

From Total Order to Database Replication

Yair Amir and Ciprian Tutu
Johns Hopkins University
Department of Computer Science
3400 N.Charles Street, Baltimore, MD 21218
{yairamir, ciprian}@cnds.jhu.edu

Abstract

This paper presents in detail an efficient and provably correct algorithm for database replication over partitionable networks. Our algorithm avoids the need for end-to-end acknowledgments for each action while supporting network partitions and merges and allowing dynamic instantiation of new replicas. One round of end-to-end acknowledgments is required only upon a membership change event such as a network partition. New actions may be introduced to the system at any point, not only while in a primary component. We show how performance can be further improved for applications that allow relaxation of consistency requirements. We provide experimental results that demonstrate the efficiency of our approach.

1 Introduction

Database replication is quickly becoming a critical tool for providing high availability, survivability and high performance for database applications. However, to provide useful replication one has to solve the non-trivial problem of maintaining data consistency between all the replicas.

The state machine approach [25] to database replication ensures that replicated databases that start consistent will remain consistent as long as they apply the same deterministic actions (transactions) in the same order. Thus, the database replication problem is reduced to the problem of constructing a global persistent consistent order of actions. This is often mistakenly considered easy to achieve using the Total Order service (e.g. ABCAST, Agreed order, etc) provided by group communication systems.

Early models of group communication, such as Virtual Synchrony, did not support network partitions and merges. The only failures tolerated by these models were process crashes, without recovery. Under this model, total order is sufficient to create global persistent consistent order.

When network partitions are possible, total order service does not directly translate to a global persistent consistent order. Existing solutions that provide active replication either avoid dealing with network partitions [27, 23, 22] or require additional end-to-end acknowledgements for every action after it is delivered by the group communication and before it is admitted to the global consistent persistent order (and can be applied to the database) [16, 12, 26].

In this paper we describe a complete and provably correct algorithm that provides global persistent consistent order in a partitionable environment without the need for end-to-end acknowledgments on a per action basis. In our approach, end-to-end acknowledgments are only used once for every network connectivity change event (such as network partition or merge) and not per action. The basic concept was first introduced as part of a PhD thesis [2]. This paper presents our newly developed insight into the problem and goes beyond [2] by supporting online additions of completely new replicas and complete removals of existing replicas while the system executes.

Our algorithm builds a generic replication engine which runs outside the database and can be seamlessly integrated with existing databases and applications. The replication engine supports various semantic models, relaxing or enforcing the consistency constraints as needed by the application. We implemented the replication engine on top of the Spread toolkit [4] and provide experimental performance results, comparing the throughput and latency of the global consistent persistent order using our algorithm, the COREL algorithm introduced in [16], and a two-phase commit algorithm. These results demonstrate the impact of eliminating the end-to-end acknowledgments on a per-action basis.

The rest of the paper is organized as follows. The following subsection discusses related work. Section 2 describes the working model. Section 3 introduces a conceptual solution. Section 4 addresses the problems exhibited by the conceptual solution in a partitionable system and introduces the Extended Virtual Synchrony model as a tool to provide global persistent order. Section 5 describes the de-

tailed replication algorithm and extends it to support online removals and additions to the set of participating replicas. Section 6 shows how the global persistent order guarantees of the algorithm can be used to support various relaxed consistency requirements. Section 7 evaluates the performance of our prototype and Section 8 concludes the paper.

1.1 Related Work

Two-phase commit protocols [12] remain the main technique used to provide a consistent view in a distributed replicated database system over an unreliable network. These protocols impose a substantial communication cost on each transaction and may require the full connectivity of all replicas to recover from some fault scenarios. Three-phase-commit protocols [26, 17] overcome some of the availability problems of two-phase-commit protocols, paying the price of an additional communication round.

Some protocols optimize for specific cases: limiting the transactional model to commutative transactions [24]; giving special weight to a specific processor or transaction [28]. Explicit use of timestamps enables other protocols [6] to avoid the need to claim locks or to enforce a global total order on actions, while other solutions settle for relaxed consistency criteria [11]. Various groups investigated methods to implement efficient lazy replication algorithms by using epidemic propagation [8, 14] or by exploiting application semantics [21].

Atomic Broadcast [13] in the context of Virtual Synchrony [7] emerged as a promising tool to solve the replication problem. Several algorithms were introduced [27, 23] to implement replication solutions based on total ordering. All these approaches, however, work only in the context of non-partitionable environments.

Keidar [16] uses the Extended Virtual Synchrony (EVS) [20] model to propose an algorithm that supports network partitions and merges. The algorithm requires that each transaction message is end-to-end acknowledged, even when failures are not present, thus increasing the latency of the protocol. In section 7 we demonstrate the impact of these end-to-end acknowledgements on performance by comparing this algorithm with ours. Fekete, Lynch and Shvartsman [9] study both [16] and [2] (which is our static algorithm) to propose an algorithm that translates View Synchrony, another specification of a partitionable group service defined in the same work, into a global total order.

Kemme, Bartoli and Babaoglu [19] study the problem of online reconfiguration of a replicated system in the presence of network events, which is an important building block for a replication algorithm. They propose various useful solutions to performing the database transfer to a joining site and provide a high-level description of an online reconfiguration method based on Enriched Virtual Synchrony allow-

ing new replicas to join the system if they are connected with the primary component. Our solution can leverage from these database transfer techniques and adds the ability to allow new sites to join the running system without the need to be connected to the primary component.

