

Replication

CSE593 Transaction Processing

Philip A. Bernstein
Copyright ©2001 Philip A. Bernstein

3/7/01

1

Outline

1. Introduction
2. Primary-Copy Replication
3. Multi-Master Replication
4. Other Approaches

3/7/01

2

1. Introduction

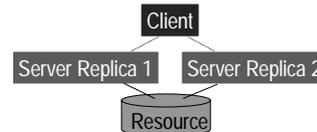
- Replication - using multiple copies of a server (called replicas) for better availability and performance.
- If you're not careful, replication can lead to
 - worse performance - updates must be applied to all replicas and synchronized
 - worse availability - some algorithms require multiple replicas to be operational for any of them to be used

3/7/01

3

Replicated Server

- Can replicate servers on a common resource
 - Data sharing - DB servers communicate with shared disk



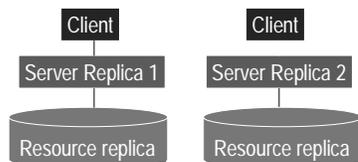
- Helps availability in primary-backup scenario
- Requires replica cache coherence mechanism ...
- Hence, this helps performance only if
 - little conflict between transactions at different servers or
 - loose coherence guarantees (e.g. read committed)

3/7/01

4

Replicated Resource

- To get more improvement in availability, replicate the resources (too)
- Also increases potential throughput
- This is what's usually meant by replication
- It's the scenario we'll focus on

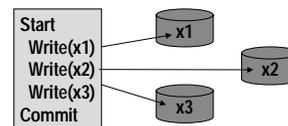


3/7/01

5

Synchronous Replication

- Replicas function just like non-replicated servers
- Synchronous replication - transaction updates all replicas of every item it updates



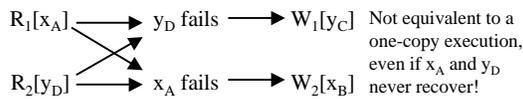
- Issues
 - Too expensive for most applications, due to heavy distributed transaction load (2-phase commit)
 - Can't control when updates are applied to replicas

3/7/01

6

Synchronous Replication - Issues

- If you just use transactions, availability suffers.
- For high-availability, the algorithms are complex and expensive, because they requires heavy-duty synchronization of failures.
- ... of failures? How do you synchronize failures?



- DBMS products support it only in special situations

3/7/01

7

Asynchronous Replication

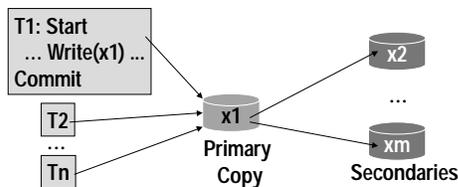
- Asynchronous replication
 - Each transaction updates one replica.
 - Updates are propagated later to other replicas.
- Primary copy: All transactions update the same copy
- Multi-master: Transactions update different copies
 - Useful for disconnected operation, partitioned network
- Both approaches ensure that
 - Updates propagate to all replicas
 - If new updates stop, replicas converge to the same state
- Only primary copy ensures serializability
 - Details later ...

3/7/01

8

2. Primary-Copy Replication

- Designate one replica as the primary copy (publisher)
- Transactions may update only the primary copy
- Updates to the primary are sent later to secondary replicas (subscribers) in the order they were applied to the primary



3/7/01

9

Asynchronous Update Propagation

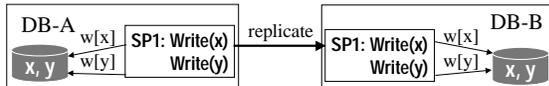
- Collect updates at primary using triggers or the log
- Triggers (Oracle, Rdb, SQL Server, DB2, ...)
- On every update at the primary, a trigger fires to store the update in the update propagation table.
- Post-process (“sniff”) the log to generate update propagations (SQL Server, DB2, Tandem Non-Stop SQL)
 - Off-line, so saves trigger and triggered update overhead, though R/W log synchronization also has a cost
 - Requires admin (what if the log sniffer fails?)
- Optionally identify updated fields to compress log
- Most DB systems support this today.
 - First in IBM IMS, Tandem NS SQL, DEC/Rdb, & ad hoc

3/7/01

10

Request Propagation

- Replicate a request rather than the updates produced by the request (e.g., a stored procedure call).



