

The hierarchical daisy architecture for causal delivery

This article has been downloaded from IOPscience. Please scroll down to see the full text article.

1999 Distrib. Syst. Engng. 6 71

(<http://iopscience.iop.org/0967-1846/6/2/302>)

View [the table of contents for this issue](#), or go to the [journal homepage](#) for more

Download details:

IP Address: 38.107.179.212

The article was downloaded on 20/02/2012 at 21:05

Please note that [terms and conditions apply](#).

The hierarchical daisy architecture for causal delivery*

Roberto Baldoni[†], Roberto Beraldi[†], Roy Friedman[‡] and Robbert van Renesse[§]

[†] Dipartimento di Informatica e Sistemistica, Università di Roma 'La Sapienza', Via Salaria 113, Roma, Italy

[‡] Department of Computer Science, Technion—Israel Institute of Technology, Haifa 32000, Israel

[§] Department of Computer Science, Cornell University, Ithaca, NY 14850, USA

E-mail: baldoni@dis.uniroma1.it, beraldi@dis.uniroma1.it, roy@cs.technion.ac.il and rvr@cs.cornell.edu

Received 15 October 1998

Abstract. In this paper, we propose the *hierarchical daisy architecture*, which provides causal delivery of messages sent to any subset of processes. The architecture provides fault tolerance and maintains the amount of control information within a reasonable size. It divides processes into *logical groups*. Messages inside a logical group are sent directly, while messages that need to cross logical groups' boundaries are forwarded by servers. We prove the correctness of the daisy architecture, discuss possible optimizations, and present simulation results.

1. Introduction

The asynchrony of communication channels is one of the major sources of nondeterminism in distributed systems, which in turn is a major cause of problems when implementing distributed applications. *Causal ordering* [4] reduces much of this asynchrony by guaranteeing that whenever a message is delivered to a process, all causally prior messages that were sent to the same process have already been delivered to it. This abstraction yields simplified solutions to many fundamental problems in distributed computing, such as *atomic snapshot* [2], *management of replicated data* [4], and *monitoring of distributed applications* [16].

Several protocols for implementing causal ordering have been proposed [4, 6, 18, 17]. These protocols mainly differ by the assumptions they make on the communication patterns, the topological structure of the underlying network, the amount of control information used to enforce causal ordering, and in how fast the protocol is in sending and delivering messages. For protocols that optimize delivery time, and do not make any assumptions on the topological structure of the network, the following bounds are known:

- If processes only broadcast messages to a group of size n , then the amount of control information that is required

to add to messages is $\Theta(n)$ [6].

- In order to support multicasts to a small number of groups g , each of size n , the amount of required control information is $\Theta(n \cdot g)$ [6].
- In order to allow each process to send each message to an arbitrary set of recipients, taken from a pool of n processes, the amount of required control information is $\Theta(n^2)$ [16].

The bound of $\Theta(n^2)$ for sending messages to arbitrary sets of recipients is prohibitively high for large values of n . Several suggestions have been made to reduce the amount of control information in the common case. For example, for protocols that use matrix timestamps, e.g., [18], it is possible to send only the difference between the previous matrix timestamp and the current one. It is assumed that in most cases the difference matrix would be very sparse, and therefore can be represented very efficiently. However, in the worst case, a full matrix must be sent, which means $O(n^2)$ integers. Prakash *et al* [15] proposed to piggyback on every message explicit information regarding all causally prior messages that are not known to have been delivered to all their destinations, under the assumption that usually there will not be too many such messages. In practice, it is unclear how large this information will be, and even using complex optimizations, it remains $O(n^2)$ in the worst case. Finally, Horus [19] provides the FILTER layer, which translates every send downcall into a broadcast; at the receivers' side, this layer filters out messages that are not intended for the local

|| This control information actually corresponds to a vector of timestamp [13].

* Based on 'The hierarchical daisy architecture for causal delivery' by Roberto Baldoni, Roy Friedman and Robbert van Renesse which appeared in Proceedings of 17th International Conference on Distributed Computing Systems (ICDCS '97); Baltimore, MD, May 27–30, 1997; pp 570–77. ©1997 IEEE.

process. This allows us to use only a single vector of size n , but is naturally only feasible for relatively small groups equipped with hardware multicast capability.

One of the main problems with most existing protocols for causal ordering is that *a single failure of a process combined with an omission of a single message can prevent the entire system from delivering any additional messages*[†]. Worse yet, such failure scenarios are not uncommon in real systems. To overcome this problem, most existing systems adopt conservative techniques that delay transmissions of some messages until previous messages have been stored by at least one additional process, so these messages can be retrieved if the sender fails. This is usually done in a manner unrelated to the control information, which prohibits exploiting these delays in order to reduce the amount of control information.

In this paper we propose a hierarchical architecture that attacks both the problem of reducing the amount of control information and the problem of fault-tolerance at the same time. Our solution splits the participating processes into local groups; each message sent to a member of a local group is translated into a broadcast message to all the members of the local group. The architecture utilizes *causal servers* to disseminate messages across local groups through the *causal servers group*, to which all causal servers belong. (i.e., each causal server is a member of both a local group and the causal servers group.) Causal servers add a delay to messages that need to cross local groups' boundaries, and broadcasting messages in the local group rather than sending them point-to-point only to the intended recipients may increase the total message count. However, these delays combined with the usage of intra-group broadcast communication guarantee fault-tolerance, since every member of the local group has a copy of the message that can be retrieved in case of a failure. Moreover, we exploit these facts to reduce the amount of control information added to messages to $O(m)$, where m is the size of the largest local group or causal server group in which a message traverses. This reduction in control information offsets the added message count and extra hops for large groups. Also, being hierarchical, this solution scales to large numbers of processes, while allowing processes to send each message to any arbitrary set of processes.

Our protocol utilizes a group communication system, e.g., Horus [19], ISIS [4], Transis [8], Totem [14], Phoenix [12], or Relacs [3]. The services of the group communication system we rely on are failure detection[‡] and automatic reconfiguration in the event of a failure or a join of a new member, reliable fifo point-to-point delivery, stability detection, and automatic re-issuing of messages from failed members that were received by only some members (but not by all of them). This functionality is supported by all the systems we have mentioned. Note that it is possible to implement these functions from scratch. However,

[†] Some papers assume reliable delivery of messages and justify this assumption by the use of a point-to-point reliable delivery protocol. As we explain in section 2.3, unless new messages are not sent until it is guaranteed that older messages have been received, this is not a valid assumption.