Kemme and Alonso [18] present and prove the correctness for a family of replication protocols that support different application semantics. The protocols are introduced in a failure-free environment and then enhanced to support server crashes and recoveries. The model does not allow network partitions, always assuming that disconnected sites have crashed. In their model, the replication protocols rely on external view-change protocols that provide uniform reliable delivery in order to provide consistency across all sites. Our work shows that the transition from the group communication uniform delivery notification to the strict database consistency is not trivial, we provide a detailed algorithm for this purpose and prove its correctness.

2 System Model

The system consists of a set of nodes (servers) $S = \{S_1, S_2, \dots, S_n\}$, each holding a copy of the entire database. Initially we assume that the set S is fixed and known in advance. Later, in Section 5.1, we will show how to deal with online changes to the set of potential replicas¹.

2.1 Failure and Communication Model

The nodes communicate by exchanging messages. The messages can be lost, servers may crash and network partitions may occur. We assume no message corruption and no Byzantine faults.

A server that crashes may subsequently recover retaining its old identifier and stable storage. Each node executes several processes: a database server, a replication engine and a group communication layer. The crash of any of the components running on a node will be detected by the other components and treated as a global node crash.

The network may partition into a finite number of disconnected components. Nodes situated in different components cannot exchange messages, while those situated in the same component can continue communicating. Two or more components may subsequently merge to form a larger component.

We employ the services of a *group communication layer* which provides reliable multicast messaging with ordering guarantees (FIFO, causal, total order). The group communication system also provides a membership notification service, informing the replication engine about the nodes that

¹Note that these are changes to the system setup, not view changes caused by temporary network events.

can be reached in the current component. The notification occurs each time a connectivity change, a server crash or recovery, or a voluntary join/leave occurs. The set of participants that can be reached by a server at a given moment in time is called a *view*. The replication layer handles the server crashes and network partitions using the notifications provided by the group communication. The basic property provided by the group communication system is called Virtual Synchrony [7] and it guarantees that processes moving together from one view to another deliver the same (ordered) set of messages in the former view.

2.2 Service Model

A *Database* is a collection of organized, related data. Clients access the data by submitting *transactions*, consisting of a set of commands that follow the ACID properties. A replication service maintains a replicated database in a distributed environment. Each server from the server set maintains a private copy of the database. The initial state of the database is identical at all servers. Several models of consistency can be defined, the strictest of which is *one-copy serializability* that requires that the concurrent execution of transactions on a replicated data set is equivalent to a serial execution on a non-replicated data set. We focus on enforcing the strict consistency model, but we also support weaker models (see Section 6).

An *action* defines a transition from the current state of the database to the next state; the next state is completely determined by the current state and the action. We view actions as having a *query* part and an *update* part, either of which can be missing. Client *transactions* translate into actions that are applied to the database. The basic model best fits one-operation transactions, but in Section 6 we show how other transaction types can also be supported.

3 Replication Algorithm

In the presence of network partitions, the replication layer identifies at most a single component of the server group as a *primary component*; the other components of a partitioned group are *non-primary components*. A change in the membership of a component is reflected in the delivery of a view-change notification by the group communication layer to each server in that component. The replication layer implements a symmetric distributed algorithm to determine the order of actions to be applied to the database. Each server builds its own knowledge about the order of actions in the system. We use the coloring model defined in [1] to indicate the knowledge level associated with each action. Each server marks the actions delivered to it with one of the following colors:

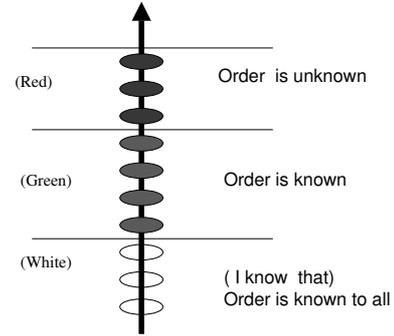


Figure 1. Action coloring

Red Action An action that has been ordered within the local component by the group communication layer, but for which the server cannot, as yet, determine the global order.

Green Action An action for which the server has determined the global order.

White Action An action for which the server knows that all of the servers have already marked it as *green*. These actions can be discarded since no other server will need them subsequently.

At each server, the *white* actions precede the *green* actions which, in turn, precede the *red* ones. An action can be marked differently at different servers; however, no action can be marked *white* by one server while it is missing or is marked *red* at another server.

The actions delivered to the replication layer in a primary component are marked green. Green actions can be applied to the database immediately while maintaining the strictest consistency requirements. In contrast, the actions delivered in a non-primary component are marked red. The global order of these actions cannot be determined yet, so, under the strong consistency requirements, these actions cannot be applied to the database at this stage.

3.1 Conceptual Algorithm

The algorithm presented in this section should, intuitively, provide an adequate solution to the replication problem. This is not actually the case, as the algorithm is not able to deal with some of the more subtle issues that can arise in a partitionable system. We present this simplified solution to provide a better insight into some of the problems the complete solution needs to cope with and to introduce the key properties of the algorithm.

