

Tree Topologies for Causal Message Delivery

Sebastian Blessing
Department of Computing
Imperial College London
sebastian.blessing12@imperial.ac.uk

Sylvan Clebsch
Microsoft Research
sylvan.clebsch@microsoft.com

Sophia Drossopoulou
Department of Computing
Imperial College London
s.drossopoulou@imperial.ac.uk

CCS Concepts • Computer systems organization → Distributed architectures; • Theory of computation → Operational semantics;

Keywords network topologies; causal messaging

ACM Reference Format:

Sebastian Blessing, Sylvan Clebsch, and Sophia Drossopoulou. 2017. Tree Topologies for Causal Message Delivery. In *Proceedings of 7th ACM SIGPLAN International Workshop on Programming Based on Actors, Agents, and Decentralized Control (AGERE'17)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3141834.3141835>

Abstract

Causal message delivery, i.e. the requirement that messages are delivered in an order respecting their causal (logical) dependencies, is often mandated in the distributed setting. So far, causal message delivery has been implemented by augmenting messages with meta data information that allows the receiver (or the platform) to re-order, and if necessary hold back, messages upon receipt before processing.

We propose that causal message delivery can be achieved by construction, simply by organizing the nodes of the distributed application into a tree topology, and without the need for any meta data in the messages.

We present our ideas informally through an example application and then develop a formal model and prove that causal message delivery is preserved in tree-based networks.

1 Introduction

The ordering of events in distributed applications is an important artifact. The ability to observe it requires a relationship between messages that induces the idea of “happens-before” [15]. This ordering can be either *total* [15] or *partial*. Partial orderings can be categorized in *causal dependencies* [?] or *explicitly stated relationships* [5].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org. *AGERE'17, October 23, 2017, Vancouver, Canada*

© 2017 Copyright held by the owner/author(s). Publication rights licensed to Association for Computing Machinery.

ACM ISBN 978-1-4503-5516-2/17/10...\$15.00

<https://doi.org/10.1145/3141834.3141835>

In the context of causal messaging, we say that each message is an *effect* and every message that was received or sent prior to that message is a *cause* of that effect. Causal ordering is preserved if every cause is enqueued before its effect and is therefore less restrictive than a total order, cheaper to achieve but strong enough for most applications [?]. A substantial amount of the research in the field of distributed systems focused on *detecting* [15, 19] causality violations or *enforcing* [6?] ordered message delivery. Merely being able to detect causality violations puts more complexity into the application and eventually requires recovery protocols to continue processing. Hence, there is a case for preferring enforcement over detection.

Motivation In current practice, programmers are required to design an application to tolerate uncoordinated and inconsistent orderings of messages or need to implement a distributed algorithm to explicitly prevent such re-orderings from ever occurring. These techniques, however, either come with linear (or even quadratic) space complexity over the set of involved processes/nodes or employ some central coordinator that, depending on the application, may be highly contended and performance critical. That is, there is a trade-off between tracking causal dependencies and achieving high availability and/or scalability.

Since *enforcing* causal message delivery is potentially expensive, and increasingly so in large distributed systems, reducing the overhead that comes with it can have a substantial impact on scalability, without being forced to give up as much on availability as traditional approaches. Central to our work is a discussion on to which extend messages need to be augmented with meta data (if at all) and whether running a distributed algorithm to ensure causality is a *necessity*.

Background Observing the order of events in distributed systems is a complex and powerful concept that forms the basis for a wide range of applications, such as distributed snapshots [4, 16], version vectors and conflict detection in distributed databases [11, 24], causal consistency [1, 2] as well as causal broadcast [6]. Causality is a generalization of FIFO ordering [26]. The subtle difference between FIFO and causal ordering is that the former only applies to messages having the same sender, whereas causality applies to messages from any process or node [15]. Besides having provided a definition of the “happens-before” [15] relationship and being able to observe it, one of the most important findings of Lamport’s work on ordering events in distributed

systems [15] is that reasoning about “happens-before” does not require synchronous coordination between processes. Moreover, one of the most interesting properties is that once causal ordering is established, many other concepts and their implementation become much simpler [7].

Causality means that processes must decide whether a message can be received or if it needs to be held back as there might be dependent messages in-flight but out of order. Detecting that receiving a message is safe is a complex task. This decision and collecting the required meta data can be space and time consuming [5, 8, 13, 20]. Moreover, inspecting physical timestamps is not enough, since time is not necessarily moving forward uniformly on all nodes in a distributed system.

A classic approach to this problem is to implement a message service based on Lamport clocks, using logical time to order messages on every node [15, 19]. Lamport clocks record causal dependencies by tagging every message sent with a *single* increasing integer value. As a result, Lamport clocks have a linear space requirement (over all involved processes) and do not allow for completely independent messages. The mechanism described in [15] induces a total order over all messages by making sure that every site of the system has received at least one message from every other site before progressing. That is, $O(n)$ acknowledgment messages. This is immensely expensive and implies more dependencies than truly exist. Moreover, a total ordering is too strict for certain applications [11], where detecting that some message cannot possibly be preceded by another message is key for application correctness. This issue can be addressed by vector clocks [12, 15], where every process maintains an integer value for every other process of the system. However, the price to pay is quadratic space complexity and the amount of nodes involved in a system is static and needs to be known a priori. Dynamic topologies are addressed by interval tree clocks in [3].

Charron-Bost provided a surprising result and theoretical argument [8], saying that the amount of space required for the necessary meta data to detect causal dependencies in distributed systems is at least linear. That is, there is no better time stamping algorithm based on time vectors smaller than $O(n)$ that “truly characterizes causality” [8, 13, 20]. The fundamental concepts have been developed in the late 1970s, and from then on research has mainly focused on developing protocols to tackle the amount of (global) space required as well as decoupling the effects causal ordering has on unrelated messages in terms of messaging performance. There is no free lunch and there are only a few strategies to reduce these overheads.

The most intuitive idea would be to restrict the number of participants. This is exactly the strategy distributed databases use in the context of version vectors [11] by intelligently picking the items for which causal dependencies are tracked, e.g. data item replicas instead of all servers. Since

the number of replicas is usually substantially smaller than the number of involved machines in a distributed system, this is a powerful optimization in practice [17, 18].

Inducing full causality assumes that all messages matter equally. Generally, for most types of applications, this is not the case. Hence, another optimization would be to explicitly specify relationships between items [5]. In the best case, this can reduce the local space required to $O(1)$. In the context of applications where a total order of events is not paramount, such as Twitter, this is an effective and powerful trade-off between performance and the non-deterministic order of messages. The advantage of these approaches is that they have no negative impact on availability and only give up on causal dependencies that are not significant to the specific application [22, 23].

Reducing the number of concurrent participants in a vector clock based system [12] is possible and saves space, but sacrifices some amount of availability. A system with only one concurrent participant implies a total order of events [14]. This is even the case in data race free multi-core environments [9].