[‡] The basic mechanism for failure detection in most group communication systems consists of having all processes exchange 'I-am-alive' messages periodically such that if a process p does not receive an 'I-am-alive' message from a process q for sufficiently long, then p suspects q .

since these are well studied problems, assuming a group communication system that provides them simplifies the discussion, and allows us to concentrate on the new things in our architecture and protocol.

In a recent paper, Rodrigues and Verissimo describe an approach which is based on causal separators for reducing the amount of control information in systems that span several network domains [17]. Their approach, however, yields an architecture that is not hierarchical. In particular, it does not reduce the amount of control information used within the same subnet domain, even though current LANs can have as many as several hundreds of machines in the same subnet. Our solution does not assume anything about the network topology, although such knowledge can sometimes be used to improve the performance of the system, e.g., by mapping causal servers to routers. Also, Rodrigues and Verissimo do not address the issue of omission failures in [17], so their solution is less fault-tolerant than ours, as described in section 2.3.

Adly and Nagi proposed using a logical hierarchy in a causal delivery protocol that also provides fault tolerance and reduced message sizes [1]. Their solution, however, requires processes to log their state on stable storage, while our solution does not have this restriction. Moreover, by relying on a group communication service, our solution is simpler, while their protocol has to explicitly deal with problems like reliable delivery that are already provided by group communication systems.

Finally, IP-multicast uses overlapping groups for large scale dissemination of information. However, IP-multicast does not provide causal ordering, nor reliable delivery.

The rest of this paper is structured as follows. The model and main definitions are introduced in section 2. The architecture is presented in section 3 and is shown to be correct in section 4. We present simulations results in section 5, and conclude with a discussion in section 6.

2. Model and definitions

2.1. Asynchronous distributed systems

A distributed system consists of a finite set P of n processes, connected by a *communication network*, or simply *network*, and communicating with each other only by sending and receiving messages through the network. We assume that the network is well connected, but unreliable and asynchronous. That is, the network can delay messages for an arbitrarily long time and may occasionally drop a message completely. However, the probability of a message to be delivered after some finite time (but unknown) is nonzero while the sender is nonfaulty. In particular, each process can send a message to any set of processes, and this message will be delivered after some finite time to all its recipients with some positive probability, although each recipient may receive the message at a different time. Note that we allow the network to occasionally 'drop' a message; we refer to this as an *omission failure*. We assume that processes do not have access to a global clock and there is no bound on their relative speed. In addition to omission failures, processes may fail by crashing [7].

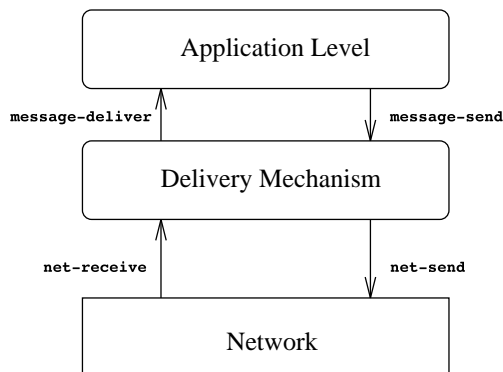


Figure 1. The structure of a process.

As is done in most papers about causal ordering, we assume that each process consists of an *application level* and a *delivery mechanism*, or *DM* for short, as illustrated in figure 1. The application level at each process can generate *message-send* events to the DM, and can accept *message-deliver* events from the DM. The DM can generate *net-send* events to the network, and accept *net-receive* events from the network. We also assume that the application level can accept a crash event, in which case it is the last event this application level receives.

We say that a message is *sent* when the corresponding *message-send* event is generated; a message is *received* when the corresponding *net-receive* is generated; a message is *delivered* when a *message-deliver* event is generated. We assume that the DM can generate *net-send* events only for messages it accepted *message-send* events for, and can generate *message-deliver* event only for messages it accepted *net-receive* events for. In case the DM generates one *net-send* event for several messages, also known as *piggybacking* or *packing*, we assume, for simplicity reasons, that each of these messages generates a separate *net-receive* event.

We denote a *message-send* event of process p_i sending message m to a subset P_k of processes by $send_i(m, P_k)$ and the *message-deliver* event of a message m at process p_j by $del_j(m)$.

An *execution* of a distributed system P is a collection of *message-send* and *message-deliver* events with the following partial order, also known as the *happens before* relation [11], defined on them as follows.

Definition 2.1 (happens before relation). An event $a \rightarrow b$ iff:

- (i) a and b are two events of the same process and a occurs before b .
- (ii) a is a *message-send* event $send_i(m, P_k)$, for some p_i, m , and P_k , and b is the corresponding *message-deliver* event $del_j(m)$, for some $p_j \in P_k$.
- (iii) There exists an event c such that $a \rightarrow c$ and $c \rightarrow b$.

We say that a process crashes in an execution σ , if the application level of this process accepts a *crash* event during σ . If neither $a \rightarrow b$ nor $b \rightarrow a$, a and b are said to be *concurrent events*.

Daisy architecture for causal delivery

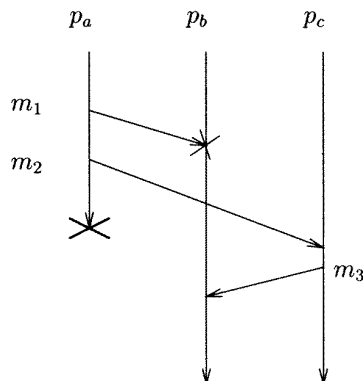


Figure 2. A failure scenario which causes causal protocols to block.

2.2. Causal ordering

Introduced by Birman and Joseph in [5], causal ordering states that the order in which messages are delivered to the application must be consistent with the *happened-before* relation of the corresponding sending events. More formally, we have the following definition.

Definition 2.2. An execution of a distributed system σ respects causal ordering if:

- (i) for any two messages m_1 and m_2 , sent by p_i and p_j , with the same destination p_k such that $send_i(m_1, P_\alpha) \rightarrow send_j(m_2, P_\beta)$ ($p_k \in P_\alpha \cap P_\beta$), $del_k(m_1) \rightarrow del_k(m_2)$;
- (ii) for any $send_i(m, P_\alpha)$ event of a process p_i that does not crash in σ and for every process $p_j \in P_\alpha$ that does not crash in σ , the application level of p_j in σ accepts an event $del_j(m)$ from its DM.