Figure 2 presents the state machine associated with the conceptual algorithm. A replica can be in one of the following four states:

- **Prim State.** The server belongs to the primary component. When a client submits a request, it is multicast

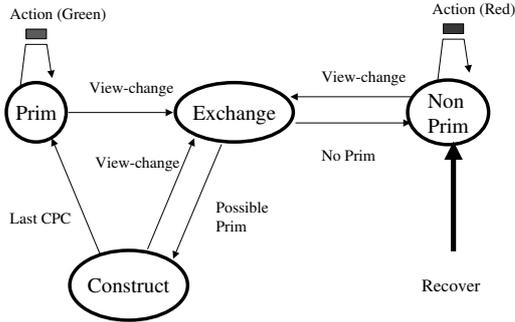


Figure 2. Conceptual Replication Algorithm

using the group communication to all the servers in the component. When a message is delivered by the group communication system to the replication layer, the action is immediately marked green and is applied to the database.

- **NonPrim State.** The server belongs to a non-primary component. Client actions are ordered within the component using the group communication system. When a message containing an action is delivered by the group communication system, it is immediately marked red.
- **Exchange State.** A server switches to this state upon delivery of a view change notification from the group communication system. All the servers in the new view will exchange information allowing them to define the set of actions that are known by some of them but not by all. These actions are subsequently exchanged and each server will apply to the database the green actions that it gained knowledge of. After this exchange is finished each server can check whether the current view has a quorum to form the next primary component. This check can be done locally, without additional exchange of messages, based on the information collected in the initial stage of this state. If the view can form the next primary component the server will move to the Construct state, otherwise it will return to the NonPrim state.
- **Construct State.** In this state, all the servers in the component have the same set of actions (they synchronized in the Exchange state) and can attempt to install the next primary component. For that they will send a Create Primary Component (CPC) message. When a server has received CPC messages from all the members of the current component it will transform all its red messages into green, apply them to the database and then switch to the Prim state. If a view change occurs before receiving all CPC messages, the server returns to the Exchange state.

In a system that is subject to partitioning we must ensure that two different components do not apply contradic-

tory actions to the database. We use a quorum mechanism to allow the selection of a unique primary component from among the disconnected components. Only the servers in the primary component will be permitted to apply actions to the database. While several types of quorums could be used, we opted to use *dynamic linear voting* [15]. Under this system, the component that contains a (weighted) majority of the last primary component becomes the new primary component.

In many systems, processes exchange information only as long as they have a direct and continuous connection. In contrast, our algorithm propagates information by means of *eventual path*: when a new component is formed, the servers exchange knowledge regarding the actions they have, their order and color. This exchange process is only invoked immediately after a view change. Furthermore, all the components exhibit this behavior, whether they will form a primary or non-primary component. This allows the information to be disseminated even in non-primary components, reducing the amount of data exchange that needs to be performed once a server joins the primary component.

4 From Total Order to Database Replication

Unfortunately, due to the asynchronous nature of the system model, we cannot reach complete common knowledge about which messages were received by which servers just before a network partition occurs or a server crashes. In fact, it has been proven that reaching consensus in asynchronous environments with the possibility of even one failure is impossible [10]. Group communication primitives based on Virtual Synchrony do not provide any guarantees of message delivery that span network partitions and server crashes. In our algorithm it is important to be able to tell whether a message that was delivered to one server right before a view change, was also delivered to all its intended recipients.

A server p cannot know, for example, whether the last actions it delivered in the Prim state, before a view-change event occurred, were delivered to all the members of the primary component; Virtual Synchrony guarantees this fact only for the servers that will install the next view together with p . These messages cannot be immediately marked *green* by p , because of the possibility that a subset of the initial membership, big enough to construct the next primary component, did not receive the messages. This subset could install the new primary component and then apply other actions as green to the database, breaking consistency with the rest of the servers. This problem will manifest itself in any algorithm that tries to operate in the presence of network partitions and remerges. A solution based on Total Order cannot be correct in this setting without further enhancement.

4.1 Extended Virtual Synchrony

In order to circumvent the inability to know who received the last messages sent before a network event occurs we use an enhanced group communication paradigm called *Extended Virtual Synchrony* (EVS) [20] that reduces the ambiguity associated with the decision problem. Instead of having to decide on two possible values, as in the consensus problem, EVS will create three possible cases. To achieve this, EVS splits the view-change notification into two notifications: a *transitional* configuration change message and a *regular* configuration change message. The transitional configuration message defines a reduced membership containing members of the next regular configuration coming directly from the same regular configuration. This allows the introduction of another form of message delivery, *safe delivery*, which maintains the total order property but also guarantees that every message delivered to any process that is a member of a configuration is delivered to every process that is a member of that configuration, unless that process fails. Messages that do not meet the requirements for safe delivery, but are received by the group communication layer, are delivered in the transitional configuration. No new messages are sent by the group communication in the transitional configuration.

The safe delivery property provides a valuable tool to deal with the incomplete knowledge in the presence of network failures or server crashes. We distinguish now three possible cases:

1. A safe message is delivered in the regular configuration. All guarantees are met and everyone in the configuration will deliver the message (either in the regular configuration or in the following transitional configuration) unless they crash.
2. A safe message is delivered in the transitional configuration. This message was received by the group communication layer just before a partition occurs. The group communication layer cannot tell whether other components that split from the previous component received and will deliver this message.
3. A safe message was sent just before a partition occurred, but it was not received by the group communication layer in some detached component. The message will, obviously, not be delivered at the detached component.

