

A Pragmatic Protocol for Database Replication in Interconnected Clusters

J. Grov
U. Oslo

L. Soares
U. Minho

A. Correia Jr.
U. Minho

J. Pereira
U. Minho

R. Oliveira
U. Minho

F. Pedone
U. Lugano

Abstract

Multi-master update everywhere database replication, as achieved by protocols based on group communication such as DBSM and Postgres-R, addresses both performance and availability. By scaling it to wide area networks, one could save costly bandwidth and avoid large round-trips to a distant master server. Also, by ensuring that updates are safely stored at a remote site within transaction boundaries, disaster recovery is guaranteed. Unfortunately, scaling existing cluster based replication protocols is troublesome.

In this paper we present a database replication protocol based on group communication that targets interconnected clusters. In contrast with previous proposals, it uses a separate multicast group for each cluster and thus does not impose any additional requirements on group communication, easing implementation and deployment in a real setting. Nonetheless, the protocol ensures one-copy equivalence while allowing all sites to execute update transactions. Experimental evaluation using the workload of the industry standard TPC-C benchmark confirms the advantages of the approach.

1. Introduction

Database replication is an attractive concept both to increase fault tolerance and to improve scalability by enabling several database sites to serve the same queries. The main challenge of such systems is that coordinating updates among the participating servers inevitably delays the execution of update-transactions. A particularly promising approach is taken by replication protocols based on group communication such as DBSM [21, 12] and Postgres-R [17, 31]. By taking advantage of optimistic concurrency control allowed by transactional semantics and of atomic multicast provided by group communication, it provides performance and scalability even in face of demanding workloads such as the industry standard TPC-C benchmark [27].

Unfortunately, scaling existing cluster based replication protocols to a wide area network is troublesome. Notably, the latency of uniform atomic (or safe) delivery required to ensure fault tolerance has a profound impact in optimistic

concurrency protocols leading to increased abort rate [11]. This wastes resources and endangers the ability to commit long lived transactions in a busy server. Although optimistic delivery can mitigate this limitation [26], using it requires an in-depth rewrite of existing protocol implementations. In fact, the only generally available group communication toolkit supporting it is Appia [20, 23].

Furthermore, although research has been addressing group communication in wide area networks for a long time, industrial deployment is far more common in clusters. Therefore one should expect wide area features to be far less tested and optimized, if implemented at all. The overhead of maintaining automatic management of membership spanning multiple geographically apart sites is also not negligible. Finally, the practicality of group communication over wide area networks is also compromised by network configuration and security issues, such as firewalls, tunnels and NAT gateways. In particular, using true multicast for efficiency is often not an option.

In this paper we present WICE, a protocol targeted at multiple clusters interconnected by a wide area network. In contrast with lazy replication protocols, such as Oracle Streams [30], WICE ensures that no globally committed transaction (i.e., a transaction that has been acknowledged to clients) is lost. On the other hand, by allowing all replicas to be fully on-line and executing update transactions, it improves resource efficiency and performance when compared to volume replication [28], often the only choice for disaster recovery in mission critical applications. In detail, the contributions of this paper are:

- Introduces the protocol providing 1-copy equivalence of the native database consistency criterion, even in the presence of faults, while confining group communication within LANs and improving practicality.
- Takes advantage of directly implementing stabilization across wide-area directly on TCP/IP to greatly reduce the likelihood of a transaction being aborted during the certification phase, which is the single greatest obstacle to the scalability of previous proposals [11].
- Provides an experimental evaluation of the protocol applied to a multi-version database when running the workload of the industry standard TPC-C bench-

mark [29], thus verifying the previous claim.

2. Background

In outline, the responsibilities of a replica control protocol are (1) the overall concurrency control, ensuring that transactions executing concurrently on different sites do not violate consistency, (2) the efficient propagation of updates among the replicas, and (3) atomic commitment, that is, ensuring that agreement is reached among all replicas on the outcome of transactions.

In a wide-area setting, a major challenge is to minimize the number of messages in order to reduce overall latency. Previous studies [7, 3, 17, 21, 11, 31] advocate the appropriateness of *optimistic* [18] concurrency control to such settings.

In these protocols, atomic commitment is handled through *dictatorial* commit [1]: No individual site can veto a commit decision as long as a remote transaction request is validated by the distributed validation procedure. It is then common [17, 21], to combine update propagation and dictatorial-commit using an atomic broadcast primitive [15] which ensures that all non faulty replicas receive updates in the same total order. This ordering is then used as a global timestamp [6], and is used to decide whether a transaction is allowed to commit or abort.

One challenge arising in all protocols that combine dictatorial-commit protocols with failover guarantees is that before a transaction can commit, it must be ensured that the updates are present at *all* non faulty sites before the transaction is allowed to commit. This is necessary to avoid *dirty reads*, as illustrated by the following example: Assume a transaction T is allowed to commit at site A before T 's updates are installed at site B , and assume A fails immediately after, such that T is never known at site B . This can clearly lead to inconsistencies, as any client who has already read T 's updates from A will assume that the updates of T are also present in cluster B . To handle this, we must ensure that the updates are *stable*, i.e., that they will be known by at least one remaining site after any kind of disaster from which we want to recover, before we allow the transaction to commit.

If in cluster settings the modular use of a primitive providing uniform atomic delivery [15] is advantageous by offloading all the group communication complexity and optimization to the communication toolkit, the high and uneven latencies in wide area networks require a differentiated treatment of local and remote replicas. This allows us to optimize the replication protocol by masking the high latency and make a conscious use of the bandwidth of wide area links.