Some *eventually consistent* applications may drop “happens-before” ordering entirely. Systems belonging to this application domain are not subject to the ideas presented in this work.

Contributions Our technique *enforces* (rather than *detecting*) causal message delivery by arranging the nodes in a tree topology, resulting in a mechanism that eliminates the need to collect meta data for detecting causal relationships entirely, and that can enforce causal order rather than merely detecting violations after the fact. We either physically connect individual nodes to a tree network or put a networking service in place to arrange sockets accordingly. During runtime of the system, no main memory or CPU time is occupied with the task of deciding whether or not a message can be delivered. Instead, they can be processed *immediately* upon receipt. However, structured networks, such as trees, do limit the choices that can be made between network configurations. Moreover, for the mechanism presented in this paper to work, we do assume that sending messages between actors on a single node respects causal order.

We expect that substantial code reductions are possible when developing distributed applications, as software systems can implicitly rely on causally ordered events. In other words, in tree networks, causality is an invariant and an inherent property of a systems behavior.

The key contribution of this paper are:

- A zero-space mechanism to enforce causal ordering by topological arrangement.
- A formal model to describe message passing in tree networks, described through an operational semantics.
- A sketch of the proof that causality is preserved.

Outline We present an informal view of tree-based causality in section 2, formalize it in section 3, and provide a sketch for a correctness proof in section 4. We conclude and discuss further work in section 5.

2 Tree-based Causality

In this section we describe what we mean by causality, by causal message delivery, and the TCP/IP delivery, and explain informally how a tree-based topology ensures causal message delivery.

Preliminaries We say that a message **causes** another message, if one of the following three rules applies:

1. If an actor (or process) receives a message and later on sends another message, then the first message is a cause of the second message.
2. If an actor (or process) sends a message and later on sends another message, then the first message is a cause of the second message.
3. Causality is transitive.

Causal message delivery mandates that if two messages are in a causal relationship and have the same recipient, then the causing message will be delivered before the other one. In Lamport's terminology, [15], if two messages are in a causal relationship, the message that *logically precedes* some other messages is to be delivered *before* the message it is related to. We expect that the actors are placed on nodes in a *tree*. The tree topology does not restrict which actors may communicate with each other – in fact, actor references can form an arbitrary graph. But it means that messages to actors in different nodes must travel via the node which is the most common ancestor of the nodes holding the two actors. In this work we are not concerned with how the tree is built, nor with how the actors are placed on the nodes – an algorithm to arrange nodes in a tree topology with a networking service is presented in [7].

We assume that communication across nodes is consistent with **TCP/IP**: packages (or messages) sent across one network connection arrive in the order they were sent – *i.e.* there is no overtaking of messages between adjacent nodes.

Example Consider three actors (or processes): a *Customer*, a *Shop* and a *Bank*: The *Customer* intends to buy an item from the *Shop* by electronic cash. She makes sure that the balance on the bank account is high enough by sending a "credit" message to the *Bank*, and then informs the *Shop*, by sending it a "buy" message. Upon receipt of the "buy" message, the *Shop* sends a "debit" message to the *Bank*. Obviously, it is crucial that the *Bank* receives the "credit" message before it receives the "debit" message.

Causal message delivery guarantees exactly that! Namely, "credit" causes "buy" (rule 2), "buy" causes "debit" (rule 1), and thus "credit" causes "debit" (rule 3). Since "credit" and "debit" have the same recipient, causality guarantees that

"credit" will arrive first! On the other hand, even though "credit" causes "buy" there is no guarantee as to whether "credit" will arrive at the *Bank* before "buy" arrives at the *Shop*.

In this paper we argue that in a distributed system, a tree topology of the nodes on which the actors are running implicitly guarantees causal message delivery.

To continue with our example, consider a system consisting of four nodes, i_1 , i_2 , i_3 , and i_4 , arranged in a tree as shown in figure 1, and where *Customer*, *Shop* and *Bank*, are scheduled on nodes i_1 , i_3 , and i_4 respectively.

In figure 1 we show four snapshots from one possible scenario. In the first snapshot, *Customer* sends "credit" to the *Bank*: since i_1 has no direct connection to i_3 , the message is first sent to i_2 . In the second snapshot, i_2 forwards "credit" to i_3 , and *Customer* sends "buy" through i_2 . Note that, a scenario where i_2 first sends "buy" and then "credit" is impossible, because i_2 will forward the messages in the order they arrive, and because TCP/IP guarantees that the message do not overtake each other on the wire. In the third snapshot, "buy" has been forwarded from i_2 to i_3 , and from i_3 to i_4 . In the last snapshot, *Shop* has received "buy", and as a result it sends "debit" to the *Bank*. In this scenario, causal message delivery is respected!

But there are more possible scenarios on the same topology: other messages may be sent interleaved with the messages from our example, and "buy" may be sent only *after* "credit" has been received. Nevertheless, even though the timing of a message may vary, its route is always the same. In this topology, the route of "credit" is a prefix of the route of "buy", and because "credit" starts its journey before "buy", we have the guarantee, by TCP-IP, that "buy" cannot overtake "credit". Therefore, "credit" will arrive at the *Bank* before "buy" arrives at the *Shop*. And because "buy" caused "debit", we also have that "credit" will arrive at the *Bank* before "debit".

But will causal message delivery be guaranteed on other topologies? In figure 2 we see two further topologies. We can see that causal message delivery for our example is guaranteed for these topologies too, but for slightly different reasons:

In the left topology in figure 2, "buy" may arrive at the *Shop* before "credit" arrives at the *Bank*, even though "credit" causes "buy" – remember that causal message delivery is agnostic about delivery on different nodes. Nevertheless, "debit"'s route has "credit"'s route as a suffix, and as "debit" will start after "credit", it cannot overtake "credit". Thus, "debit" will be delivered after "credit" has been delivered.

In the right topology in figure 2, we also have that "buy" may arrive at the *Shop* before "credit" arrives at the *Bank*. However, "credit" is guaranteed to arrive at i_4 before "buy". Therefore, "debit" will be placed on i_1 's queue after "credit", and will arrive at *Bank* after "credit".