The power of this differentiation lies in the fact that, with respect to the same message, it is impossible for one server to be in case 1, while another is in case 3. To illustrate the use of this property consider the Construct phase of our algorithm: If a server p receives all CPC messages in the regular configuration, it knows that every server in that configuration will receive all the messages before the next regular

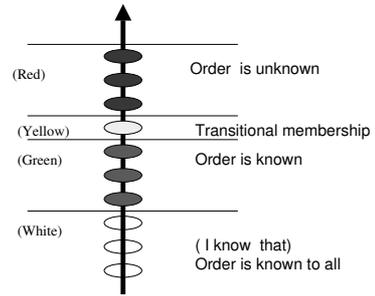


Figure 3. Updated coloring model

configuration is delivered, unless they crash; some servers may, however, receive some of the CPC messages in a transitional configuration. Conversely, if a server q receives a configuration change for a new regular configuration before receiving all of the CPC messages, then no server could have received a message that q did not receive as safe in the previous configuration. In particular, no server received all of the CPC messages as safe in the previous regular configuration. Thus q will know that it is in case 3 and no other server is in case 1. Finally, if a server r received all CPC messages, but some of those were delivered in a transitional configuration, then r cannot know whether there is a server p that received all CPC messages in the regular configuration or whether there is a server q that did not receive some of the CPC messages at all; r does, however, know that there cannot exist both p and q as described.

5 Replication Algorithm

Based on the above observations the algorithm skeleton presented in Section 3.1 needs to be refined. We will take advantage of the Safe delivery properties and of the differentiated view change notification that EVS provides. The two delicate states are, as mentioned, Prim and Construct.²

In the **Prim** state, only actions that are delivered as safe during the regular configuration can be applied to the database. Actions that were delivered in the transitional configuration cannot be marked as green and applied to the database before we know that the next regular configuration will be the one defining the primary component of the system. If an action a is delivered in the transitional membership and is marked directly as green and applied to the database, then it is possible that one of the detached components that did not receive this action will install the next primary component and will continue applying new actions

²While the same problem manifests itself in any state, it is only these two states where knowledge about the message delivery is critical, as it determines either the global total order (in Prim) or the creation of the new primary (Construct).

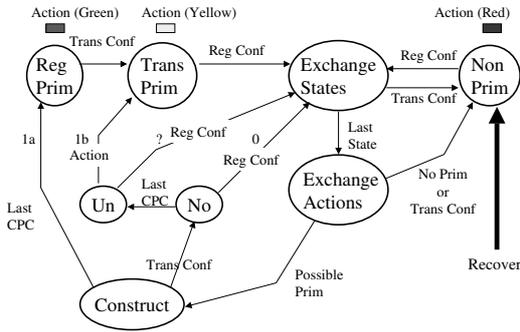


Figure 4. Replication State Machine

to the database, without applying a , thus breaking the consistency of the database. To avoid this situation, the **Prim** state was split into two states: **RegPrim** and **TransPrim** and a new message color was introduced to the coloring model:

Yellow Action An action that was delivered in a transitional configuration of a primary component.

A *yellow* action becomes *green* at a server as soon as this server learns that another server marked the action *green* or when this server becomes part of a primary component. As discussed in the previous section, if an action is marked as *yellow* at some server p , then there cannot exist r and s , in this component, such that one marked the action as *red* and the other marked it *green*.

In the presence of consecutive network changes, the process of installing a new primary component can be interrupted by another configuration change. If a transitional configuration is received by a server p while in the **Construct** state, before receiving all the CPC messages, the server will not be able to install the new primary and will switch to a new state: **No**. In this state p expects to receive the delivery of the new regular configuration which will trigger the initiation of a new exchange round. However, if p receives all the rest of the CPC messages in **No** (i.e. in the transitional configuration), it means that it is possible that some server q has received all the CPC messages in **Construct** and has moved to **RegPrim**, completing the installation of the new primary.

To account for this possibility, p will switch to another new state: **Un** (undecided). If an action message is received in this state then p will know for sure that there was a server q that switched to **RegPrim** and even managed to generate new actions before noticing the network failure that caused the cascaded membership change. Server p , in this situation (1b), has to act as if installing the primary component in order to be consistent, mark its old yellow/red actions as green, mark the received action as yellow and switch to **TransPrim**, “joining” q who will come from **RegPrim** as it will also eventually notice the new configuration change. If the regular configuration message is delivered without any

message being received in the **Un** state (transition marked ? in Figure 4), p remains uncertain whether there was a server that installed the primary component and will not attempt to participate in the formation of a new primary until this dilemma is cleared through exchange of information with one or, in the worst case, all of the members that tried to install the same primary as p .

Figure 4 shows the updated state machine. Aside from the changes already mentioned, the Exchange state was also split into **ExchangeStates** and **ExchangeActions**, mainly for clarity reasons. From a procedural point of view, once a view change is delivered, the members of each view will try to establish a maximal common state that can be reached by combining the information and actions held by each server. After the common state is determined, the participants proceed to exchange the relevant actions. Obviously, if the new membership is a subset of the old one, there is no need for action exchange, as the states are already synchronized.