While for wide area specific database replication we opt for the optimistic execution approach, two recent studies of group communication based replication recommend the

use of pessimistic execution [2, 19] and follow the active replication approach [24, 14]. In [2], each update request is atomically broadcast and executed by every replica. In [19], transactions are broadcast at once greatly reducing the previous (per request) message overhead. To avoid executing transactions sequentially, in [19] it is assumed that transactions are a priori annotated with conflict classes so that coarse-grained locks can be acquired before starting execution and thus non-conflicting transactions are allowed to execute concurrently.

3. System Model

We assume the page model for a database [6]: A collection of named data items which have a value. The combined values of the data items at any given moment is the database state. We do not make any assumptions on the granularity of data items.

A database site is modeled as a sequential process. In detail, the execution of each site is modeled as a sequence of steps that may change the site's state. Namely, the database state is manipulated by executing $READ(x)$ and $WRITE(x)$ steps, where x represents a database tuple. A transaction is a sequence of read and write operations followed by a $COMMIT(t)$ or $ABORT(t)$ operation. Each site contains a complete copy of the database and is responsible for ensuring local concurrency control.

We consider a finite set of database sites that communicate through a fully connected network. Both computation and communication are asynchronous. Sites may fail only by crashing and do not recover, thus stopping to execute database operations, or send or deliver further messages.

Database sites are organized in clusters. Within a cluster we assume a primary component group membership service that provides current and consistent views of the sites believed to be up [9]. This service is intended to allow, at any moment, the deterministic identification of a distinguished site as the cluster's *delegate* as well as providing a view-synchronous multicast primitive (Section 3.2). The availability of a primary component group membership service implicitly assumes that consensus is solvable in our system model [13]. The assumed failure patterns and failure detection capabilities of our model are thus indirectly determined by the actual solution adopted for consensus.

Among clusters, we assume that the failure of an entire cluster is reliably detected at the other sites. That is, if all sites in a cluster fail then the fact is eventually noticed by the other clusters' delegates. Otherwise, the cluster is never suspected to have failed.¹ At each cluster, the set of clusters believed to be up is given by a function `remoteClusters()`.

¹This assumption is equivalent to have a perfect failure detector among the clusters [8]. In a wide area setting, its provision would require the use of a specially dedicated communication infrastructure among the clusters or rely on human intervention to declare the unavailability of all cluster sites.

3.1. Database Interface

The replication protocol presented in Section 4 uses a replication interface with the database engine that is part of the API being defined in the context of the GORDA project [10]. The interface has been implemented in a number of DBMS, notably in PostgreSQL [16] and Derby [4]. The interested reader can find its detailed definition in [22]. Basically, it allows the inspection of a transaction's execution at three specific points: (1) at the beginning of the transaction's execution, (2) at the end of the transaction's execution, just before it starts committing updates or rolls back, and (3) when the local database system has committed the transaction and is ready to reply to the client. Furthermore, the database engine provides an update function executed with priority over any other running transactions that allows to update the values of a given set of items.

More precisely, we assume that the replicated database engine allows to register four callback functions as follows:

onExecuting(tid) invoked before a transaction is about to enter the executing state, i.e., before it starts execution.

The transaction is identified by *tid*.

onCommitting(tid, rs, ws, wv) invoked when the transaction *tid* succeeds and is about to enter the commit phase. The database provides the set of tuples read (*rs*) and written (*ws*) by the transaction, as well as the written values (*wv*). At this point the transaction has all its updates buffered and all write locks still acquired.

onAborting(tid) invoked when the transaction *tid* fails and is about to abort.

onCommitted(tid) invoked after the transaction has completed making all updates persistent, released locks, entered the committed state and is ready to reply to the client.

When it invokes any of the above functions, the database engine suspends the execution of the transaction until the protocol replies by invoking the database functions `continueExecuting(tid)`, `continueCommitting(tid)`, `continueAborting(tid)` and `continueCommitted(tid)`, respectively.

Replica updates are submitted to the database using the `db_update(tid, ws, wv)` function which applies the values in *wv* to the tuples in *ws* by means of a high priority transaction. A transaction submitted through `db_update` only triggers the `onCommitted(tid)` event. High priority means that any regular (i.e., non high priority) transaction holding locks on any item in *ws* will be aborted. Moreover, high priority transactions are serialized when requesting locks and then executed concurrently.

3.2. Communication Primitives

Among sites within the same cluster, a group communication toolkit is available providing reliable point-to-point communication and FIFO uniform view-synchronous multicast [9]. These primitives rely on the existence of a (primary component) group membership service that tracks the

membership of the cluster through a sequence of *views* satisfying the following property:

Agreement on the view history Let v_i^p denote the i^{th} view at site p . If p installs v_i^p and if q installs v_i^q , then we have $v_i^p = v_i^q$.

The agreement property allows us to denote a view simply by v_i without mentioning a particular site. The specification of a group membership service includes additional properties that, for the sake of simplicity, we intentionally omit here.

Point-to-point reliable communication is defined by two primitives `r_send` and `r_deliver`, and satisfies the following properties:

Delivery Integrity A message m is delivered at most once and only if m has been previously sent.

Reliable Delivery: If sites p and q both belong to $v_i \cap v_{i+1}$ and p sends m to q in v_i then q delivers m before v_{i+1} .

Uniform view-synchronous multicast is defined through primitives `u_vscast` and `u_vsdeliver` and satisfies the following two properties:

Delivery Integrity A message m is delivered at most once and only if m has been previously vscast.

Sending View Delivery A message m vscast in view v_i is delivered in view v_i .

Uniform View Synchrony If site p belongs to view v_i and delivers m before v_{i+1} , then every site q in $v_i \cap v_{i+1}$ delivers m before v_{i+1} .

FIFO uniform view-synchronous multicast is invoked through primitive `fifo_u_vscast` and in addition to the above properties satisfies:

FIFO Delivery If site p vscasts two consecutive messages m_1 and m_2 , then no site delivers m_2 before m_1 .