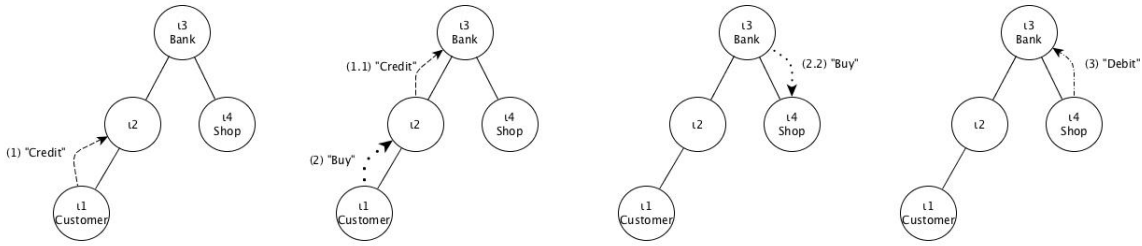


Figure 1. One possible scenario on one tree topology

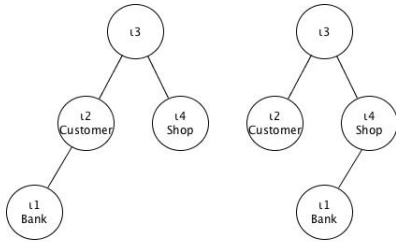


Figure 2. Two further topologies

So far, we considered three different topologies with four nodes, and in these topologies, for three actors and three messages, causal message delivery is guaranteed. We used the argument that the route taken by the “causing” message is included in some way in the route of the combined messages it “caused”.

But does this argument apply over any topologies with any number of nodes, any number of actors, and any number of messages involved in a causal relationship, and any number of messages unrelated to the causal chain?

The answer is yes. However, even though at an intuitive level the guarantee that a tree topology guarantees causal message delivery seemed convincing, it became less so when considering more cases. For this reason, we developed a formal model and a proof, which we will present in the next section.

3 Formal Model

In this section we develop a formal model, expressed as an operational semantics for sending and receiving messages in a tree-based network. The runtime configuration \mathcal{T}, Σ does not reflect the causal relationship of messages. This can only be deduced from the “history” ghost state Δ . The core of our argument is the “arrives-before” relation between messages, $<$, which is defined so that if $m < m'$ can be observed in the current snapshot \mathcal{T}, Σ , then m will arrive ι before m' arrives at ι , where ι is m 's destination. Interestingly, the relation is not transitive, and even allows in some cases that $m < m'$ as well as $m' < m$. We prove that causality, which can be observed in Δ , implies the arrives-before relation. We

were surprised that we were able to prove our result without needing to reflect time – or even logical time – in the model.

3.1 Runtime entities and operational semantics

Figure 12 presents the runtime entities and shows the operational semantics describing the behavior of a tree-based distributed system.

Runtime entities Both node and message identifiers, referred to as ι and μ , are (uniquely) picked out of the set of natural numbers. A message m consists of an ID and its origin node. At any point in time during execution of the system a message is contained in exactly one message queue until it was received. That is, message delivery is immediate. We model a message queue q as an unbounded sequence of messages. The tree topology \mathcal{T} is modeled as a mapping from node identifier to node identifier. Every node within the tree holds a pointer to its parent. The root node points to itself. Every node n holds separate in- and out-queues for each of its children and its parent node. In-queues contain messages being received from the node related to that neighbor, and out-queues contain messages that need to be delivered to the next hop towards the destination of a message, respectively. We separated sender from receiver queues in order to model unbounded point-to-point delay, resulting in ordered but arbitrarily delayed distributed network semantics. The state of every node, i.e. the contents of its message queues, can be observed in Σ . Furthermore, events are defined as either sending or receiving a message m at some node ι , and are recorded as shadow state in the history Δ . Paths are sequences of node identifiers.

Note that we are not modeling actors and their state as part of the system, since sending messages between actors on a single distributed participant is causal [9] and therefore we are only interested in inter-node communication. Auxiliary helper functions to mutate runtime entities such as queues and the history to allow for a succinct presentation of the operational semantics are briefly discussed in section ?? and shown in figure ?? and ??.

Notation In this paper, we make use of some additional notation for convenience.

- We use set operations on the domains of mappings:

- $x \in \text{map} \iff x \in \text{dom}(\text{map})$
- We use set operations on sequences:
 - $e \in \text{seq} \iff \text{seq} = _ \cdot e \cdot _$
- We use set operations to express the existence of queues and messages:
 - $q \in \Sigma \iff \exists t, t'. \text{In}(\Sigma, t, \nu) = q \vee \text{Out}(\Sigma, t, t') = q$
 - $m \in \Sigma \iff \exists q \in \Sigma. m \in \Sigma$
 - $m \in n \iff \exists t. [m \in n(t) \downarrow_1 \vee m \in n(t) \downarrow_2]$
 - $m \in q \iff \exists q' \in \Sigma. q' = _ \cdot m \cdot _ \wedge q = q'$
- We use index operations to examine sequences:
 - $\text{Seq}[k]$ is the k^{th} item in Seq

Operational semantics In a distributed system, execution progresses by sending and receiving messages. For any message to reach a node from its origin within a tree network, a routing and forwarding scheme is required. Our operational semantics consist of four rules, discussed in more detail in the following.

M-Send Sending a new message. The newly created message at node t_1 with a previously unknown identifier μ is pushed to the out-queue towards its destination node t_2 . Sending this message is recorded in Δ .

M-Route Routing a message to the next node. A message m is at the top of the in-queue at node t_1 and was forwarded from node t_2 . However, t_1 is not the final destination of m . Consequently, the message is dispatched from that in-queue to an out-queue towards a node closer to $\text{Dest}(m)$.

M-Forward Forwarding a message to the next node. The message m has been newly created or routed at node t_1 and was pushed to the out-queue for t_2 . Messages on out-queues are always delivered to the next node on the path from t_1 to the destination of m . Note that both routing and forwarding are not recorded within the history Δ .

M-Receive Receiving a message. The message m meant for node t_1 has reached the head of some in-queue at t_1 . Therefore, execution for delivering m is finished. The message is consumed by popping it from the respective queue and its receipt is recorded in Δ .

To illustrate introduced runtime entities, we refer to the example application discussed in section 2, figure 1. The initial configuration, before forwarding “credit” (m_0) looks like the following:

$$\begin{aligned} \mathcal{T} &= [t_1 \rightarrow t_2, t_2 \rightarrow t_3, t_3 \rightarrow t_3, t_4 \rightarrow t_3] \\ n_1 &= [t_2 \rightarrow (\epsilon, m_0 \cdot \epsilon)] \\ n_2 &= [t_1 \rightarrow (\epsilon, \epsilon), t_3 \rightarrow (\epsilon, \epsilon)] \\ n_3 &= [t_2 \rightarrow (\epsilon, \epsilon), t_4 \rightarrow (\epsilon, \epsilon)] \\ n_4 &= [t_3 \rightarrow (\epsilon, \epsilon)] \\ \Sigma &= [t_1 \rightarrow n_1, t_2 \rightarrow n_2, t_3 \rightarrow n_3, t_4 \rightarrow n_4] \\ \Delta &= \text{Snd}(m_0, t_1) \end{aligned}$$