Hence, the problem of implementing causal ordering is the problem of designing a protocol for the DM that will always obey the requirements of definition 2.2. We say that a message that was received by the DM of a process p_j is *causally deliverable* if all causally prior messages that were sent to p_j have been delivered to the application level of p_j .

The main obstacles facing implementations of causal ordering are the amount of control information required to ensure causal ordering, and overcoming failures of processes and message omissions. The amount of control information used by the protocol is important for the scalability of the solution; if it is too large, then the protocol becomes infeasible for large groups. Also, as we discuss in the following section, if the protocol does not explicitly handle message omissions and crash failures, then a single failure scenario can block the protocol from delivering messages.

2.3. Fault tolerance

The problem of fault tolerance in causal ordering protocols that allow sending messages to overlapping, but different, sets of recipients has been pointed out by Birman *et al* [6].

Figure 2 illustrates an example of problems caused by a crash of a process combined with an omission failure. Message m_1 is lost by process p_b due to an omission failure. Assuming that the causal ordering protocol releases messages to the network as soon as it receives them from

the application, message m_2 is sent from p_a to p_c before p_a notices the omission of m_1 . Note that m_2 causally follows m_1 . As soon as m_2 is delivered to p_c , it sends message m_3 to p_b . Following this, p_a crashes before noticing that m_1 was lost. Now, when p_b receives m_3 , it cannot deliver it, since m_3 causally depends on m_1 , which has not arrived yet. On the other hand, since p_a has failed, there is no way to retrieve m_1 , and p_b is blocked forever from delivering p_c 's messages.

Note that this problem occurs even if processes employ a point-to-point reliable delivery mechanism, such as positive acknowledgements. For example, since m_1 is sent to p_b and m_2 is sent to p_c , such a mechanism would still allow m_2 to be sent and be delivered before the faith of m_1 is determined. The only known solutions to these problems are:

- (a) Wait for the stability of all previously sent and received messages whenever the subset of recipients changes. This is the solution currently used by ISIS, but it requires (i) delaying such messages for arbitrary long periods of time before sending them, and, (ii) an amount of control information of $O(n^2)$.
- (b) Broadcast every message to all processes, so there is always a way to retrieve m_1 . This solution is currently supported by Horus, but it may create too many messages in the network, and is therefore not feasible for large groups.
- (c) Piggyback on every message all previously unstable messages. This solution was used by early versions of ISIS. It does not impose delay to messages while sending. However, it may generate extremely large messages. Note that, given a message, all previously unstable messages represents actually the control information of that message.

Previous solutions show a tradeoff between the control information, the delay imposed on messages and the number of messages exchanged. For example, reducing the control information piggybacked by messages requires either to block messages, which adds a delay to messages while sending, as in solution (a), or to add messages to the execution as in solution (b). Our solution to this problem is described in the next section.

3. The proposed architecture

In this section, we describe a hierarchical architecture for implementing causal ordering by a collection of DMs. We start by describing the simple case, in which there are only two levels in the hierarchy, and then discuss how this can be generalized to more levels. Finally, we discuss several optimizations to the general architecture, which can optimize the performance of the system.

3.1. The base case

We divide the processes into *local groups*, such that each process is a member of only one local group, as illustrated in figure 3. Groups are managed by a group communication system, e.g., Horus. In each group, a single process is chosen (deterministically) to be the *causal server* for that group.

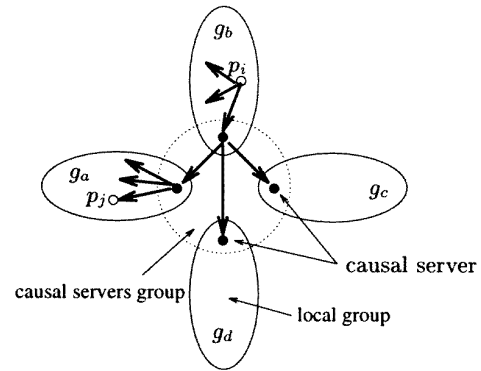


Figure 3. An example of a daisy configuration with a message sent from process p_i to p_j .

All causal servers are also members of another group simply called the causal servers group.

We refer to this as the *daisy configuration*. We assume that the group communication system can detect process failures, and report them to members of the group, and allow new members to join the group. We also assume that inside a group, the system provides both point-to-point and multicast fifo, reliable message delivery, for example, by employing a positive acknowledgement protocol or a negative acknowledgement protocol with periodic updates. Furthermore, we assume that (i) each message delivered within a group is buffered by the group communication system until it is known to have been received everywhere (in such a case we say that the message is *stable*), and that (ii) in the event of a failure, the system automatically relays omitted messages that were originated by failed members to members that have not received these messages. All these are standard features in all group communication systems we have mentioned.

We also assume a *state-transfer* mechanism, i.e., a mechanism that allows processes to exchange information during view changes, and guarantees that this information is delivered before any other message of a newly formed view. Finally, we assume that the system supports causal broadcast within a group, by employing, for example, the conventional causal broadcast protocol, which requires only a vector of integers of the size of the local group [18].

In order to send a message, a process does a causal broadcast of this message to all members of its local group, using the causal broadcast protocol we assumed. Each member that receives a message that is intended for it, delivers the message locally as soon as the message becomes causally deliverable. If the message is also intended for processes that are not members of the local group, the causal server waits until the message becomes causally deliverable in the local group, and then does a causal broadcast of this message in the causal servers group. Each causal server that receives such a message, waits until it becomes causally deliverable in the causal servers group, and then if some of the intended recipients of this message belong to its local group, does a causal broadcast of this message inside its local group. In this case, again, as soon as the message becomes causally deliverable, each process that was supposed to receive this message delivers it, and all other processes discard it[†].

[†] Note that the DM buffers messages in case they need to be retransmitted after a crash.

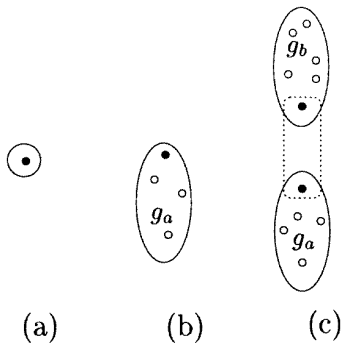


Figure 4. Evolving daisy configuration.