Among clusters, messages are exchanged through a point-to-point FIFO reliable channel. They use primitives `fifo_r_send` and `fifo_r_deliver` satisfying the following properties:

Delivery Integrity A message m is delivered at most once and only if m has been previously sent.

Reliable Delivery: If cluster p and q are both correct and p sends m to q then q eventually delivers m .

FIFO Delivery If cluster p sends two consecutive messages m_1 and m_2 to cluster q , then q does not deliver m_2 before m_1 .

A cluster is said to be correct if it does not fail entirely.

4. The WICE Protocol

The WICE protocol adopts an optimistic concurrency control policy. Transactions are executed optimistically at any site and then, just before commit, certified against concurrent transactions. WICE borrows from protocols such as Postgres-R [17] and DBSM [21] often called *certification based* protocols. These protocols share two fundamental characteristics: (1) each database site is assumed to store

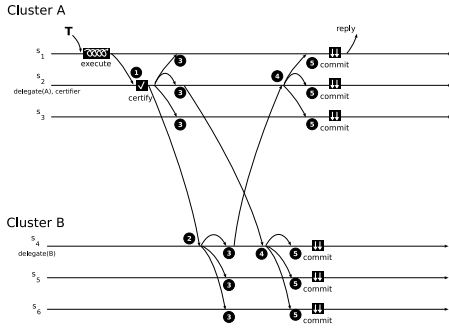


Figure 1: WICE: example of handling of transaction T

the whole database and transactions can be executed at any site, and (2) all update transactions are certified and, if valid, committed in the same order at all sites.

This total order allows the use of a simple optimistic certification procedure to provide one-copy serializability (1-SR)[18]: The execution of a transaction T is valid if and only if, for each item x read by T , the value read by T was written by the most recent updater of x among all transactions ordered before T . As discussed in [12] and [31], we can gain scalability in multi-version database systems if we instead apply one-copy snapshot isolation (1-SI). The certification procedure is the same, but read-operations are never considered as each transaction T is assumed to read from a *snapshot* defined by all committed transaction when T entered execution[5].

The fundamental difference of the three protocols is when and where the certification is performed. Both Postgres-R and DBSM use a total order broadcast primitive and certify each transaction once the totally ordered message is delivered. In Postgres-R, each transaction is certified at the site that executed it and the outcome of the certification is then sent to all the other sites. In the DBSM, the read and write sets of the transaction are sent to all sites allowing the certification to be carried at all sites avoiding the last communication step of Postgres-R.

WICE does not make use of a total order communication primitive, instead ordering is explicitly handled by the protocol. In WICE, one of the sites plays the role of certifier, it totally orders and certifies all transactions. Each valid transaction is then broadcast together with its commit order and updates. This allows to leverage the knowledge about the system's topology and to make optimizations that would not be possible otherwise.

The WICE algorithm is exemplified in Figure 1. In a nutshell, the handling of a transaction proceeds as follows. Consider a system consisting of two clusters A and B. Each cluster has a designated delegate responsible for handling the communication with the other cluster. The delegate of cluster A, site s_2 is also responsible for certifying all exe-

cuted transactions. When an update transaction T is submitted to site s_1 (T 's initiator), it is readily executed and sent to the certifier. If it succeeds, then the certifier propagates T 's updates and commit order, both locally and to cluster's B delegate. The latter, in turn, propagates T locally. Once a delegate is certain that all sites in its cluster delivered T 's data it acknowledges the fact to the other cluster's delegate. This acknowledgement is multicast locally by each delegate. Once a database site knows T 's data has been delivered everywhere and all previous transactions had been committed or aborted, then it commits T . The initiator of T can then reply to its client.

Note that the algorithm discussed here only applies to update transactions, as read-only transactions do not need such a validation. Nevertheless we cannot allow any transaction to read and expose updates before the updating transactions become stable, i.e., committed. For clarity, we omit this from the protocol and assume it to be handled by the local DBMS by blocking the commit of a read-only transaction until all updaters from which it has read from become stable.

4.1. Algorithm

We now consider the protocol algorithm in detail (Figure 2). It is composed by a set of handlers that deal with events triggered by the database engine ("Events at the initiator" and "Transaction commit") and with message delivery. We assume that every database site knows the current system's certifier through a function `certifier()`. The local concurrency control strategy of a given site, which we admit to be either snapshot isolation (SI) or strict two-phase locking (S2PL), is given by the function `localCC()`. Each cluster delegate can find the other participating clusters through a function `remoteClusters()` as well as identifying some delegate's cluster through function `cluster()`. Further, the function `delegate()` is used to determine whether the current site is the delegate of its cluster or not.

Global site variables Each database site manages four sets containing transactions known to be certified, locally updated, locally committed and remotely stable. It keeps track of the number of locally executed transactions in variable `lts`. The certifier keeps track of the number of certified transactions in variable `gts`.

Events at the initiator Before a transaction `tid` executes its first operation, the `onExecuting` handler is invoked. The version of the database seen by `tid` is required for the validation procedure, and for sites running snapshot isolation, this is equal to the number of committed transactions when `tid` begins execution. For sites using two-phase locking, the version must instead be recorded at the end of the execution, i.e., in the `onCommitting` handler.

If the transaction at any time aborts locally, `onAborting()` is invoked and the transaction is simply forgotten by the

protocol. On the contrary, if *tid* succeeds execution then `onCommitting()` is invoked. If local consistency is S2PL, the database version is recorded here. Then, *tid*'s read set, write set and written values (*rs*, *ws* and *wv*) provided by the database are reliably sent to the certifier along with the version of the database on which the transaction executed. The transaction's execution is left suspended until it is certified and its outcome known. If *tid* ends up committing then `continueCommitting(tid)` will be called, otherwise the initiator receives a (ABORT, *tid*) message from the certifier and forces the transaction to abort locally.