5.1 Dynamic Replica Instantiation and Removal

As mentioned in the description of the model, the algorithm that we presented so far works under the limitation of a fixed set of potential replicas. It is of great value, however, to allow for the dynamic instantiation of new replicas as well as for their deactivation. Moreover, if the system does not support permanent removal of replicas, it is susceptible to blocking in case of a permanent failure or disconnection of a majority of nodes in the primary component.

However, dynamically changing the set of servers is not straightforward: the set change needs to be synchronized over all the participating servers in order to avoid confusion and incorrect decisions such as two distinct components deciding they are the primary, one being the rightful one in the old configuration, the other being entitled to this in the new configuration. Since this is basically a consensus problem, it cannot be solved in a traditional fashion. We circumvent the problem with the help of the persistent global total order that the algorithm provides.

When a replica r wants to permanently leave the system, it will broadcast a PERSISTENT_LEAVE message that will be ordered as if it was an action message. When this message becomes *green* at replica s , s can update its local data structures to exclude r from the list of potential replicas. The PERSISTENT_LEAVE message can also be administratively inserted into the system to signal the permanent removal, due to failure, of one of the replicas. The message will be issued by a site that is still in the system and will contain the server id of the dead replica.

A new replica r that wants to join the replicated system will first need to connect to one of the members (s of the system, without joining the group). s will act as a representative for the new site to the existing group by creating a

PERSISTENT_JOIN message to announce r 's intention to join the group. This message will be ordered as a regular action, according to the standard algorithm. When the message becomes *green* at a server, that replica will update its data structures to include the newcomer's server id and set the green line (the last globally ordered message that the server has) for the joining member as the action corresponding to the PERSISTENT_JOIN message. Basically, from this point on the servers acknowledge the existence of the new member, although r is still not connected to the group. When the PERSISTENT_JOIN message becomes green at the peer server (the representative), the peer server will take a snapshot of the database and start transferring it to the joining member. If the initial peer fails or a network partition occurs before the transfer is finished, the new server will try to establish a connection with a different member of the system and continue its update. If the new peer already ordered the PERSISTENT_JOIN message sent by the first representative, it will know about r and the state that it has to reach before joining the system, therefore will be able to resume the transfer procedure. If the new peer has not yet ordered the PERSISTENT_JOIN message it will issue another PERSISTENT_JOIN message for the r . PERSISTENT_JOIN messages for members that are already present in the local data structures are ignored by the existing servers, so only the first ordered PERSISTENT_JOIN will define the entry point of the new site into the system. Since the algorithm guarantees global total ordering, this entry point is uniquely defined. Finally, when the transfer is complete, r will set the action counter to the last action that was ordered by the system and will join the group of replicas. This will be seen as a view change by the existing members and they will go through the EXCHANGE states and continue according to the algorithm.

Another method for performing online reconfiguration is described in [19]. This method requires the joining site to be permanently connected to the primary component while being updated. We maintain the flexibility of the engine and we allow joining replicas to be connected to non-primary components during their update stage. It can even be the case that a new site is accepted into the system without **ever** being connected to the primary component, due to the eventual path propagation method. The insertion of a new replica into the system in a non-primary component, can be useful to certain applications as is shown in Section 6.

The static algorithm code was presented in [2], while the complete algorithm code, including the dynamic capabilities can be found in the extended version of this paper [5].

5.2 Proof of Correctness

The algorithm in its static form was proven correct in [2]. The correctness properties that were guaranteed were

liveness, FIFO order and Total global order. Here, we prove that the enhanced dynamic version of the algorithm still preserves the same guarantees.

Lemma 1 (Global Total Order (static)) *If both servers s and r performed their i th actions, then these actions are identical.*

Lemma 2 (Global FIFO Order (static)) *If server r performed an action a generated by server s , then r already performed every action that s generated prior to a .*

These are the two properties that define the **Safety** criterion in [2]. These specifications need to be refined to encompass the removal of servers or the addition of new servers to the system.

Theorem 1 (Global Total Order (dynamic)) *If both servers s and r performed their i th action, then these actions are identical.*

Proof: Consider the system in its start-up configuration set. Any server in this configuration will trivially maintain this property according to Lemma 1. Consider a server s that joins the system. The safety properties of the static algorithm guarantee that after ordering the same set of actions, all servers will have the same consistent database. This is the case when a PERSISTENT_JOIN action is ordered. According to the algorithm s will set its global action counter to the one assigned by the system to the PERSISTENT_JOIN action. From this point on the behavior of s is indistinguishable from a server in the original configuration and the claim is maintained as per Lemma 1. \square

Theorem 2 (Global FIFO Order (dynamic)) *If server r performed an action a generated by server s , then r already performed every action that s generated prior to a , or it inherited a database state which incorporated the effect of these actions.*

Proof: According to Lemma 2, the theorem holds true from the initial starting point until a new member is added to the system. Consider r , a member who joins the system. According to the algorithm, the joining member transfers the state of the database as defined by the action ordered immediately before the PERSISTENT_JOIN message. All actions generated by s and ordered before the PERSISTENT_JOIN will be incorporated in the database that r received. From Theorem 1, the PERSISTENT_JOIN message is ordered at the same place at all servers. All actions generated by s and ordered after the PERSISTENT_JOIN message will be ordered similarly at every server, including r , according to Theorem 1. Since Lemma 2 holds for any other member, this is sufficient to guarantee that r will order all other actions generated by s prior to a , and ordered after r joined the system. \square

Lemma 3 (Liveness (static)) *If server s orders action a and there exists a set of servers containing s and r , and a*

time from which on that set does not face any communication or process failures, then server r eventually orders action a .

