure detector for quiescent reliable communication”. In: *Distributed Algorithms*. Ed. by M. Mavronicolas and P. Tsigas. Berlin, Heidelberg: Springer Berlin Heidelberg, 1997, pp. 126–140. ISBN: 978-3-540-69600-1.
- [2] L. Bettini et al. “Global Progress in Dynamically Interleaved Multiparty Sessions”. In: *CONCUR 2008 - Concurrency Theory*. Ed. by F. van Breugel and M. Chechik. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 418–433. ISBN: 978-3-540-85361-9.
- [3] L. Bocchi et al. “A Theory of Design-by-Contract for Distributed Multiparty Interactions”. In: *CONCUR 2010 - Concurrency Theory*. Ed. by P. Gastin and F. Laroussinie. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 162–176. ISBN: 978-3-642-15375-4.
- [4] L. Caires and H.T. Vieira. “Conversation types”. In: *Theoretical Computer Science* 411.51 (2010). European Symposium on Programming 2009, pp. 4399–4440. ISSN: 0304-3975. DOI: <https://doi.org/10.1016/j.tcs.2010.09.010>. URL: <https://www.sciencedirect.com/science/article/pii/S0304397510004895>.
- [5] T.D. Chandra and S. Toueg. “Unreliable Failure Detectors for Reliable Distributed Systems”. In: *J. ACM* 43.2 (Mar. 1996), pp. 225–267. ISSN: 0004-5411. DOI: [10.1145/226643.226647](https://doi.org/10.1145/226643.226647). URL: <https://doi.org/10.1145/226643.226647>.
- [6] M. Coppo et al. “A Gentle Introduction to Multiparty Asynchronous Session Types”. In: *Formal Methods for Multicore Programming: 15th International School on Formal Methods for the Design of Computer, Communication, and Software Systems, SFM 2015, Bertinoro, Italy, June 15-19, 2015, Advanced Lectures*. Ed. by M. Bernardo and Einar B. Johnsen. Cham: Springer International Publishing, 2015, pp. 146–178. ISBN: 978-3-319-18941-3. DOI: [10.1007/978-3-319-18941-3_4](https://doi.org/10.1007/978-3-319-18941-3_4). URL: https://doi.org/10.1007/978-3-319-18941-3_4.
- [7] G. Coulouris, J. Dollimore, and T. Kindberg. *Distributed Systems: Concepts and Design (3rd Edition)*. Addison-Wesley, 2001, p. 452.
- [8] K. Honda, N. Yoshida, and M. Carbone. “Multiparty Asynchronous Session Types”. In: *Proceedings of the 35th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. POPL ’08. San Francisco, California, USA: Association for Computing Machinery, 2008, pp. 273–284. ISBN: 9781595936899. DOI: [10.1145/1328438.1328472](https://doi.org/10.1145/1328438.1328472). URL: <https://doi.org/10.1145/1328438.1328472>.
- [9] K. Honda, N. Yoshida, and M. Carbone. “Multiparty Asynchronous Session Types”. In: *J. ACM* 63.1 (Mar. 2016). ISSN: 0004-5411. DOI: [10.1145/2827695](https://doi.org/10.1145/2827695). URL: <https://doi.org/10.1145/2827695>.
- [10] L. Lamport. “Lower Bounds for Asynchronous Consensus”. In: *Distrib. Comput.* 19.2 (Oct. 2006), pp. 104–125. ISSN: 0178-2770. DOI: [10.1007/s00446-006-0155-x](https://doi.org/10.1007/s00446-006-0155-x). URL: <https://doi.org/10.1007/s00446-006-0155-x>.

-
- [11] L. Lamport. “Paxos Made Simple”. In: *ACM SIGACT News (Distributed Computing Column)* 32, 4 (Whole Number 121, December 2001) (Dec. 2001), pp. 51–58. URL: <https://www.microsoft.com/en-us/research/publication/paxos-made-simple/>.
 - [12] R. Milner, J. Parrow, and D. Walker. “A calculus of mobile processes, I”. In: *Information and Computation* 100.1 (1992), pp. 1–40. ISSN: 0890-5401. DOI: [https://doi.org/10.1016/0890-5401\(92\)90008-4](https://doi.org/10.1016/0890-5401(92)90008-4). URL: <https://www.sciencedirect.com/science/article/pii/0890540192900084>.
 - [13] K. Peters, U. Nestmann, and C. Wagner. “Fault-Tolerant Multiparty Session Types”. Provided by K. Peters. 2021.
 - [14] A. Scalas and N. Yoshida. “Multiparty session types, beyond duality”. In: *Journal of Logical and Algebraic Methods in Programming* 97 (2018), pp. 55–84.
 - [15] A.S. Tanenbaum and M. van Steen. *Distributed Systems: principles and paradigms*. Pearson Prentice Hall, 2017.
 - [16] M. Viering et al. “A Typing Discipline for Statically Verified Crash Failure Handling in Distributed Systems”. In: *Programming Languages and Systems*. Ed. by A. Ahmed. Cham: Springer International Publishing, 2018, pp. 799–826. ISBN: 978-3-319-89884-1.
 - [17] N. Yoshida et al. “Parameterised Multiparty Session Types”. In: *Foundations of Software Science and Computational Structures*. Ed. by L. Ong. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 128–145. ISBN: 978-3-642-12032-9.