Certification Upon delivering an update transaction to certify — (CERTIFY, *tid*, *ts*, *rs*, *ws*, *wv*) — from some initiator site the certifier performs the certification of *tid* against its concurrent transactions. For every certified transaction (but not necessarily committed yet) *ctid* with timestamp equal or greater than *tid*'s, a certification function is called with *ctid*'s write set and *tid*'s read and write sets. When preserving 1-SR the certification function checks *tid*'s read and write sets against *ctid*'s write set. If 1-SI is the adopted consistency criterion then only the write sets of both transactions are compared. In both cases, if there is a non empty intersection then the certification fails and an abort message is sent back to *tid*'s initiator.

When *tid*'s passes the certification test then the certifier's sequence number is incremented and *tid* added to its set of certified transactions. The transaction's id, commit order, write set and written values are then sent to all other replicas. Locally, *tid* is sent using the FIFO uniform view-synchronous multicast primitive as a (UPDATE_LOC, *tid*, *gts*, *ws*, *wv*) message. Remotely, it is sent using the FIFO reliable point-to-point primitive to each remote cluster as a (UPDATE_REM, *tid*, *gts*, *ws*, *wv*) message.

Remote delivery of updates Once a cluster delegate delivers a transaction from the certifier it simply forwards the message to the local replicas using the FIFO uniform view-synchronous multicast primitive.

Local delivery of updates When a replica delivers a transaction *tid* it signals the fact adding it to its set of updated transactions. The use of a uniform primitive ensures that once the transaction is delivered at the current replica it is eventually delivered at all non faulty replicas in the cluster. Therefore, if the replica is a cluster delegate it acknowledges the fact that *tid* became stable at the cluster to all clusters. The just delivered updates are applied. If the replica is the *tid*'s initiator then it just needs to proceed with `continueCommitting(tid)`. Although *tid* does not hold high priority locks at the initiator, the fact that it passed certification means that between its execution and the given commit order, no other certified transaction conflicted with it, and consequently, *tid* will not be aborted by another transaction requesting high-priority locks at *tid*'s initiator. For all other sites, `db.update` is invoked.

Delivery of remote acks Each time a delegate delivers a stability acknowledgment for transaction *tid* from some cluster, the pair (*tid*, *cluster*) is added to its acks set. When *tid* has been acknowledged by all remote clusters, then the delegate locally declares the transaction remotely stable using the (non- uniform) view-synchronous multicast primitive — (STABLE_REM, *tid*). When this message is delivered each replica adds *tid* to its remotestable set.

Transaction commit Here, each site handles the `onCommitted` callback. When `onCommitted` (*tid*) is invoked the site just increments its local database version *lts* and adds *tid* to its committed set. Since all *tid* locks have been released then any new transaction can read from *tid* and therefore from a more recent version of the database. When *tid* is known to be committed locally and stable everywhere the database is then allowed to reply to the client, which happens after `continueCommitted(tid)`.

4.2. Failure Handling

The WICE algorithm tolerates both the failure of single database sites as well as the failure of whole clusters. In this section we present and explain the recovery procedures in both cases.

Locally, each cluster is governed by a group membership service and local communication rests on view-synchronous multicast primitives. This definitely eases failure handling locally. In the event of a site been expelled from the group (because it was taken down, has failed, became unreachable, etc.) every other site in the group eventually becomes aware of the fact by installing a new view of the group. This allows each site to deterministically determine the cluster's delegate should the former failed. Moreover, view-synchrony ensures that all sites surviving the previous view delivered the same set of messages, thus not requiring any synchronization among them. As a result, no particular procedure is required on the failure or an ordinary site. In the next two sections we examine the failures of a cluster's delegate and of the system's certifier. Then, we consider the failure of an entire cluster. For the sake of simplicity and lack of space, we assume that no sites are added to a cluster and that once a site is expelled from the group, whatever was the reason for this, it is no longer readmitted.

4.2.1 Delegate Failover

In Figure 3a, we sketch a protocol for recovering from a site failure when this site was the cluster's delegate. On a view change, site *d* becomes aware it is the new cluster's delegate. To ensure that no transactions are blocked, *d* must re-run all transaction updates and acknowledgements received from remote clusters that may have been incompletely processed by the previous delegate.

Global site variables

```

1 local = ts = {}
2 certified = updated = {}
3 committed = remotestable = acks = {}
4 gts = lts = 0

```

Events at the initiator

```

5 upon onExecuting(tid)
6   if localCC() == S1 then
7     local[tid]=lts
8     continueExecuting(tid)
9   end

```

```

10 upon onCommitting(tid, rs, ws, wv, type)
11   if localCC() == S2PL then
12     local[tid]=lts
13     rsend(CERTIFY, tid, local[tid], rs, ws, wv) to certifier()
14   end

```

```

15 upon onAborting(tid)
16   continueAborting(tid)
17 end

```

```

18 upon rdeliver(ABORT, tid) from i
19   db.abort(tid)
20 end

```

(1) Certification

```

21 upon rdeliver(CERTIFY, tid, ts, rs, ws, wv) from initiator
22   foreach (ctid, clts, cws, cwv) in certified do
23     if clts ≥ ts and !certification(cws, rs, ws) then
24       r_send(ABORT, tid) to initiator
25       return
26     gts = gts + 1
27     enqueue (tid, gts, ws, wv) to certified
28     fifo.u_vscast(UPDATE.LOC, tid, gts, ws, wv)
29     foreach cluster in remoteClusters() do
30       fifo.r_send(UPDATE.REM, tid, gts, ws, wv) to
31         cluster
32   end