Assuming the messages “buy” and “debit” are given the identifiers m_1 and m_2 , a (possible) configuration right before the *Shop* asks the *Bank* for money is:

$$\begin{aligned} \mathcal{T} &= [t_1 \rightarrow t_2, t_2 \rightarrow t_3, t_3 \rightarrow t_3, t_4 \rightarrow t_3] \\ n_1 &= [t_2 \rightarrow (\epsilon, \epsilon)] \\ n_2 &= [t_1 \rightarrow (\epsilon, \epsilon), t_3 \rightarrow (\epsilon, \epsilon)] \\ n_3 &= [t_2 \rightarrow (\epsilon, \epsilon), t_4 \rightarrow (\epsilon, \epsilon)] \\ n_4 &= [t_3 \rightarrow (\epsilon, m_2 \cdot \epsilon)] \\ \Sigma &= [t_1 \rightarrow n_1, t_2 \rightarrow n_2, t_3 \rightarrow n_3, t_4 \rightarrow n_4] \\ \Delta &= \text{Snd}(m_0, t_1) \cdot \text{Snd}(m_1, t_1) \cdot \text{Rec}(m_0) \cdot \text{Rec}(m_1) \cdot \text{Snd}(m_2, t_4) \end{aligned}$$

We can observe for this instance of the applications execution that both m_0 and m_1 were already received before m_2 was sent at t_4 . Moreover m_0 was received before m_1 was. Consequently, it is impossible for m_2 to arrive at t_3 after the next configuration that would follow from the application of the operational semantics. That is, causal message delivery is preserved. Note that the delay between sending messages can be arbitrary. Message m_0 could also have been received before m_1 was even sent. In fact, for this particular example scenario, there are *two* possible instances of Δ :

$$\begin{aligned} \Delta_1 &= \text{Snd}(m_0, t_1) \cdot \text{Rec}(m_0) \cdot \text{Snd}(m_1, t_1) \cdot \text{Rec}(m_1) \cdot \text{Snd}(m_2, t_4) \cdot \text{Rec}(m_2) \\ \Delta_2 &= \text{Snd}(m_0, t_1) \cdot \text{Snd}(m_1, t_1) \cdot \text{Rec}(m_0) \cdot \text{Rec}(m_1) \cdot \text{Snd}(m_2, t_4) \cdot \text{Rec}(m_2) \end{aligned}$$

The only difference between the two histories is the point in time at which m_0 is received, as m_0 is sent asynchronously and there is no upper bound on when to send m_1 after m_0 has been put on the out-queue towards t_1 . However, the order of $\text{Snd}(m_0, t_1)$ and $\text{Snd}(m_1, t_1)$ cannot be different, as m_0 causes m_1 . The same applies to the order of receiving m_1 and sending m_2 . This leaves us with only Δ_1 and Δ_2 as possible histories for the scenario illustrated in figure 1.

We now introduce several definitions to form the basis of our formal argument. We show that the order of messages delivery in our example application is causal in *any* well-formed tree, state and configuration and can be observed in any possible recording of Δ .

3.2 Well-formed trees, paths and next hops

In this section we introduce several definitions we will be using throughout the discussion of the formal system and proofs. At first, we discuss the basics of children and descendants as well as our definition of paths and next hops in tree networks.

Definition 3.1. (Children and descendants). Given $cfg = (\mathcal{T}, \Sigma)$, we define:

- The children of t in \mathcal{T} is the set of nodes which are in the range of $\mathcal{T}(t)$, except t itself:
 - $\text{children}(t, \mathcal{T}) = \{t' \mid \mathcal{T}(t') = t \wedge t' \neq t\}$

ι	\in	NodeId	::=	\mathbb{N}
μ	\in	MessageId	::=	\mathbb{N}
m	\in	Message	::=	MessageId \times NodeId
q	\in	Queue	::=	Message*
\mathcal{T}	\in	Tree	::=	NodeId \rightarrow NodeId
n	\in	Node	::=	NodeId \rightarrow Queue \times Queue
Σ	\in	State	::=	NodeId \rightarrow Node
e	\in	Event	::=	Snd(m, ι) Rec(m)
Δ	\in	History	::=	Event*
p	\in	Path	::=	$\iota \cdot p$

M-SEND	
$\text{Snd}((\mu, _), _) \notin \Delta$	
$\frac{}{\mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow \Delta \cdot \text{Snd}((\mu, \iota_2), \iota_1), \Sigma[\iota_1 \mapsto \text{PushOut}(\Sigma, \mathcal{T}, \iota_1, (\mu, \iota_2))]}$	
M-ROUTE	
$\text{TopIn}(\Sigma, \iota_1, \iota_2) = m$	
$\frac{}{\mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow \Delta, \Sigma[\iota_1 \mapsto \text{PopIn}(\Sigma, \iota_1, \iota_2)], [\iota_2 \mapsto \text{PushOut}(\Sigma, \mathcal{T}, \iota_1, m)]}$	
M-FORWARD	
$\text{TopOut}(\Sigma, \iota_1, \iota_2) = m$	
$\frac{}{\mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow \Delta, \Sigma[\iota_1 \mapsto \text{PopOut}(\Sigma, \iota_1, \iota_2), \iota_2 \mapsto \text{PushIn}(\Sigma, \iota_1, \iota_2, m)]}$	
M-RECEIVE	
$\text{TopIn}(\Sigma, \iota_1, \iota_2) = (\mu, \iota_1)$	
$\frac{}{\mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow \Delta \cdot \text{Rec}((\mu, \iota_1)), \Sigma[\iota_1 \mapsto \text{PopIn}(\Sigma, \iota_1, \iota_2)]}$	

Figure 3. Runtime entities and operational semantics

$\text{next} ::= \text{NodeId} \times \text{NodeId} \times \mathcal{T} \rightarrow \text{NodeId}$

$$\text{next}(\iota_1, \iota_2, \mathcal{T}) = \begin{cases} \iota_2 & \text{if } \iota_2 \in \text{children}(\iota_1, \mathcal{T}) \\ \mathcal{T}(\iota_1) & \text{if } \iota_2 = \mathcal{T}(\iota_2) \\ \text{next}(\iota_1, \mathcal{T}(\iota_2), \mathcal{T}) & \text{otherwise} \end{cases}$$

Figure 4. Next hop towards a destination node

$\text{path} ::= \text{NodeId} \times \text{NodeId} \times \mathcal{T} \rightarrow \text{NodeId}^*$

$$\text{path}(\iota_1, \iota_2, \mathcal{T}) = \begin{cases} \iota_1 \cdot \text{path}(\text{next}(\iota_1, \iota_2, \mathcal{T}), \iota_2, \mathcal{T}) & \text{if } \iota_1 \neq \iota_2 \\ \iota_1 & \text{if } \iota_1 = \iota_2 \end{cases}$$

Figure 5. Paths in a tree

$$\frac{\bar{\iota} \neq \epsilon}{p \sqsubset \bar{\iota} \cdot p} \text{ P-INCL}$$

Figure 6. Path inclusion

- The descendants of a node ι in \mathcal{T} are described as the transitive closure of the child set:
 - $\text{descendants}(\iota, \mathcal{T}) = \{\iota' \mid \exists \iota'' . \iota'' \in \text{children}(\iota, \mathcal{T}) \wedge \iota' \in \text{descendants}(\iota'', \mathcal{T})\} \cup \text{children}(\iota, \mathcal{T})$

Figure 5 shows how to compute the next hop towards a node ι_2 from ι_1 . The next hop is trivial to decide if ι_2 is in the child set of ι_1 . Similarly, if ι_2 is the root node, the next hop is the parent of ι_1 . In the general case, the next node can be decided by determining the next hop from ι_1 to the parent of ι_2 .