3.1.1. Evolving daisies. An important point to note about this architecture lies in the fact that it dynamically changes according to the number of processes running. For example, to get the daisy architecture depicted in figure 3, there are some intermediate steps to do which are depicted in figure 4. At the beginning there is just one causal server, i.e., a singleton group (figure 4(a)). At this point, new processes join the group g_a (figure 4(b)). When the control information becomes too large, or too many messages are received by everyone, a new causal server is created and the processes split into two local groups g_a and g_b (figure 4(c))[†]. If the number of processes in the system continues to grow, we get the configuration of figure 3. So, it turns out that in this architecture the number of server processes is usually much smaller than the number of processes in a group. If the number of processes decreases, then two groups can be merged back into a single group.

In this work we assume that the application is responsible for instructing the system when to split or merge local groups. In particular, during splits the application has to inform the DM which processes are to depart from the corresponding local group. In the discussion below, we call such processes *departing processes*, while the rest of the processes in the local group that splits are called *remaining processes*.

In order to ensure that all messages are delivered to all live processes during splits, we use the following scheme:

- (i) Once instructed to depart from a local group g_o , the departing processes stop sending new messages, and form a new local group g_d without leaving g_o (at this point).
- (ii) One of the departing processes, denoted by p_d is designated to be a causal server for g_d , and joins the causal servers group.
- (iii) After joining the causal servers group, p_d broadcasts a *group-init* message in p_d . Following this, p_d starts forwarding messages that were received in the causal servers group according to the same mechanism we described above. (i.e., only messages with intended recipients inside g_d that are causally deliverable in the causal servers group.)
- (iv) Every process in g_d that receives the *init-group* message in g_d , leaves g_o , and at this point is again allowed to send new messages, but only in g_d .

[†] New group memberships can be determined according to the amount of information exchanged among processes: if two processes communicate frequently they will become members of the same group.

- (v) If a departing process p_a joins g_d after p_d sent the *init-group* message, then p_d forwards to that process all messages that are still unstable in the causal servers group and are intended (also) to p_a , and a copy of *init-group*, as part of the state-transfer mechanism of view changes.

Note that during the transition phase, a departing process may receive a message twice, once in g_o and once in g_d . For this, we use a simple duplicate suppression mechanism that filters duplicate messages from being delivered twice to a process[‡]. Also note that (v) above guarantees that processes will not miss any message due to the asynchrony of message delivery in the causal servers group.

In the case of join, processes simply join an existing group, and after they have joined that group, they leave their original group. In this case too the state-transfer mechanism is also used to guarantee that joining processes do not lose messages due to this transfer.

3.1.2. Stability detection. Recall that for fault-tolerance purposes, processes buffer messages they send and receive. However, there is no point to continue buffering a message after it was delivered to all of its intended recipients (i.e., the message is stable). Thus, in order to keep the amount of local storage from overflowing, the causal delivery service must employ a stability detection mechanism.

A stability detection mechanism is usually based on the fact that processes inform each other of messages they receive, either by piggybacking this information on messages, or occasionally generating specific messages for this purpose. For efficiency purposes, many protocols employ ‘gossiping’ in the implementation of the stability detection protocol, which shortens the time required for a process to learn that a message is stable, and reduces the message overhead.

All group communication systems we have mentioned in section 1 provide a stability detection mechanism within a local group. Note, however, that our hierarchical mechanism requires stability information to traverse across groups, and that the sender of a message should not discard it until it is stable w.r.t. every recipient, and not just within the group. The stability detection mechanisms of some group communication systems (e.g., Horus) require the application at a process to explicitly inform the stability mechanism when it is OK to report that it has received the message in order to achieve end-to-end semantics. We therefore assume this capability as well, and describe how cross-group stability detection is implemented using the intra-group stability mechanism provided by the group communication system.

A causal server does not report a message as stable in the servers group, until it becomes stable in its local group. (If there are no recipients in the local group, the causal server

[‡] A duplicate suppression mechanism consists of adding a unique identifier to each message, which includes the originator’s unique identifier, and the sequential number of this message at this process. Since our protocol guarantee causal ordering, and since causal ordering also implies fifo ordering, each process needs to maintain a single vector with one entry for each process in the system. However, since this vector is kept in memory, and is not added to messages, it does not impose a scalability problem.

can immediately report the message as stable.) A causal server that belongs to a local group where a message was initiated, does not report the message as stable to its local group until it is stable in the causal servers group. This way, once a message is declared stable in the group that included its original sender, it is guaranteed to have arrived at all its destinations, and it is therefore safe to discard it.

3.1.3. Handling failures. Note that if a nonsender member of the group fails, then the group communication system takes care of resending all messages that originated by that member and were received by only part of the members of the same local group (but not by all of them). This, combined with the fact that every message is broadcast to all members in the local group, and the point-to-point reliable delivery mechanism of the group communication system, guarantee that such a failure will not block the protocol from making progress, as in the examples discussed in section 2.3. Thus, the only thing that need be addressed explicitly by our protocol is a failure of a causal server.

A failure of a causal server is problematic since it is a member of two groups, and messages in one group do not propagate automatically to the other by the group communication system. In particular, when a causal server fails, the local group it belongs to is reconfigured and another member is elected to be the causal server for that group. However, this server must be updated regarding the messages that were delivered in the servers group, but have not made it to the local group when the previous causal server crashed.

Recall that a causal server does not report a message as stable until it becomes stable in its local group. Thus, whenever a causal server fails, all messages that were still not reported by it as stable in its local group are recorded. When the new causal server joins the servers group, a state update message is sent to it, which includes copies of all such messages. The new causal server can then forward these messages in its local group. A duplicate suppression mechanism to the one discussed in section 3.1.1, is used to prevent duplicate messages from being delivered in the local group.

3.2. The general case

If we have to cope with a huge number of processes, the basic architecture becomes inadequate, since the size of the control information and the number of messages generated becomes very large. In this case, we can add another daisy configuration and divide processes between the two daisies according, for example, to their communication activity[†]. Processes that communicate rarely can be put on distinct daisies. At this point either one of the causal servers of one of the daisies or a new ‘*ad hoc*’ causal server becomes responsible for forwarding messages from one daisy to another. We call this server a *hyper-server*. All hyper-servers form a *hyper-servers group*. Figure 5 shows a system with three daisies. Each one has selected one causal server for inter-daisy communication.