```

(2) Remote delivery of updates

```

32 upon fifo.r_deliver(UPDATE.REM, tid, ts, ws, wv) from certifier
33   fifo.u_vscast(UPDATE.LOC, tid, ts, ws, wv)
34 end

```

(3) Local delivery of updates

```

35 upon fifo.u_vsdelay(UPDATE.LOC, tid, ts, ws, wv)
36   ts[tid] = ts
37   enqueue (tid, ts, ws, wv) to updated
38   if delegate() then
39     foreach cluster in remoteClusters() do
40       r_send(ACK.REM, tid) to cluster
41   if local[tid] then
42     continueCommitting(tid)
43   else
44     db.update(tid, ws, wv)
45   end

```

(4 and 5) Delivery of remote acks

```

46 upon r_deliver(ACK.REM, tid) from cluster
47   acked = {}
48   add (tid, cluster) to acks
49   foreach (tid, c) in acks do
50     add c to acked
51   if remoteClusters() ⊆ acked then
52     u_vscast(STABLE.REM, tid)
53   end

```

```

54 upon vsdelay(STABLE.REM, tid)
55   add (tid) to remotestable
56 end

```

Transaction commit

```

57 upon onCommitted(tid) and ts[tid] = lts + 1
58   lts = lts + 1
59   add tid to committed
60 end

```

```

61 upon (tid) in committed and (tid) in remotestable
62   continueCommitted(tid)
63 end

```

Figure 2: WICE protocol

New delegate: Synchronization request When initialized, the new delegate d sends a message (DELEGATE_SYNC, lts) to the certifier in order to ensure that all transactions certified since lts are delivered in its local cluster. The lts value corresponds to the latest transactions updated in d 's cluster. The new delegate also contacts each remote cluster with (ACK_SYNC, lts, TRUE) acknowledging the local stability of all transactions up to lts, requesting similar action from the recipients (argument TRUE of the message).

Certifier: Handle synchronization request When delivering this message, the certifier resends (in order) each certified transaction with a certification timestamp larger than d 's lts value.

All delegates: Synchronize ACK's When the message (ACK_SYNC, clts, reply) from a cluster is delivered in a remote cluster C , the delegate of C regards all its updated transactions with $ts \leq clts$ as acknowledged by cluster. It then just checks whether these transactions became stable in every cluster and proceeds accordingly. If reply was set to TRUE a similar message (now with reply set to FALSE) is sent back to the initializing delegate (just elected) so it can also update the respective acknowledgements.

4.2.2 Certifier Failover

The most serious single server failure is when the current system's certifier becomes unavailable. When initialized, the new certifier advertises itself to all delegates. There may be previously certified transactions not yet known to new certifier so a state synchronization is due. Figure 3b shows our synchronization protocol in pseudocode. The code assumes two existing functions, blockCertification() and unblockCertification(). Their implementation is not shown,

but they state whether all arriving certification requests should be buffered, awaiting the synchronization protocol to finish.

New certifier: Synchronization request The new certifier c starts by invoking blockCertification() and requesting from each cluster all the transactions they might have delivered and updated after the last one updated by c .

Each delegate: Send missing transactions When a (CERTSYNC_REQUEST, clts) is received by the delegate of a cluster C , it replies with a list of its updated transactions (tid, ts, ws, wv) such that $ts > clts$, that is, transactions not yet seen by the new certifier.

Certifier: Missing updates When processing a (CERTSYNC_REPLY, clts, missing) from remote cluster C , the new certifier c then checks each member of the missing list whether it has already received this transaction from another cluster. This will happen if two or more remote clusters both know about a transaction which is unknown to c . If not, the transaction is enqueued in c 's certified queue. As soon as all replies from remoteCluster() are delivered, c sets the certifiers counter gts to lts and starts distributing from its certified queue (1) locally transactions with $ts > lts$ and (2) remotely according to each cluster's last updated transaction. The certifier's gts counter is updated for each transaction distributed locally. Finished the update, certification is unblocked.

4.2.3 Multiple Failures

The WICE protocol shall tolerate situations where multiple servers or entire clusters can fail abruptly. Most failure scenarios can be handled using a combination of the procedure for single servers. To avoid blocking during synchronization, we assume that all running synchronization routines

are restarted if a delegate fails.

The only scenario which requires special treatment is the loss of an entire cluster. In that case, the other clusters must be informed as soon as possible to allow blocking current and future transactions to become stable. A handler for this event is illustrated in Figure 3c.

5. Evaluation

In replication protocols that rely on a system-wide uniform atomic broadcast, updates cannot be applied before their carrier message has been delivered (and acknowledged) by all sites. This means that a full round-trip to the most distant site $2 \cdot t_{max}$ is required before updates can be installed, regardless of the location of the initiator. As the probability of two concurrent transactions conflicting depends on the latency, this has a profound impact in the abort rate of certification based protocols such as DBSM and Postgres-R [11].

In WICE, and considering two clusters C_A and C_B , total ordering of messages is performed using a sequencer sited, say, in cluster C_A , also referred to as the primary cluster. The updates of each update transaction can be installed as soon as the certification result is known but they are made visible to clients only after stabilization. Thus, it makes sense to distinguish between *install-interval* and *commit-interval*. Commit-interval denotes the time elapsed from the end of execution until the transaction gets committed at the originating site and is still lower bounded by $2 \cdot t_{max}$. The install-interval is the time elapsed from the moment the transaction finishes its optimistic execution until some site installs the incoming updates. Ignoring intra-cluster latency, and considering transactions originated at C_A , the install-interval is negligible for servers in cluster C_A and close to t_{max} in cluster C_B . On the other hand, for transactions originating in cluster C_B , the install-interval will be close to t_{max} and $2 \cdot t_{max}$, for C_A and C_B respectively.