The concatenation of next hops from ι_1 to ι_2 is called the *path* from ι_1 to ι_2 , as defined in figure 6. As mentioned in section 2, path inclusion is one of the important pieces to solve the puzzle why message delivery in tree networks is causal. Formalized in figure 9, throughout the discussion of this paper, we refer to it as a path p is *totally included* in a path p' if p is the suffix of the non-empty sequence p' .

3.3 Well-formed configuration

Definition 3.2. (A well-formed tree). We say that a tree \mathcal{T} is well-formed, formally $\models \mathcal{T}$, if there are no cycles and if there is exactly one root node:

- $\forall \iota \in \text{dom}(\mathcal{T}). \iota \notin \text{descendants}(\iota, \mathcal{T})$
- $\exists! \iota \in \text{dom}(\mathcal{T}). \mathcal{T}(\iota) = \iota$

Definition 3.3. (A well-formed state). We say that the state Σ is well-formed, formally $\models \Sigma$, if a message only ever appears in exactly one message queue and if the destination of every message is a valid node in Σ :

- $\forall m \in \Sigma. \exists! q \in \Sigma. m \in q$
- $\forall \mu, \iota [(\mu, \iota) \in \Sigma \rightarrow \iota \in \text{dom}(\Sigma)]$

$<_{\Delta} \subseteq \text{Event} \times \text{Event}$	$<_{\Delta} \subseteq \text{Message} \times \text{Message}$	$\ll_{\Delta} \subseteq \text{Message} \times \text{Message}$
$e <_{\Delta} e'$	$m <_{\Delta} m'$	$m \ll_{\Delta} m'$
iff	iff	iff
$\Delta = _ \cdot e \cdot _ \cdot e' \cdot _$	$\text{Rec}(m) <_{\Delta} \text{Snd}(m', \text{Dest}(m))$	$m <_{\Delta} m'$
	\vee	\vee
	$\exists i. \text{Snd}(m, i) <_{\Delta} \text{Snd}(m', i)$	$\exists m''. (m \ll_{\Delta} m'' \wedge m'' <_{\Delta} m')$

Figure 7. Causally dependent events and messages

$\mathcal{T}, \Delta, \Sigma \vdash m < m'$	iff
1. $\exists q \in \Sigma. q = _ \cdot m \cdot _ \cdot m' \cdot _$	
\vee	
2. $\exists i, i_1, i_2 \in \Sigma. [m \in \text{Out}(\Sigma, i, i_1) \wedge m' \in \text{In}(\Sigma, i, i_2)]$	
\vee	
3. $\exists i_1, i_2 \in \Sigma [m \in \Sigma(i_1) \wedge m' \in \Sigma(i_2) \wedge i_1 \neq i_2 \wedge \text{path}(i_1, m, \mathcal{T}) \sqsubset \text{path}(i_2, \text{Dest}(m), \mathcal{T})]$	
\vee	
4. $\text{Dest}(m) \neq \text{Dest}(m') \wedge [\text{Rec}(m) \in \Delta \vee \text{Rec}(m') \in \Delta]$	
\vee	
5. $\text{Dest}(m) = \text{Dest}(m') \wedge [\text{Rec}(m') \in \Delta \rightarrow \text{Rec}(m) <_{\Delta} \text{Rec}(m')]$	

Figure 8. Message arrives at its destination before another message

Definition 3.4. (A well-formed state for a given tree). We say that the state Σ is well-formed with respect to \mathcal{T} , formally $\mathcal{T} \models \Sigma$, if all queues point to valid nodes, and if the state only contains valid nodes. All messages reside in a queue which is on a node along the route a message has to travel in order to reach its destination (figures 5 and 6). That is, formally:

- $\text{dom}(\Sigma) = \text{dom}(\mathcal{T})$
- $\forall i, i' \in \Sigma. [\Sigma(i, i')$ is defined
 $\rightarrow i' \in (\text{children}(i, \mathcal{T}) \cup \text{children}(i', \mathcal{T}))]$
- $\forall i_1, i_2 \in \Sigma. \forall m [m \in \text{Out}(\Sigma, i_1, i_2)$
 $\rightarrow \text{next}(i_1, \text{Dest}(m), \mathcal{T}) = i_2]$
- $\forall i_1, i_2 \in \Sigma. \forall m [m \in \text{In}(\Sigma, i_1, i_2)$
 $\rightarrow \text{next}(i_2, \text{Dest}(m), \mathcal{T}) = i_1]$

Definition 3.5. (A well-formed history). We say that the history Δ is well-formed, formally $\models \Delta$, if a message is exactly sent and received once:

- $\forall \mu [\#\{k \mid \Delta[k] = \text{Rec}(\mu, _)\} \leq 1 \wedge$
 $\#\{k \mid \Delta[k] = \text{Snd}(\mu, _, _)\} \leq 1]$

Causally dependent messages. We now introduce a formal counterpart, the relation $<_{\Delta}$, to express that m caused m' . With $<_{\Delta}$ of figure 7 we apply an order on the events in Δ and combine all three relations to give a transitive version of $<_{\Delta}$, namely \ll_{Δ} (figure 7). Hence, a message m causes another message m' if any of the three cases apply:

- m was received prior to m' being sent from the node where m was received.

- m and m' have been sent from the same node, in that order. That is, the sending of m preceded the sending of m' .
- There exists a message m'' such that $<_{\Delta}$ can be established transitively using m'' .