[†] Processes can also be divided according to the physical topology of the communication network. Processes on the same LAN segment could be members of the same group.

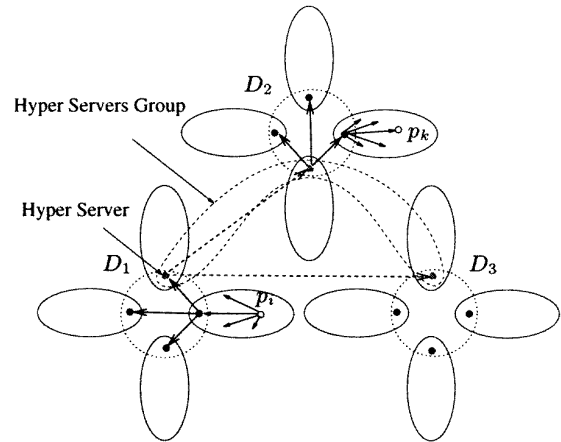


Figure 5. A system with three daisies in which a message is sent from process p_i to p_k .

In this general architecture, we assume that inter-daisy communication among causal servers is causally ordered. When a causal server receives a message m from p_i to be forwarded to a set of destinations in other daisies, it broadcasts m in the causal servers group of that daisy. The hyper-server of this daisy then broadcasts m in the hyper-servers group. Hyper-servers whose daisies include recipients of m broadcast m in the servers group of their daisy, so that causal servers in those daisies whose local groups include recipients of m can broadcast m there.

This can be done several times, to achieve several levels of hierarchy. Of course, in practice, we assume that 2–4 levels would be enough for most applications, since the number of levels grows logarithmically with the total number of processes.

3.3. Performance issues

The number of messages required to deliver a message depends on the number of different groups in which there are recipients; we need one broadcast in each group where there are recipients, one broadcast in the group that the originator of the message belongs to, and one broadcast for each servers group the message has to go through. Using packing techniques, as described in [9], it is possible to greatly reduce the number of actual (physical) messages being sent.

Since within a group we are using the causal broadcast protocol presented in [18], a message always carries a single vector of length m , where m is the size of the (local or servers) group the message is being sent in.

3.4. Possible optimizations

We now outline several optimizations that can reduce either the amount of control information, the number of messages, or the delays. However, each of these optimizations involves a tradeoff. Decreasing the amount of control information and the number of messages yield longer delays for messages, and a lower degree of fault-tolerance. On the other hand, reducing message delays requires more control information and slightly more complex rules for delivering messages.

- (i) A causal server that receives a message from its local group can forward the message to the servers group as soon as it receives the message, and not wait for it to become causally deliverable in the local group. Similarly, a causal server that receives a message in the servers group can forward it immediately inside its local group. However, this requires the message to carry one vector for each group it traverses in, and requires complex rules for deciding when a message becomes causally deliverable at its recipients.
 - (ii) It is possible to send messages inside a group only to their intended recipients, instead of broadcasting them to the entire group. This requires the use of matrices of size m^2 for groups of size m , instead of just a vector of length m . Also, in order to have some degree of fault tolerance, messages cannot be delivered until it is known that they were received by at least k more processes, in order to be able to survive k failures.
 - (iii) It is possible to send all messages inside a local group point-to-point to the causal server, and have the causal server forward the messages to all recipients within the group on a point-to-point basis. (Recall that we assumed that the group communication system employs a reliable fifo mechanism, such as negative acknowledgements, for messages sent within the group.) Forwarding message in the servers group occurs in this case using the regular causal broadcast algorithm.
- This approach reduces the number of messages required to send a message, and the amount of control information carried on messages, since delivering point-to-point messages reliably and in fifo order requires only a single integer to be added to each message. The disadvantage of this approach is that messages sent within the group need to pass through an extra hop. Also, this creates an additional load on the causal server, since it has to forward all messages sent in the local group to their recipients inside this group.

Note that all these ideas can also be implemented at higher levels of the hierarchy. The same trade-offs apply at these levels as well.

4. Correctness

Here we only show that this architecture guarantees causal delivery (*safety property*), that each message received will be eventually delivered (*liveness property*), and that all messages are eventually delivered once to every intended recipient (*reliability*) in a system composed of a two-level daisy. The proof for higher levels is a trivial extension of the proof we give for two levels, and is therefore omitted.

4.1. Safety

Note that omission failures and crash failures cannot alter the safety of our architecture. They can prevent the protocol from delivering messages, but this is a liveness issue. Liveness is shown in the next section, so in the rest of this section we can ignore failures.

Theorem 4.1. *The protocol presented in section 3 delivers messages in causal order.*

Proof. Assume, by way of contradiction, that there exists an execution σ of the protocol in which some messages are delivered in an order that violates causal ordering. Let m and m' be two such messages. Denote by p_i the process that sends m , p_j the process that sends m' , g_a the local group that p_i is a member of, and g_b the local group that p_j is a member of. Assume, w.l.o.g., that $send_i(m, P_k) \rightarrow send_j(m', P_{k'})$, but $del_q(m') \rightarrow del_q(m)$ for some $p_q \in P_k \cap P_{k'}$.

Clearly, this cannot happen if $g_a = g_b$. Thus, we must assume that $g_a \neq g_b$. Since $send_i(m, P_k) \rightarrow send_j(m', P_{k'})$, m is causally deliverable at the causal server of g_b , denoted by p_b , by the time it forwards m' to the causal servers group. Thus, in the causal servers group, m' gets a vector timestamp that indicates that it causally follows m .

Let g_c be the group of p_q , and denote by p_c the causal server of p_q . Thus, p_q sends m in g_c before it sends m' in g_c , meaning that the vector timestamp of m in g_c indicates that it has to be delivered after m' . Therefore, p_q cannot deliver m' unless it has delivered m , a contradiction to the assumption that $del_q(m') \rightarrow del_q(m)$. \square

4.2. Liveness

Before proving liveness, let us introduce the notion of *persistent causal hole* (PCH). By the definition of causal order, a message cannot be delivered to a destination p_i until all causally prior messages also sent to p_i have been delivered there. Hence, if a message m sent to a process p_i is lost and cannot be retrieved somewhere in the distributed system, all messages m' sent to p_i such that $m \rightarrow m'$ cannot be delivered to p_i . Due to the acyclicity of the happened-before relation, we assume, w.l.o.g., that all messages that are causally prior to m and sent to p_i have been delivered by p_i . In this case, we say that m creates a PCH in p_i . Let us introduce the following lemmas.