The most significant advantage of the WICE protocol when compared to DBSM in a wide area network should therefore be its impact on the abort rate due to early delivery of updates. In this section, we experimentally verify this claim.

5.1. Experimental Environment

Experimental evaluation is conducted by running an actual implementation of the protocol within a simulated environment. By profiling real components with CPU cycle counters, the technique captures the actual overhead introduced by protocols [25]. By fine tuning the simulation components to accurately reproduce real components, it realistically reproduces results of real distributed systems [27]. When compared to testing in a real setting, this allows a tight control over experimental conditions, with advantages in repeatability and observability. The approach has been

Transaction Name	Empirical Distribution	Estimators	
Delivery	normal	mean=143.70	sd=2.33
Neworder	uniform	min=6.45	max=16.83
Orderstatus	normal	mean=1.66	sd=0.83
Stocklevel	uniform	min=1.85	max=2.33
Payment	normal	mean=2.26	sd=0.21

Table 1: CPU Times distributions (milliseconds).

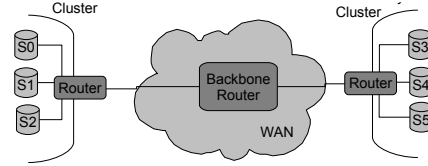


Figure 4: Network Topology.

previously used to evaluate database replication protocols both in LANs and WANs [11]. In detail, we use simulated database clients, database engines and networks, and real implementations of replication and group communication protocols.

The workload generator is configured according to the industry standard on-line transaction processing benchmark TPC-C [29]. Briefly, a wholesale supplier with a number of geographically distributed sales districts and associated warehouses. This workload is update intensive, as 92% of the transactions perform updates. It is also varied, as the *delivery* transaction takes a considerable amount of CPU time and has a very large read-set. The *payment* transaction is likely to produce Write-Write conflicts. The *neworder* transaction is short-lived and with higher locality.

The results thus vary according to the platform used for calibration of the simulated environment [27]. Results presented in this paper refer to the following hardware configuration: Each server has a single CPU AMD Opteron 250 running at 2.4GHz, 4GB RAM and a RAID 5 SATA disk array with fibre attachment. Transaction processing engines and overheads are configured according to PostgreSQL 8.0. Storage throughput as measured at the transaction log is 40MBps. CPU overheads are presented in Table 1 along with the corresponding generator distribution and estimators parameters. With properly configured indexes and within the range of presented results, it was verified that these are independent of the size of the database, as dictated by TPC-C scaling rules. Note also that these values do not include contention, as when blocked waiting for a resource processes are not scheduled. Also according to PostgreSQL 8.0, transaction processing engines use a multi-version concurrency control approach.

In our target scenario, 3 database servers are positioned at each of two different sites, as shown in Figure 4. The network simulator is configured as a pair of switched 1Gbps Ethernet local area networks, connected by a dedicated T3 link (45Mbps) with 400ms round-trip latency, representative of an inter-continental satellite link. As a baseline, we present also results obtained when configuring all 6 servers

Figure 3a: Delegate failover

```

New delegate: Synchronization request
1 upon site is initialized as new delegate
2 rsend(DELEGATE_SYNC, lts) to certifier()
3 foreach cluster in remoteClusters() do
4   rsend(ACK_SYNC, lts, TRUE) to cluster
5 end

Certifier: Handle synchronization request
6 upon rdeliver(DELEGATE_SYNC, clts) from cluster
7   foreach (ctid, cts, cws, cwv) in certified do
8     if cts > clts then
9       fifo_rsend(UPDATE_REM, ctid, cts, cws,
10        cwv) to cluster
11 end

All delegates: Synchronize ACK's
11 upon rdeliver(ACK_SYNC, clts, reply) from cluster
12   foreach (utid, uts, uws, uwv) in updated do
13     acked = {}
14     if clts ≥ uts then
15       add (utid, cluster) to acks
16       foreach (utid, c) in acks do
17         add c to acked
18       if remoteClusters() ⊆ acked then
19         u_vscast(STABLE_REM, utid)
20     if reply == TRUE then
21       rsend(ACK_SYNC, lts, FALSE) to cluster
22 end

```

Figure 3b: Certifier failover

```

Global site variables
1 synch = []

New certifier: Synchronization request
2 upon site is initialized as the new certifier
3   blockCertification()
4   foreach cluster in remoteClusters() do
5     rsend(CERTSYNC_REQUEST, lts) to cluster
6   end

All delegates: Send missing transactions
7 upon rdeliver(CERTSYNC_REQUEST, clts) from certifier
8   missing = []
9   foreach (tid, ts, ws, wv) in updated do
10    if ts > clts then
11      enqueue (tid, ts, ws, wv) to missing
12    rsend(CERTSYNC_REPLY, lts, missing) to certifier
13 end

```

Certifier: Missing updates

```

14 upon rdeliver(CERTSYNC_REPLY, clts, missing) from cluster
15   synched = {}
16   foreach (tid, ts, ws, wv) in missing do
17     if (tid, ts, ws, wv) ∉ certified then
18       enqueue (tid, ts, ws, wv) to certified
19   add (cluster, clts) to synch;
20   foreach (c, ts) in synch do
21     add c to synched;
22   if remoteClusters() ⊆ synched then
23     gts = lts
24     foreach (tid, ts, ws, wv) in certified do
25       if (ts > lts) then
26         gts = gts + 1
27         fifo_u_vscast(UPDATE_LOC, tid,
28          ts, ws, wv)
29     foreach (cluster, clts) in synch do
30       if ts > clts then
31         fifo_rsend(UPDATE_REM,
32          tid, ts, ws, wv) to cluster
33 end

```

Figure 3c:

All delegates: On failure of remote cluster

```

1 upon failure notification of cluster C
2   foreach (tid, ts, ws, wv) in updated do
3     acked = {}
4     foreach (tid, c) in acks do
5       add c to acked
6     if remoteClusters() ⊆ acked then
7       u_vscast(STABLE_REM, utid)
8 end

```

Figure 3: Failover handlers

within the same local area network.