Definition 3.6. (A well-formed state and path for a given history). We say that the state Σ is well-formed for a given history Δ , formally $\Delta \models \Sigma$, if for every message in the state there exists a Snd event recorded in the history, and, if a message was received, the corresponding Snd event must have been observed in Δ earlier with respect to $<_{\Delta}$. Also, messages can still be in transit:

- $\forall m [m \in \Sigma \rightarrow \exists i' \text{Snd}(m, i') \in \Delta]$
- $\forall m \in \Sigma [\text{Rec}(m) \in \Delta \rightarrow \text{Snd}(m, _) <_{\Delta} \text{Rec}(m)]$
- $\forall m \in \Sigma [\text{Snd}(m, _) \in \Delta \wedge \text{Rec}(m) \notin \Delta]$

The path of a message is well-formed, formally $\mathcal{T}, \Delta, \Sigma \models m$, if it does not take any unnecessary detours towards its destination, with respect to path inclusion (figure 9 and 6):

- $\forall i, i'. [m \in \Sigma(i) \wedge \text{Snd}(m, i') \in \Delta$
 $\rightarrow \text{path}(i, \text{Dest}(m), \mathcal{T}) \sqsubset \text{path}(i', \text{Dest}(m), \mathcal{T})]$

Definition 3.7. (A well-formed configuration). We say that a configuration $\text{cfg} = (\mathcal{T}, \Delta, \Sigma)$ is well-formed, formally $\text{WF}(\text{cfg})$, if the properties described in definitions 3.2 to 3.6 hold:

1. $\models \mathcal{T}$
2. $\models \Sigma$
3. $\models \Delta$
4. $\mathcal{T} \models \Sigma$

5. $\Delta \models \Sigma$
6. $\forall m \in \Sigma. \mathcal{T}, \Delta, \Sigma \models m$

3.4 Cause and arrival

We now define what it means for a message to be a cause of another message (figure 7), and what it means for a message to arrive before another one (figure 8).

We define the relation \ll_{Δ} in figure 7, where $m \ll_{\Delta} m'$ means that message m is a cause of message m' . The definition reflects the informal description of causality given at the beginning of section 2. It is easy to see that \ll_{Δ} is transitive and that $m \ll_{\Delta} m$ never holds. In a Δ_0 describing the example from section 2, we would have that $Credit \ll_{\Delta_0} Buy \ll_{\Delta_0} Debit$.

We then define a relation to express when a message m is guaranteed to “arrive at its destination before” another message m' , formally $m < m'$. The relation depends on the tree topology \mathcal{T} , the history Δ and state Σ . As described in in figure 8, the relation $\mathcal{T}, \Delta, \Sigma \vdash m < m'$ holds if any of the following cases apply:

1. m and m' are in the same queue, and m precedes m' in that queue.
2. The two messages are on the same node, but m already resides on an out-queue towards the destination of m' , whereas m' is still in an in-queue from some child node.
3. The two messages are on two different nodes, but the path of m is fully included in the path of m' towards their, potentially individual, destinations.
4. The two messages have different destinations, and one of them has already been received.
5. The two messages have the same destination, and if m' has been received, then m was received before it.

To be precise, $m_1 < m_2$ means that m_1 is guaranteed to arrive at its destination before m_2 will arrive at m_1 's destination. Therefore, the relation is applicable even if the two messages do not have the same destination, even if the second message will never visit the first message's destination, even if they should not share any part of the route. For example, assuming the topology of figure 2, and that the “Credit” message was in the out-queue of ι_2 , and the “Buy” message was at the in-queue of ι_4 . Then we would have that “Credit” and “Buy” have different destinations, and their routes have nothing in common – nevertheless, we have that $Credit < Buy$. Interestingly, we here also have that $Buy < Credit$.

Notice that while $m < m$ is never possible, the relation $<$ is neither symmetric, nor antisymmetric, nor transitive. To see an example for the latter, consider the topology on figure 2, and messages msg_A , msg_B , and msg_C , where msg_A has destination ι_3 and is on the out-queue of node ι_2 towards ι_3 , while msg_B has destination ι_1 and is on the out-queue of node ι_2 towards ι_1 , while msg_C has destination ι_2 and is on the out-queue of node ι_4 towards ι_3 . Then, we have that

$msg_A < msg_B$ and $msg_B < msg_C$, but we do not have that $msg_A < msg_C$.

We require that causality implies an earlier arrival – *c.f.*, the definition of consistent configurations in the next paragraph.

Consistent configurations We say that a configuration $(\mathcal{T}, \Delta, \Sigma)$ is *consistent* if it is well-formed (definition ??), and for any causally related messages m and m' m is guaranteed to arrive at its destination before m' arrives at the destination of the former.

Definition 3.8.

- Consistent $(\mathcal{T}, \Delta, \Sigma)$ iff
- WF $(\mathcal{T}, \Delta, \Sigma) \wedge$
 - $\forall m, m'. [m \ll_{\Delta} m' \rightarrow \mathcal{T}, \Delta, \Sigma \vdash m < m']$

4 Consistency implies Causality

Theorem ?? expresses the key property of our system: It guarantees that indeed in a tree topology messages are delivered in causal order. Namely, in a consistent runtime configuration, if after a number of execution steps, a message m is a cause of another message m' , and m' has already been delivered, and the two messages have the same destination, then m will also have been delivered.

Theorem 4.1 (Message delivery is causal).

- $$\begin{array}{l}
 \text{Consistent}(\mathcal{T}, \Delta, \Sigma) \quad \wedge \\
 \mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow^* \Delta', \Sigma' \quad \wedge \\
 m \ll_{\Delta'} m' \quad \wedge \\
 \text{Dest}(m) = \text{Dest}(m') \quad \wedge \\
 \text{Rec}(m) \in \Delta' \\
 \rightarrow \\
 \text{Rec}(m) <_{\Delta'} \text{Rec}(m')
 \end{array}$$

Proof sketch Theorem ?? follows directly from lemma ??, which guarantees that in a consistent configuration, messages which have the same destination will be delivered in causal order, and from lemma ??, which guarantees that execution preserves consistency. Both lemmas are stated below. In a companion technical report in preparation we have longer proofs – here we give proof sketches.

Lemma 4.2. Consistency implies Causality

- $$\begin{array}{l}
 \text{Consistent}(\mathcal{T}, \Delta, \Sigma) \quad \wedge \\
 \text{Dest}(m) = \text{Dest}(m') \quad \wedge \\
 m \ll_{\Delta} m' \quad \wedge \\
 \text{Rec}(m) \in \Delta \\
 \rightarrow \\
 \text{Rec}(m) <_{\Delta} \text{Rec}(m')
 \end{array}$$

Proof. By application of definition Consistent $(\mathcal{T}, \Delta, \Sigma)$, and because $m \ll_{\Delta} m'$, we know that $\mathcal{T}, \Delta, \Sigma \vdash m < m'$. Because $\text{Rec}(m) \in \Delta$, only case 5 of the definition of $\mathcal{T}, \Delta, \Sigma \vdash m < m'$ is applicable, and this gives that $\text{Rec}(m) <_{\Delta} \text{Rec}(m')$

□

We will now prove that execution preserves consistency of configurations. Lemma ?? guarantees that execution starting from a well-formed configuration leads to another well-formed configuration. Lemma ?? guarantees that execution starting from a well-formed configuration preserves the ordering of messages.