Lemma 4.2. *A group communication system (e.g., HORUS) and our protocol ensure that in a single group of processes there cannot be a PCH.*

Proof. Assume, by way of contradiction, that a message m , sent by p_j to p_i , creates a PCH in some execution of the protocol. According to the protocol of section 3.1, each message is actually sent to each member of the local group, and every such member buffers m until it becomes stable (section 3.1.2). By definition of stability, if m creates a PCH, it cannot be stable.

In the absence of crash failures, the liveness of the causal broadcasting protocol adopted [18] and the reliable fifo point-to-point delivery mechanism of the group communication system guarantees that m will eventually be delivered to p_i . If some process crashes, the group communication system automatically relays all unstable messages among the nonfaulty processes. Therefore, m will be eventually received by p_i . A contradiction to the assumption that m creates a PCH. \square

Lemma 4.3. *A group communication system (e.g., HORUS) and our protocol ensure that there cannot be a PCH in a group of processes g_a and in the associated causal server group g_c due to a message sent by g_a (resp. g_c) and to be relayed to g_c (resp. g_a).*

Proof. Assume, by way of contradiction, that there is a PCH in one of two groups due to a message that should have been relayed from g_a (resp. g_c) to g_c (resp. g_a), and denote by m the message that caused this PCH. In the absence of a crash failure of the causal server of g_a , the liveness of the causal broadcasting protocol adopted [18], lemma 4.2, the reliable fifo point-to-point delivery mechanism of the group communication system, and the relay mechanism done by the causal server of g_a (section 3.1) guarantees that any message m broadcast in the group g_a (resp. g_c) and with some recipient out of g_a (resp. in g_a) will be eventually received by any process of g_c (resp. g_a).

Therefore, we must assume that p_a , the causal server of g_a , crashes, m is destined to some recipients outside of g_a (resp. inside of g_a) and was not yet relayed by p_a to g_c (resp. g_a). Hence, the intergroup stability mechanism of section 3.1.2 considers m to be unstable. According to our protocol (section 3.1.3), another member of g_a , p'_a is elected (in a deterministic way) to be g_a 's new causal server. Due to lemma 4.2, p'_a will eventually receive all unstable messages of g_a . Similarly, as soon as p'_a is elected, it will join g_c 's group and will receive all unstable messages in g_c , and will relay unstable messages from g_a to g_c and vice versa.

Thus, m will eventually be received by processes in group g_a (resp. g_c), a contradiction to the assumption that m causes a PCH. \square

Theorem 4.4. *Every message that is received by the DM at a process p_a that never crashes, and is intended for p_a is eventually delivered by the DM of that process.*

Proof. Since we have shown that the causality relation among messages is guaranteed by our architecture (theorem 4.1) and we use a causal ordering protocol in each group that has been proven to be live in [18], we will prove liveness of our protocol showing that there cannot be a PCH in any execution of our protocol.

Hence, assume, by way of contradiction, that there exists a message sent by a member of group g_a that creates a PCH in a member of group g_b , and let g_c be the causal server group. Two cases are possible:

- $g_a = g_b$. Lemma 4.2 rules out a PCH in this case.
- $g_a \neq g_b$ (it is possible that $g_b = g_c$). By lemma 4.2 and by applying lemma 4.3 repeatedly, first, to the pair g_a and g_c , and then to the pair of groups g_c and g_b , there cannot be a PCH in this case either.

Therefore, in both cases there cannot be a PCH, a contradiction to the assumption that m causes a PCH. \square

4.3. Reliability

Reliable delivery of messages sent within a local group or a causal servers group is guaranteed by the use of group communication. As we indicated in the discussion above, we employ an additional duplicate suppression mechanism to ensure that duplicates that might occur during splits, joins, or failures are eliminated.

The only thing that needs to be proved is that processes do not lose messages during splits or joins. A process p_a

might lose a message m during a split or join if m is delivered in the new group g_d before p_a joins it, and is delivered in the original group g_o after p_a has left it. However, since p_a first join g_d , then by the assumption that m was not delivered in g_o , m is not stable in the causal servers group. Thus, our stability detection scheme guarantees that p_d (as a member of the causal servers group) holds a copy of m and will therefore send a copy of m to p_a as part of the state-transfer mechanism. Thus, p_a will receive m . Hence, we have the following theorem.

Theorem 4.5. *Every message that is sent by a process that does not crash is delivered to all its recipients that do not crash.*

5. Simulations

Causal ordering can be employed in a large number of network topologies, with variable number of processes, and each of these configurations is likely to yield somewhat different results. On the other hand, it is impractical to simulate every possible configuration. Since we are mainly targeting large scale communications, we have eventually decided to simulate what we consider a typical setting in which wide area causal delivery could be employed, and can still be simulated with the computer power we had access too. In this setting we compare the daisy architecture with a flat one.

5.1. The simulation model

The system. We assume a system of 60 processes located in a small geographic area such as a campus. Processes within the same local group are connected by a high-speed switched LAN. Causal servers groups are connected using a building backbone, such that each causal server is connected to both a LAN for its local group and to the corresponding causal servers' backbone. A campus backbone is used to connect the *hyper causal servers*. Thus, we have a three-level hierarchy.

In our simulations, each process generates messages according to a Poisson process with a mean rate of p messages per second. This means no correlation between messages and no bursty message traffic.

As for stability information, we assume that each time a message is received by a destination, a reply message is immediately generated. A more elaborate study of stability detection in large area settings is reported in [10].

The network. To model the network layer we define the message latency, T , as the time elapsed between the `net-send` event of a message m and the m 's `net-receive` event (see section 2.1). Each *network packet* consists of a control information part and a data part; the latter is provided by the application. The length of the control information field depends on the architecture used and on the size of the system while the length of the data field is variable with mean 300 bytes.

T is modelled as an exponentially distributed random variable with mean s plus a constant value proportional to the ratio between the control information size and the network

speed where s is 2 ms and 10 ms in the case of building backbone and of a campus backbone, respectively. Moreover, for simplicity, we assume (i) T does not depend on the network load, and, (ii) switched LANs incur only a negligible constant delay (i.e., we do not model collisions to media access), and do not drop messages. In particular, this means that they do not introduce any message reordering.