In all scenarios, we vary the number of simulated clients from 60 to 6000, equally spread by all servers. We also take advantage of the locality in TPC-C: Clients associated with the same warehouse are connected to the same server to exploit locality, as suggested by the TPC-C specification. Note however, that with a small probability any client updates records associated with any warehouse.

5.2. Performance Results

The performance of the WICE protocol is evaluated by observing the throughput, latency and abort rate achieved when compared with plain DBSM. As a baseline, we present results obtained by grouping all 6 servers in the same cluster (DBSM CLUSTER). The results, obtained with Write-Write conflict certification (achieving 1-SI), are presented in Figure 5. Results are presented separately for each cluster.

The first interesting observation from the baseline protocol (DBSM CLUSTER) is that the capacity of the system is exhausted with 6000 clients. This shows up as throughput peaking (Figure 5(a)), increasing latency due to queuing (Figure 5(b)), and abort rate due to increased concurrency (Figure 5(c)). By examining resource usage logs one concludes that this is due to saturation of available CPU time. We should thus focus on system behavior up to 4000 clients, as a properly configured system will perform flow control to ensure operation in that range. Throughput grows linearly, latency is approximately constant and the abort rate negligible.

Then, we turn our attention to DBSM in the target scenario. Although throughput scalability is apparently close to linear, it is misleading as it corresponds to a high abort rate and a linearly increasing latency, in particular in cluster

C_B (Figures 5(e) and 5(f)). Both are explained by the same phenomenon: As locks are withheld during wide area stabilization, queuing delays arise, thus proportionally increasing the probability of later being aborted. Aborted transactions have to be resubmitted by the application, thus further loading the system. It is also important to underline that, as expected, latency and abort rate impact both clusters equally as both suffer with the same $2 \cdot t_{max}$ commit-interval.

As expected, the WICE protocol improves the performance at the primary cluster without negatively impacting secondary clusters. Namely, in the primary cluster the abort rate is negligible (Figure 5(c)), comparable only with the DBSM CLUSTER scenario. The latency is also approximately constant in the safe operating range (i.e., up to 4000 clients), although impacted by the round-trip over the wide area link (Figure 5(b)). Note however that such impact is very close to the absolute minimum of $2 \cdot t_{max}$ at 400 ms.

Also as expected, the abort rate of transactions initiated in the second cluster, which are impacted by a t_{max} to $2 \cdot t_{max}$ commit-interval, is not negligible although still offering a substantial improvement on DBSM. In the next section, we discuss the impact of this in the expected usage scenario of WICE.

5.3. Discussion

The workload assignment used in the previous section deserves some additional comments. The WICE protocol targets the global enterprise where the goal of replication is twofold. First, by providing a cluster for each region of the globe one avoids having to route all queries to a central location and thus avoid imposing the large latency on clients when no updates are performed, while at the same time balancing the load. Second, it improves availability as even catastrophic disasters can only impact the computing

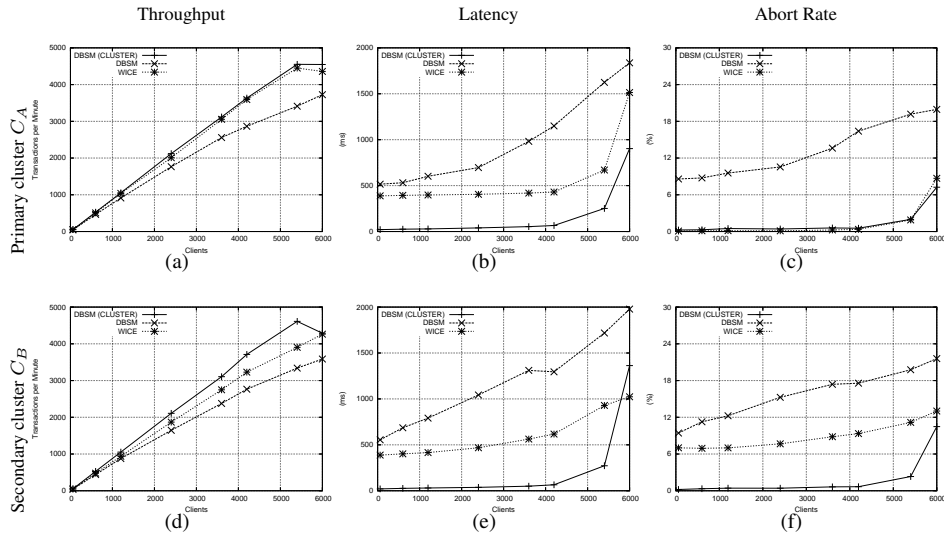


Figure 5: Performance results with 1-SI.

or communication infrastructure at a single location. One has therefore to consider clusters located in different time-zones, having distinct peak utilization periods.

This means that the evaluation scenario in the previous section, in which traffic in both clusters is exactly the same, is the worst case scenario for the proposed protocol. In reality, one should be able to migrate the centralized sequencer to the currently most loaded cluster. The additional abort rate at other locations can then be easily solved by resubmission, as these clusters are off peak and thus with underutilized resources.

We also have not assumed that resubmission can be done automatically by the database management system. However, this is true for many workloads, especially in current multi-tiered applications. By taking advantage of such option one could thus completely mask the abort rate at secondary clusters.

6. Conclusion

Eager update-everywhere database replication optimized for interconnected clusters in wide area networks is a valuable contribution to the infrastructure of the global enterprise. By providing the ability to locally serve clients it improves performance and by allowing failover ensures disaster recovery with no data loss. This is a hard problem, which existing commercial solutions address either by admitting some data loss or by centralizing update processing.