Lemma 4.3 (Execution preserves well-formedness).

$$\begin{aligned} & WF(\mathcal{T}, \Delta, \Sigma) \wedge \mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow \Delta', \Sigma' \\ & \rightarrow \\ & WF(\mathcal{T}, \Delta', \Sigma') \end{aligned}$$

Proof. By case analysis over $\mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow \Delta', \Sigma'$ and application of the definition of $WF(\mathcal{T}, \Delta, \Sigma)$ \square

Lemma 4.4 (Message order is preserved).

$$\begin{aligned} & WF(\mathcal{T}, \Delta, \Sigma) \wedge \mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow \Delta', \Sigma' \wedge \mathcal{T}, \Delta, \Sigma \vdash m < m' \\ & \rightarrow \\ & \mathcal{T}, \Delta', \Sigma' \vdash m < m' \end{aligned}$$

Proof. By case analysis over $\mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow \Delta', \Sigma'$. \square

Lemma ?? characterises the causal dependencies which can be derived from $\Delta \cdot e$ but not from Δ .

Lemma 4.5 (Extending causal dependencies).

$$\begin{aligned} & \neg(m \ll_{\Delta} m') \quad \wedge \\ & m \ll_{\Delta'} m' \quad \wedge \\ & \Delta' = \Delta \cdot e \quad \wedge \\ & \rightarrow \\ & [e = \text{Snd}(m', t) \vee e = \text{Rec}(m')] \quad \wedge \\ & [m <_{\Delta} m' \vee \exists m''. m \ll_{\Delta} m'' \wedge m'' <_{\Delta'} m'] \end{aligned}$$

Lemma ?? guarantees that any number of execution steps preserves consistency.

Lemma 4.6 (Multiple step execution preserves consistency).

$$\begin{aligned} & \text{Consistent}(\mathcal{T}, \Delta, \Sigma) \wedge \mathcal{T} \vdash \Delta, \Sigma \rightsquigarrow^* \Delta', \Sigma' \\ & \rightarrow \\ & \text{Consistent}(\mathcal{T}, \Delta', \Sigma') \end{aligned}$$

Proof. By induction on the number of steps in \rightsquigarrow^* . For 0 steps it is trivial, for 1 step by application of lemmas ??, ??, and ??, and for more than 1 step by straightforward induction. \square

5 Conclusion and Future Work

We have shown that distributed systems arranged in a tree topology are causal by construction. Specifically, we:

- Provided a mechanism to enforce causal ordering without the need for augmenting messages or storing meta data,
- Proposed a formal model to capture message passing in tree networks,
- Outlined proof sketch that causality is ensured.

In recent years we have worked on the development of an actor-based programming language called *Pony* [7, 9, 10, 25]. In the current distribution, Pony is concurrent, but our aim

is to develop it so as to provide *transparent* distributed programming: so that the notion of distribution is not exposed to the programmer. The result of this paper is central to the development of Distributed Pony: arranging the nodes into a tree topology,

In further work we want to study more dynamic distributed systems, which support actor migration, the addition (and potentially exchange) of nodes at runtime, as well as distributed snapshotting. A preliminary outlook on such mechanisms appeared in [7].

We believe that causal message delivery simplifies programming considerably, and thus our result is of wider importance. For instance, causal message delivery (along with types) has been leveraged for fully concurrent Garbage Collection [?].

On the other hand, enforcing causality by a topological arrangement of nodes is, at least for some applications, a case of trading space for time [21]. Therefore, we plan to also investigate the negative effects and runtime/space complexity tree topologies might have on messaging and routing performance, what possible optimizations are and how partial trees can be used to explicitly specify causal relationships to reduce the routing overhead on totally unrelated messages.

A significant open question is the handling of failures. What will recovery look like in a tree network? What protocols can be put in place to allow for topology changes? How does the ability to enforce causality without the ability to determine causal relations (as in our system), affect failure recovery? We plan to draw inspiration from distributed snapshots [4, 16].

References

- [1] M. Ahamad, P. W. Hutto, and R. John. 1991. Implementing and programming causal distributed shared memory. In *[1991] Proceedings. 11th International Conference on Distributed Computing Systems*. 274–281. <https://doi.org/10.1109/ICDCS.1991.148677>
- [2] Mustaque Ahamad, Gil Neiger, James E. Burns, Prince Kohli, and Phillip W. Hutto. 1995. Causal memory: definitions, implementation, and programming. *Distributed Computing* 9, 1 (01 Mar 1995), 37–49. <https://doi.org/10.1007/BF01784241>
- [3] Paulo Sérgio Almeida, Carlos Baquero, and Victor Fonte. 2008. *Interval Tree Clocks*.
- [4] Özalp Babaoğlu and Keith Marzullo. 1993. Distributed Systems (2Nd Ed.). ACM Press/Addison-Wesley Publishing Co., New York, NY, USA, Chapter Consistent Global States of Distributed Systems: Fundamental Concepts and Mechanisms, 55–96. <http://dl.acm.org/citation.cfm?id=302430.302434>
- [5] Peter Bailis, Alan Fekete, Ali Ghodsi, Joseph M. Hellerstein, and Ion Stoica. 2012. The Potential Dangers of Causal Consistency and an Explicit Solution. In *Proceedings of the Third ACM Symposium on Cloud Computing (SoCC '12)*. ACM, New York, NY, USA, Article 22, 7 pages. <https://doi.org/10.1145/2391229.2391251>
- [6] Kenneth Birman, André Schiper, and Pat Stephenson. 1991. Lightweight Causal and Atomic Group Multicast. *ACM Trans. Comput.* 9, 3 (Aug. 1991), 272–314. <https://doi.org/10.1145/128738.128742>
- [7] Sebastian Blessing. 2013. A String of Ponies - Transparent Distributed Programming with Actors. (2013). <https://www.doc.ic.ac.uk/~scb12/distinguished-projects/2013/s.blessing.pdf>