This is the liveness property defined in [2] and proven to be satisfied by the static replication algorithm. This specification needs to be refined to include the notion of servers permanently leaving the system.

Theorem 3 (Liveness (dynamic)) *If server s orders action a in a configuration that contains r and there exists a set of servers containing s and r , and a time from which on that set does not face any communication or process failures, then server r eventually orders action a .*

Proof: The theorem is a direct extension of Lemma 3, which acknowledges the potential existence of different server-set configurations. An action that is ordered by a server in one configuration will be ordered by all servers in the same configuration as a direct consequence of Theorem 1. Servers that leave the system or crash do not meet the requirements for the liveness property, while servers that join the system will order the actions generated in any configuration that includes them, unless they crash. \square

6 Supporting Various Application Semantics

The presented algorithm was designed to provide strict consistency semantics by applying actions to the database only when they are marked green and their global order is determined. In the real world, where incomplete knowledge is unavoidable, many applications would rather have an immediate answer, than incur a long latency to obtain a complete and consistent answer. Therefore, we provide additional service types for clients in a non-primary component. The result of a *weak query* is obtained from a consistent, but possibly obsolete state of the database, as reflected by the green actions known to the server at the time of the query, even while in a non-primary component. Other applications prefer getting an immediate reply based on the latest information available, although possibly inconsistent. In the primary component the state of the database reflects the most updated situation and is always consistent. In a non-primary component, however, red actions must be taken into account in order to provide the latest, though not consistent, information. We call this type of query a *dirty query*.

Different semantics can be supported also with respect to updates. In the *timestamp* semantics case, the application is interested only in the most recent information/ Location tracking is a good example of an application that would employ such semantics. Similarly, *commutative* semantics are used in applications where the order of action execution is irrelevant as long as all actions are eventually applied. In an inventory management application all operations on the stock would be commutative. For both semantics, the one-

copy serializability property is not maintained in the presence of network partitions. However, after the network is repaired and the partitioned components merge, the databases states converge.

The algorithm can be significantly optimized if the engine has the ability to distinguish a query-only action from an action that contains updates. A query issued at one server can be answered as soon as all previous actions generated by this server were applied to the database, without the need to generate and order an action message.

Modern database applications exploit the ability to execute a procedure specified by a transaction. These are called *active* transactions and they are supported by our algorithm, provided that the invoked procedure is deterministic and depends solely on the current database state. The procedure will be invoked at the time the action is ordered, rather than before the creation of the update.

Finally, we mentioned that our model best fits one-operation transactions. Actually, any non-interactive transactions that do not invoke triggers are supported in a similar way. However, some applications need to use interactive transactions which, within the same transaction, read data and then perform updates based on a user decision, rather than a deterministic procedure. Such behavior, cannot be modeled using one action, but can be mimicked with the aid of two actions. The first action reads the necessary data, while the second one is an active action as described above. This active action encapsulates the update dictated by the user, but first checks whether the values of the data read by the first action are still valid. If not, the update is not applied, as if the transaction was aborted in the traditional sense. Note that if one server “aborts”, all of the servers will abort that (trans)action, since they apply an identical deterministic rule to an identical state of the database.

7 Performance Analysis

In this section we evaluate our replication engine and compare its performance to that of two existing solutions: two-phase commit (2PC) and COREL by Keidar [16]. 2PC is the algorithm adopted by most replicated systems that require strict consistency. 2PC requires two forced disk writes and $2n$ unicast messages per action. COREL exploits group communication properties to improve on that. COREL requires one forced disk write and n multicast messages per action. In contrast, our engine only requires $1/n$ forced disk write and one multicast message per action on average (only the initiating server needs to force the action to disk).

We implemented all three algorithms and compared their performance in normal operation, without view changes. Our 2PC implementation does not perform the locking required to guarantee the unique order of transaction execution, as this is usually the task of the database. Therefore