The proposed WICE protocol shows how to scale replication protocols based on group communication to a wide area setting with increased performance, while at the same time increasing their practicality. This is achieved by restricting group communication within clusters and using a simple peer protocol over long distance links. The evaluation performed in a realistic platform illustrates the advantages of the approach, namely, linear throughput scalability,

up to 2 times less latency and a negligible abort rate at the cluster supporting the region currently generating the most traffic.

References

- [1] M. Abdallah, R. Guerraoui, and P. Pucheral. Dictatorial transaction processing: Atomic commitment without veto right. *Distributed Parallel Databases*, 2002.
- [2] Y. Amir, C. Danilov, M. Miskin-Amir, J. Stanton, and C. Tutu. Practical wide area database replication. Technical report, Center for Networking and Distributed Systems, Computer Science Department, Johns Hopkins University, 2002.
- [3] T. Anderson, Y. Breitbart, H. F. Korth, and A. Woo. Replication, consistency, and practicality: are these mutually exclusive? In *Proceedings of ACM SIGMOD international conference on Management of data*, 1998.
- [4] Apache. Apache derby. <http://db.apache.org/derby>.
- [5] H. Berenson, P. Bernstein, J. Gray, J. Melton, E. O’Neil, and P. O’Neil. A critique of ansi sql isolation levels. In *Proceedings of the 1995 ACM SIGMOD international conference on Management of data*, 1995.
- [6] P. Bernstein, V. Hadzilacos, and N. Goodman. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, 1987.
- [7] M. J. Carey and M. Livny. Conflict detection tradeoffs for replicated data. *ACM Transactions on Database Systems*, 1991.
- [8] T. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of the ACM*, 43(2), Mar. 1996.
- [9] G. V. Chockler, I. Keidar, and R. Vitenberg. Group communication specifications: a comprehensive study. *ACM Computing Surveys*, 2001.
- [10] G. Consortium. Gorda - open replication of databases. <http://gorda.di.uminho.pt/consortium>, October 2004.
- [11] A. Correia Jr., A. Sousa, L. Soares, J. Pereira, R. Oliveira, and F. Moura. Group-based replication of on-line transaction processing servers. In *Dependable Computing: Second Latin-American Symposium*, 2005.

- [12] S. Elnikety, F. Pedone, and W. Zwaenepoel. Database replication using generalized snapshot isolation. In *Proceedings of The 24th IEEE Symposium on Reliable Distributed Systems*, 2005.
- [13] M. Fischer, N. Lynch, and M. Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM*, 1985.
- [14] R. Guerraoui and A. Schiper. Software-based replication for fault tolerance. *IEEE Computer*, 1997.
- [15] V. Hadzilacos and S. Toueg. A modular approach to fault-tolerant broadcasts and related problems. Technical report, Cornell University, 1994.
- [16] P. Inc. Postgresql. <http://www.postgresql.org>.
- [17] B. Kemme and G. Alonso. Don't be lazy, be consistent: Postgres-R, A new way to implement database replication. In *Proceedings of 26th International Conference on Very Large Data Bases*, 2000.
- [18] H. T. Kung and J. T. Robinson. On optimistic methods for concurrency control. *ACM Transactions on Database Systems*, 1981.
- [19] Y. Lin, B. Kemme, M. Patiño-Martínez, and R. Jiménez-Peris. Consistent data replication: Is it feasible in wans? In *Euro-Par*, 2005.
- [20] H. Miranda, A. Pinto, and L. Rodrigues. Appia, a flexible protocol kernel supporting multiple coordinated channels. In *Proceedings of The IEEE 21st International Conference on Distributed Computing Systems*, 2001.
- [21] F. Pedone, R. Guerraoui, and A. Schiper. The database state machine approach. *Distributed Parallel Databases*, 2003.
- [22] J. Pereira, A. C. Jr., N. Carvalho, S. Guedes, R. Oliveira, and L. Rodrigues. Database interfaces for replication support. Technical report, Universidade do Minho/Faculdade de Ciências da Universidade de Lisboa, 2006.
- [23] L. Rodrigues, J. Mocito, and N. Carvalho. From spontaneous total order to uniform total order: different degrees of optimistic delivery. In *In Proceedings of the 21st ACM Symposium on Applied Computing*, 2006.
- [24] F. Schneider. Replication management using the state-machine approach. In *Distributed Systems*. ACM Press/Addison-Wesley Publishing Co., 1993.
- [25] L. Soares and J. Pereira. Experimental performability evaluation of middleware for large-scale distributed systems. In *7th International Workshop on Performability Modeling of Computer and Communication Systems*, 2005.
- [26] A. Sousa, J. Pereira, F. Moura, and R. Oliveira. Optimistic total order in wide area networks. In *Proceedings of The 21st Symposium on Reliable Distributed Systems*, 2002.
- [27] A. Sousa, J. Pereira, L. Soares, A. C. Jr., L. Rocha, R. Oliveira, and F. Moura. Testing the dependability and performance of group communication based database replication protocols. In *IEEE International Conference on Dependable Systems and Networks - Performance and Dependability Symposium*, 2005.
- [28] Symantec. Veritas backup software. <http://www.symantec.com/enterprise/veritas/index.jsp>.
- [29] T. P. P. C. (TPC). TPC benchmarkTM C standard specification revision 5.0, Feb. 2001.
- [30] M. Tamma. *Oracle Streams - High Speed Replication and Data Sharing*. Rampant TechPress, 2004.
- [31] S. Wu and B. Kemme. Postgres-r(si): Combining replica control with concurrency control based on snapshot isolation. In *International Conference on Data Engineering*, 2005.