We further assume that messages are not fragmented, so each application message corresponds to a network packet. We therefore use the terms ‘message’ and ‘packet’ interchangeably. In order to point out the behaviour of the architectures during simulations we do not consider the effects of a transport layer implementation.

The daisy architecture. We study the performance of a daisy architecture as shown in figure 5. Our system consists of 60 processes, such that there are three hyper-servers in the system, four causal servers, and each local group contains four processes. The control information associated with a message is, in the worst case, a vector of five integers whose size, assuming 32 bits integers, is 20 bytes.

Each time a message sent by a causal server is received by another causal server, it must be resequenced before being delivered. We denote as *delivery-delay* the average time elapsed between the generation of a *net-receive* event and the generation of the corresponding *message-deliver* event.

The *end-to-end delay* is the time elapsed between the generation of the *message-send* event of a message m and the time of the generation of the corresponding *message-deliver* at the destination. This time includes, in the worst case, three distinct delivery delays. The sender *stability delay* is the average time elapsed between the generation of a *net-send* event and the reception of the stability information of all the other members of the group.

The flat architecture. In the flat architecture all processes form a single group. Each process in the flat group sends messages to all other processes, and associates with the message a set of recipients. Upon receiving a message, a process filters out messages that are not intended for it. With this scheme, the control information added to application messages in order to ensure causal delivery consists of a vector of 60 integers plus a vector of 60 boolean bits that define who should receive the message[†].

The path followed by messages in the flat architecture is similar to the one of the daisy architecture, by replacing causal servers with simple routers. A message received by a router is immediately forwarded on the other network. Hence, in this case, a message is only resequenced at the destination site.

The *end-to-end delivery* delay is given by the sum of the network delays (the delay in the campus backbone plus the delay in the building backbones) and the delivery delay at the destination site. The sender *stability delay* is given by the time elapsed from the *net-send* event and the receipt of all corresponding acknowledgements.

[†] Note that in this architecture, even messages exchanged solely by processes on the same LAN must be sent and acknowledged by all processes in the system.

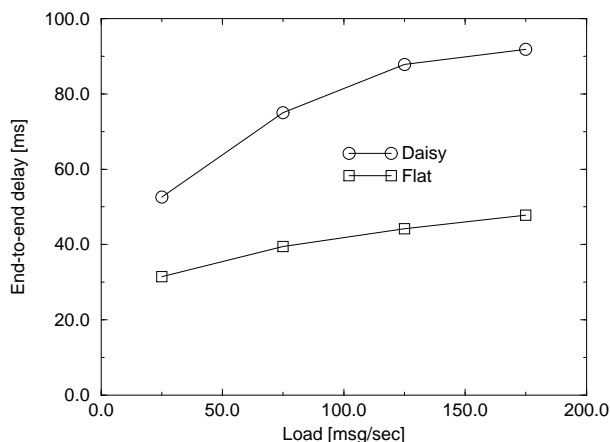


Figure 6. End-to-end delay versus load.

5.2. Simulation results

During the simulations, each run was repeated several times, each of which was executing until at least 1000 messages were delivered by all processes. A 95% confidence interval was computed on the results using the t-Student distribution. (Intervals are not reported in the figures since they are less than four per cent of the sample mean.)

In this work we only study a *broadcast setting*, in which each application message is sent to all processes. We would like to point out that this is a worst case scenario for the daisy architecture, since otherwise, causal servers do not forward messages to local groups that have no recipients among them. On the contrary, for the flat architecture this scenario is the same as any other scenario, since it always forwards all messages to all processes.

End-to-end delay. In figure 6 we consider the end-to-end delay in the case where a message experiences three hops to arrive to its destination, which is the worst case scenario. The end-to-end delay in the flat architecture is lower than in the daisy one. This is not surprising considering the fact that in the daisy architecture there are three distinct and successive reorderings due to causal servers. On the other hand, our system does not take into account collisions. It is reasonable to assume that if a nonswitched network will be used, the actual latency of the flat architecture will be higher with respect to the daisy one, since it generates more messages that would cause more collisions.

Delivery buffers. Figure 7 illustrates the average resequential buffer size as a function of the load. It is interesting to remark that the flat architecture requires more memory than the daisy one. This is mainly due to the size of the control information associated with each message, which is much bigger than what is required by the daisy architecture.

Stability delay. The stability delay in both architectures as a function of the load is reported in figure 8. As explained above, in the flat architecture, the receiver immediately sends back an acknowledgement, without waiting for the delivery

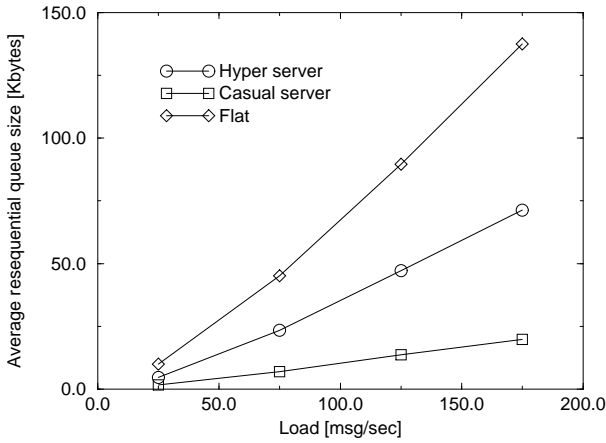


Figure 7. Resequential buffer versus load.

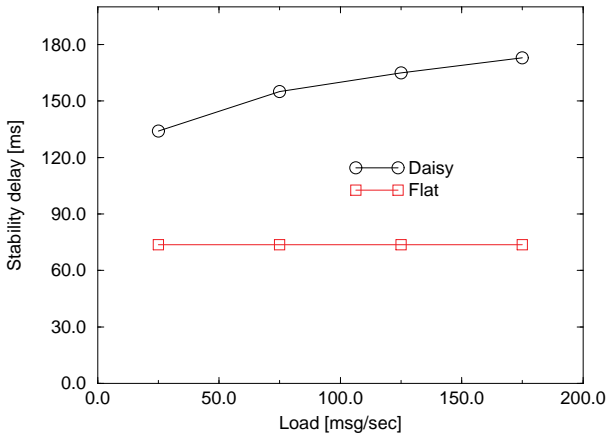


Figure 8. Stability delay versus load.

of the message. Since the network latency is not influenced by the traffic, the stability delay is constant. In the daisy architecture, the stability delay is influenced by the delivery delay of the servers (due to re-ordering) as opposed to the flat architecture.