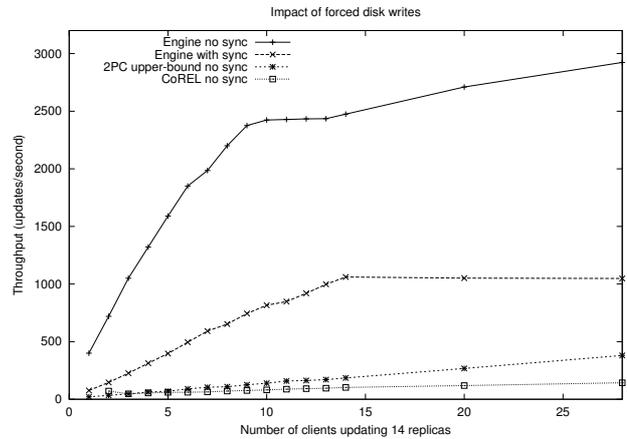
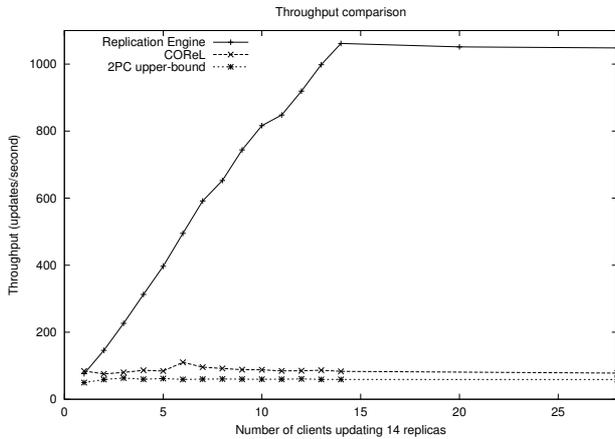


Figure 5. Throughput Comparison

a complete 2PC will perform strictly worse than our upper-bound implementation.

Since we are interested in the intrinsic performance of the replication methods, clients receive responses to their actions as soon as the actions are globally ordered, without any interaction with a database. A follow-up work [3] evaluates a complete solution that replicates a Postgres database over local and wide area networks using our engine.

All the tests were conducted with 14 replicas, each running on a dual processor Pentium III-667 with Linux connected by a 100Mbps/second local area switch. Each action is 200 bytes long (e.g. an SQL statement).

Figure 5(a) compares the maximal throughput that a system of 14 replicas can sustain under each of the three methods. We vary the number of clients that simultaneously submit requests into the system between 1 and 28, evenly spread between the replicas as much as possible. The clients are constantly injecting actions into the system, the next action from a client being introduced immediately after the previous action from that client is completed and its result reported to the client.

Our engine achieves a maximum throughput of 1050 updates/second once there are sufficient clients to saturate the system, outperforming the other methods by at least a factor of 10. COReL outperforms the upper-bound 2PC as expected, mainly due to the saving in disk writes reaching a maximum of 110 updates/second as opposed to 63 updates/second for the upper-bound 2PC.

High-performance database environments commonly use superior storage technology (e.g. flash disks). In order to estimate the performance that the three methods would exhibit in such environment, we used asynchronous disk writes instead of forced disk writes. Figure 5(b) shows that our engine tops at processing 3000 updates/second. Under the same conditions, the upper-bound 2PC algorithm achieves 400 updates/second. COReL reaches a through-

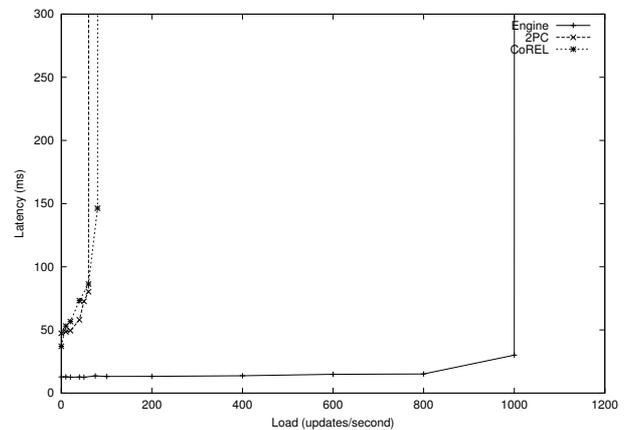


Figure 6. Latency variation with load

put of approximately 100 updates/second with 28 clients, but more clients are needed in order to saturate the COReL system due to its higher latency. With 50 clients, COReL saturates the system with about 200 updates/second. 2PC outperforms COReL in this experiment because of two reasons: the fact that we use an upper-bound 2PC as mentioned above, and the particular switch that serves our local area network that is capable of transmitting multiple unicast messages between different pairs in parallel.

We also measured the response time a client experiences under different loads (Figure 6). Our Engine maintains an average latency of 15ms with load increasing up to 800 updates/second and breaks at the maximum supported load of 1050 updates/second. COReL and 2PC experience latencies of 35ms up to 80ms under loads up to 100 updates/second with COReL being able to sustain more throughput.

8 Conclusions

We presented a complete algorithm for database replication over partitionable networks sophisticatedly utilizing

group communication and proved its correctness. Our avoidance of the need for end-to-end acknowledgment per action contributed to superior performance. We showed how to incorporate online instantiation of new replicas and permanent removal of existing replicas. We also demonstrated how to efficiently support various types of applications that require different semantics.

Acknowledgements

We thank Jonathan Stanton for his numerous technical ideas and support that helped us optimize the overall performance of the system. We also thank Michal Miskin-Amir and Jonathan Stanton for their insightful suggestions that considerably improved the presentation of this paper. This work was partially funded by grant F30602-00-2-0550 from the Defense Advanced Research Projects Agency (DARPA). The views expressed in this paper are not necessarily endorsed by DARPA.