- [8] Bernadette Charron-Bost. 1991. Concerning the Size of Logical Clocks in Distributed Systems. *Inf. Process. Lett.* 39, 1 (July 1991), 11–16. [https://doi.org/10.1016/0020-0190\(91\)90055-M](https://doi.org/10.1016/0020-0190(91)90055-M)
- [9] Sylvan Clebsch and Sophia Drossopoulou. 2013. Fully Concurrent Garbage Collection of Actors on Many-core Machines. *SIGPLAN Not.* 48, 10 (Oct. 2013), 553–570. <https://doi.org/10.1145/2544173.2509557>
- [10] Sylvan Clebsch, Sophia Drossopoulou, Sebastian Blessing, and Andy McNeil. 2015. Deny Capabilities for Safe, Fast Actors. In *Proceedings of the 5th International Workshop on Programming Based on Actors, Agents, and Decentralized Control (AGERE! 2015)*. ACM, New York, NY, USA, 1–12. <https://doi.org/10.1145/2824815.2824816>
- [11] Giuseppe DeCandia, Deniz Hastorun, Madan Jambani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Voshall, and Werner Vogels. 2007. Dynamo: Amazon's Highly Available Key-value Store. *SIGOPS Oper. Syst. Rev.* 41, 6 (Oct. 2007), 205–220. <https://doi.org/10.1145/1323293.1294281>
- [12] C. J. Fidge. February, 1988. Timestamps In Message-Passing Systems That Preserve The Partial Ordering. (February, 1988).
- [13] Alexey Gotsman, Hongseok Yang, Marek Zawirski, and Sebastian Burckhardt. 2014. Replicated Data Types: Specification, Verification, Optimality. <https://www.microsoft.com/en-us/research/publication/replicated-data-types-specification-verification-optimality/>
- [14] M. P. Herlihy and J. M. Wing. 1987. Axioms for Concurrent Objects. In *Proceedings of the 14th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages (POPL '87)*. ACM, New York, NY, USA, 13–26. <https://doi.org/10.1145/41625.41627>
- [15] Leslie Lamport. 1978. Time, Clocks, and the Ordering of Events in a Distributed System. *Commun. ACM* 21, 7 (July 1978), 558–565. <https://doi.org/10.1145/359545.359563>
- [16] Friedemann Mattern. 1988. Virtual Time and Global States of Distributed Systems. In *Parallel and Distributed Algorithms*. North-Holland, 215–226.
- [17] D. S. Parker, G. J. Popek, G. Rudisin, A. Stoughton, B. J. Walker, E. Walton, J. M. Chow, D. Edwards, S. Kiser, and C. Kline. 1983. Detection of Mutual Inconsistency in Distributed Systems. *IEEE Trans. Softw. Eng.* 9, 3 (May 1983), 240–247. <https://doi.org/10.1109/TSE.1983.236733>
- [18] Nuno M. Prego, Carlos Baquero, Paulo Sérgio Almeida, Victor Fonte, and Ricardo Gonçalves. 2010. Dotted Version Vectors: Logical Clocks for Optimistic Replication. *CoRR* abs/1011.5808 (2010). <http://arxiv.org/abs/1011.5808>
- [19] M. Raynal and M. Singhal. 1996. Logical time: capturing causality in distributed systems. *Computer* 29, 2 (Feb 1996), 49–56. <https://doi.org/10.1109/2.485846>
- [20] Reinhard Schwarz and Friedemann Mattern. 1994. Detecting Causal Relationships in Distributed Computations: In Search of the Holy Grail. *Distrib. Comput.* 7, 3 (March 1994), 149–174. <https://doi.org/10.1007/BF02277859>
- [21] Mark Stamp. Accessed July 2017. Once Upon a Time-Memory Tradeoff. (Accessed July 2017). <http://www.cs.sjsu.edu/~stamp/RUA/TMTO.pdf>
- [22] Doug Terry. 2011. *Replicated Data Consistency Explained Through Baseball*. Technical Report. <https://www.microsoft.com/en-us/research/publication/replicated-data-consistency-explained-through-baseball/>
- [23] Douglas B. Terry, Alan J. Demers, Karin Petersen, Mike Spreitzer, Marvin Theimer, and Brent W. Welch. 1994. Session Guarantees for Weakly Consistent Replicated Data. In *Proceedings of the Third International Conference on Parallel and Distributed Information Systems (PDIS '94)*. IEEE Computer Society, Washington, DC, USA, 140–149. <http://dl.acm.org/citation.cfm?id=645792.668302>

- [24] D. B. Terry, M. M. Theimer, Karin Petersen, A. J. Demers, M. J. Spreitzer, and C. H. Hauser. 1995. Managing Update Conflicts in Bayou, a Weakly Connected Replicated Storage System. In *Proceedings of the Fifteenth ACM Symposium on Operating Systems Principles (SOSP '95)*. ACM, New York, NY, USA, 172–182. <https://doi.org/10.1145/224056.224070>
- [25] The Pony Core Development Team. 2013 - today. Pony. (2013 - today). <http://www.ponylang.org>
- [26] Robbert van Renesse. 1993. Causal Controversy at Le Mont St.-Michel. *Operating Systems Review* 27, 2 (1993), 44–53. <https://doi.org/10.1145/155848.155857>

6 Appendix

Dest	::=	Message → NodeId
In	::=	State × NodeId × NodeId → Message
Out	::=	State × NodeId × NodeId → Message
path	::=	NodeId × Message × \mathcal{T} → NodeId*
next	::=	NodeId × Message × \mathcal{T} → NodeId
Dest(m)	=	$m \downarrow_2$
In(Σ, t_1, t_2)	=	$\Sigma(t_1)(t_2) \downarrow_1$
Out(Σ, t_1, t_2)	=	$\Sigma(t_1)(t_2) \downarrow_2$
path(t, m, \mathcal{T})	=	path($t, \text{Dest}(m), \mathcal{T}$)
next(t, m, \mathcal{T})	=	next($t, \text{Dest}(m), \mathcal{T}$)

Figure 9. Shorthands

TopIn	::=	State × NodeId × NodeId → Message
TopOut	::=	State × NodeId × NodeId → Message
PushIn	::=	State × NodeId × NodeId × Message → Node
PushOut	::=	State × Tree × NodeId × Message → Node
PopIn	::=	State × NodeId × NodeId → Node
PopOut	::=	State × NodeId × NodeId → Node

Figure 10. Auxiliary functions - declarations

TopIn(Σ, t_1, t_2)	=	m if In(Σ, t_1, t_2) = $m \cdot _$
TopOut(Σ, t_1, t_2)	=	m if Out(Σ, t_1, t_2) = $m \cdot _$
PushIn(Σ, t_1, t_2, m)	=	$\Sigma(t_1)[t_2 \mapsto \text{In}(\Sigma, t_1, t_2) \cdot m, \text{Out}(\Sigma, t_1, t_2)]$
PushOut($\Sigma, \mathcal{T}, t, m$)	=	$\Sigma(t)[t' \mapsto \text{In}(\Sigma, t, t'), \text{Out}(\Sigma, t, t') \cdot m]$ where $t' = \text{next}(t, m, \mathcal{T})$
PopIn(Σ, t_1, t_2)	=	$\Sigma(t_1)[t_2 \mapsto (q, \text{Out}(\Sigma, t_1, t_2))]$ where In(Σ, t_1, t_2) = $m \cdot q$
PopOut(Σ, t_1, t_2)	=	$\Sigma(t_1)[t_2 \mapsto (\text{In}(\Sigma, t_1, t_2), q)]$ where Out(Σ, t_1, t_2) = $m \cdot q$

Figure 11. Auxiliary functions - definitions