Stability buffer. The plots of the storage requirements caused by the need to buffer unstable messages are reported in figure 9.

Note that the broadcast setting is also a worst case scenario for stability delay in the daisy architecture. In fact, if there is no recipient in a local group of messages, the causal server immediately declares the message stable, reducing the stability delay and hence the buffering requirements.

Final remark. As we mentioned above, simulating all possible scenarios is impractical. Instead, we have tried to pick a typical setting in which large scale causal ordering may be used, but one that also fits within our computing resources. Even though the scenarios we were simulating were a worst case when comparing our architecture with the flat one, our measurements indicate that the proposed architecture is reasonable in terms of end-to-end delivery delay and storage requirements. While the end-to-end delay is somewhat higher, the storage requirements and the reduced

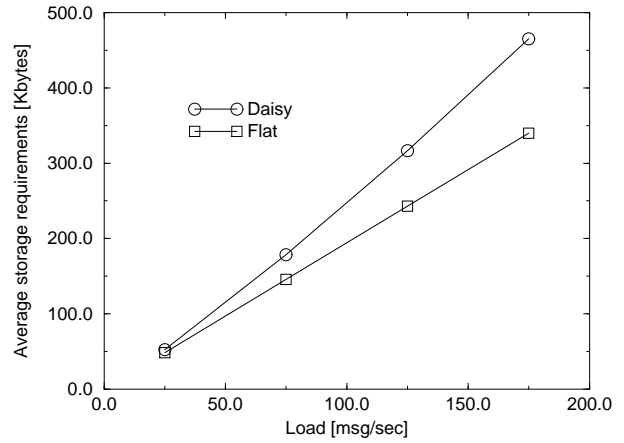


Figure 9. Storage requirements to buffer unstable messages.

control information and buffering requirements make the daisy architecture an attractive alternative to the flat one. This is true in particular when considering the fact that these properties only accentuate as the number of processes increase.

6. Conclusions

In this paper we have described the hierarchical daisy architecture for fault-tolerant causal ordering. This architecture can tolerate both omission and crash failures, while keeping the amount of control information small. In our solution, messages may need to go through one or more intermediate nodes on their way to their destination, but similar delays are needed in any case in order to guarantee fault-tolerance.

A smart utilization of our architecture can make it well suited for large area networks, by mapping processes in the same (or close) domains to the same daisies, and members of distinct domains to separate daisies. This is just a performance optimization, and is not required for the correctness of the proposed implementation.

Finally, while more elaborate simulations and an actual experience from a real implementation are the only way to determine the usefulness of a proposed architecture and its possible optimizations, our initial study indicates that the daisy architecture is a promising one.

Acknowledgments

The authors would like to thank the anonymous reviewers for their comments that improved the content and the presentation of the paper. This work was supported by DARPA/ONR grant N00014-96-1-1014.

References

- [1] Adly N and Nagi M 1995 Maintaining causal order in large scale distributed systems using a logical hierarchy *Proc. 12th IASTED Int. Conf. on Applied Informatics* pp 214–19
- [2] Alagar S and Venkatesan S 1994 An optimal algorithm for distributed snapshots with causal message ordering *Inf. Process. Lett.* **50** 311–16

- [3] Babaoğlu Ö, Davoli R, Giachini L and Baker M 1994 Relacs: a communication infrastructure for constructing reliable applications in large-scale distributed systems *Technical Report UBLCS-94-15*, Department of Computer Science, University of Bologna
- [4] Birman K and Joseph T 1987 Exploiting virtual synchrony in distributed systems *Proc. 11th ACM Symp. on Operating Systems Principles* pp 123–138
- [5] Birman K and Joseph T 1987 Reliable communication in the presence of failures *ACM Trans. Comput. Syst.* **5** 47–76
- [6] Birman K, Schiper A and Stephenson P 1991 Lightweight causal and atomic group multicast *ACM Trans. Comput. Syst.* **9** 272–314
- [7] Cristian F, Aghili H, Strong R and Dolev D 1985 Atomic broadcast: from simple diffusion to byzantine agreement *Proc. 15th Int. Conf. on Fault-Tolerant Computing (Austin, TX)*
- [8] Dolev D and Malki D 1996 The Transis approach to high availability cluster communication *Commun. ACM* **39** 64–70
- [9] Friedman R and van Renesse R 1996 Strong and weak virtual synchrony in Horus *Proc. 15th Symp. on Reliable Distributed Systems* pp 140–9
- [10] Guo K, van Renesse R, Vogels W and Birman K 1997 Hierarchical message stability tracing protocols *Technical Report TR-97-1647*, Department of Computer Science, Cornell University
- [11] Lamport L 1978 Time, clocks and the ordering of events in a distributed system *Commun. ACM* **21** 558–65
- [12] Malloth C, Felber P, Schiper A and Wilhelm U 1995 Phoenix: a toolkit for building fault-tolerant distributed application in large scale *Technical Report* Department d'Informatique, Ecole Polytechnique Federale de Lausanne
- [13] Mattern F 1989 Virtual time and global states of distributed systems in *Proc. Int. Workshop on Parallel and Distributed Algorithms* pp 215–26
- [14] Moser L, Melliar-Smith P M, Agarwal D, Budhia R and Lingley-Papadopoulos C 1996 Totem: a fault-tolerant multicast group communication system *Commun. ACM* **39** 54–63
- [15] Prakash R, Raynal M and Singhal M 1997 An efficient causal ordering algorithm for mobile computing environment *J. Parallel Distrib. Comput.* **41** 190–204
- [16] Raynal M, Schiper A and Toueg S 1991 The causal ordering abstraction and a simple way to implement it *Inf. Process. Lett.* **39** 343–50
- [17] Rodrigues L and Verissimo P 1995 Causal separators and topological timestamping: an approach to support causal multicast in large-scale systems *Proc. 15th Int. Conf. on Distributed Systems*
- [18] Schwarz R and Mattern F 1994 Detecting causal relations in distributed computing: in search of the holy grail *Distrib. Comput.* **7** 149–74
- [19] van Renesse R, Birman K and Maffeis S 1996 Horus: a flexible group communication system *Commun. ACM* **39** 76–83