References

- [1] O. Amir, Y. Amir, and D. Dolev. A highly available application in the Transis environment. *Lecture Notes in Computer Science*, 774:125–139, 1993.
- [2] Y. Amir. *Replication Using Group Communication over a Partitioned Network*. PhD thesis, Hebrew University of Jerusalem, Jerusalem, Israel, 1995. <http://www.cnds.jhu.edu/publications/yair-phd.ps>.
- [3] Y. Amir, C. Danilov, M. Miskin-Amir, J. Stanton, and C. Tutu. Practical wide-area database replication. Technical Report CNDS 2002-1, Johns Hopkins University, Center for Networking and Distributed Systems, 2002.
- [4] Y. Amir and J. Stanton. The spread wide area group communication system. Technical Report CNDS 98-4, Johns Hopkins University, Center for Networking and Distributed Systems, 1998.
- [5] Y. Amir and C. Tutu. From total order to database replication. Technical Report CNDS 2002-3, Johns Hopkins University, Center for Networking and Distributed Systems, 2002.
- [6] P. Bernstein, D. Shipman, and J. Rothnie. Concurrency control in a system for distributed databases (sdd-1). *ACM Transactions on Database Systems*, 5(1):18–51, Mar. 1980.
- [7] K. P. Birman and T. A. Joseph. Exploiting virtual synchrony in distributed systems. In *Proceedings of the ACM Symposium on OS Principles*, pages 123–138, Austin, TX, 1987.
- [8] A. Demers et al. Epidemic algorithms for replicated database maintenance. In *Proceedings of the 6th Annual ACM Symposium on Principles of Distributed Computing*, pages 1–12, Vancouver, BC, Canada, Aug. 1987.
- [9] A. Fekete, N. Lynch, and A. Shvartsman. Specifying and using a partitionable group communication service. *ACM Transactions on Computer Systems*, 19(2):171–216, May 2001.
- [10] M. H. Fischer, N. A. Lynch, and M. S. Paterson. Impossibility of consensus with one faulty process. *Journal of the ACM*, 32(2):374–382, Apr. 1985.
- [11] R. Golding. *Weak-Consistency Group Communication and Membership*. PhD thesis, UC Santa Cruz, 1992.
- [12] J. N. Gray and A. Reuter. *Transaction Processing: concepts and techniques*. Data Management Systems. Morgan Kaufmann Publishers, Inc., San Mateo (CA), USA, 1993.
- [13] V. Hadzilacos and S. Toueg. Fault-tolerant broadcasts and related problems. In S. Mullender, editor, *Distributed Systems*, chapter 5. Addison-Wesley, second edition, 1993.
- [14] J. Holliday, D. Agrawal, and A. E. Abbadi. Database replication using epidemic update. Technical Report TRCS00-01, University of California Santa-Barbara, 19, 2000.
- [15] S. Jajodia and D. Mutchler. Dynamic voting algorithms for maintaining the consistency of a replicated database. *ACM Transactions on Database Systems*, 15(2):230–280, 1990.
- [16] I. Keidar. A highly available paradigm for consistent object replication. Master’s thesis, Institute of Computer Science, The Hebrew University of Jerusalem, Israel, 1994.
- [17] I. Keidar and D. Dolev. Increasing the resilience of atomic commit at no additional cost. In *Symposium on Principles of Database Systems*, pages 245–254, 1995.
- [18] B. Kemme and G. Alonso. A new approach to developing and implementing eager database replication protocols. *ACM Transactions on Database Systems*, 25(3):333 – 379, 2000.
- [19] B. Kemme, A. Bartoli, and O. Babaoglu. Online reconfiguration in replicated databases based on group communication. In *Proceedings of the International Conference on Dependable Systems and Networks*, Göteborg, Sweden, 2001.
- [20] L. E. Moser, Y. Amir, P. M. Melliar-Smith, and D. A. Agarwal. Extended virtual synchrony. In *International Conference on Distributed Computing Systems*, pages 56–65, 1994.
- [21] M. Patino-Martinez, R. Jimenez-Peris, B. Kemme, and G. Alonso. Scalable replication in database clusters. In *Proceedings of 14th International Symposium on Distributed Computing (DISC’2000)*, 2000.
- [22] F. Pedone. *The Database State Machine and Group Communication Issues*. PhD thesis, École Polytechnique Fédérale de Lausanne, Switzerland, 1999.
- [23] F. Pedone, R. Guerraoui, and A. Schiper. Exploiting atomic broadcast in replicated databases. In *Proceedings of EuroPar (EuroPar’98)*, Sept. 1998.
- [24] C. Pu and A. Leff. Replica control in distributed systems: an asynchronous approach. *ACM SIGMOD Record*, 20(2):377–386, 1991.
- [25] F. B. Schneider. Implementing fault-tolerant services using the state machine approach: A tutorial. *ACM Computing Surveys*, 22(4):299–319, Dec. 1990.
- [26] D. Skeen. A quorum-based commit protocol. Berkley Workshop on Distributed Data Management and Computer Networks, February 1982.
- [27] I. Stanoi, D. Agrawal, and A. E. Abbadi. Using broadcast primitives in replicated databases. In *Proceedings of the 18th IEEE International Conference on Distributed Computing Systems ’98*, pages 148–155, Amsterdam, May 1998.
- [28] M. Stonebraker. Concurrency control and consistency of multiple copies of data in distributed INGRES. *IEEE Transactions on Software Engineering*, SE-5:188–194, May 1979.