- Must ensure requests run in the same order at primary and replica (same requirement as updates).
 - Log the requests or extend triggers to capture them.
- Could run request synchronously at all replicas, but commit even if one replica fails.
 - Need a recovery procedure for failed replicas.

3/7/01

11

Products

- All major DBMS products have a rich primary-copy replication mechanism
- Differences are in detailed features
 - performance
 - ease of management
 - richness of filtering predicates
 - push vs. pull propagation
 - stored procedure support
 - transports (e.g. Sybase SQLAnywhere can use email!)
 - ...
- The following summary is an incomplete snapshot of products as of July 1999.

3/7/01

12

SQL Server 7.0

- Publication - a collection of articles to subscribe to
- Article – a horiz/vertical table slice or stored proc
 - Customizable table filter (WHERE clause or stored proc)
 - Stored proc may be transaction protected (replicate on commit). Replicates the requests instead of each update.
- *Snapshot replication* makes a copy
- *Transactional replication* maintains the copy by propagating updates from publisher to subscribers
 - Post-processes log to store updates in Distribution DB
 - Distribution DB may be separate from the publisher DB
 - Updates can be pushed to or pulled from subscriber
 - Can customize propagated updates using stored procs

3/7/01

13

SQL Server 7.0 (cont'd)

- *Immediate updating subscriber*
 - Can update data, synchronizing with publisher via 2PC
 - Uses triggers to capture updates (**Not For Replication** disables trigger for updates from the publisher)
 - Subscriber sends before/after row timestamp. Publisher checks row didn't change since subscriber's current copy
 - Publisher then forwards updates to other subscribers
- Access control lists protect publishers from unauthorized subscribers
- *Merge replication*- described later

3/7/01

14

Oracle 8i

- Like SQL Server, can replicate updates to table fragments or stored proc calls at the master copy
- Uses triggers to capture updates in a deferred queue
 - Updates are row-oriented, identified by primary key
 - Can optimize by sending keys and updated columns only
- Group updates by transaction, which are propagated:
 - Either serially in commit order or
 - in parallel with some dependent transaction ordering: each read reads the "commit number" of the data item; updates are ordered by dependent commit number
- Snapshots are updated in a batch refresh.
 - Pushed from master to snapshots, using queue scheduler

3/7/01

15

DB2

- Very similar feature set to SQL Server and Oracle
- Filtered subscriber (no stored proc replication (?))
 - Create snapshot, then update incrementally (push or pull)
- Captures DB2 updates from the DB2 log
 - For other systems, captures updates using triggers
- Many table type options:
 - Read-only snapshot copy, optionally with timestamp
 - Aggregates, with cumulative or incremental values
 - Consistent change data, optionally with row versions
 - "Replica" tables, for multi-master updating
- Interoperates with many third party DBMS's

3/7/01

16

Failure Handling

- Secondary failure - nothing to do till it recovers
 - At recovery, apply the updates it missed while down
 - Needs to determine which updates it missed, just like log-based recovery
 - If down for too long, it may be faster to get a whole copy
- Primary failure – Products just wait till it recovers
 - Can get higher availability by electing a new primary
 - A secondary that detects primary's failure announces a new election by broadcasting its unique replica identifier
 - Other secondaries reply with their replica identifier
 - The largest replica identifier wins

3/7/01

17

Failure Handling (cont'd)

- Primary failure (cont'd)
 - All replicas must now check that they have the same updates from the failed primary
 - During the election, each replica reports the id of the last log record it received from the primary
 - The most up-to-date replica sends its latest updates to (at least) the new primary.
 - Could still lose an update that committed at the primary and wasn't forwarded before the primary failed ... but solving it requires synchronous replication (2-phase commit to propagate updates to replicas)

3/7/01

18

Communications Failures

- Secondaries can't distinguish a primary failure from a communication failure that partitions the network.
- If the secondaries elect a new primary and the old primary is still running, there will be a reconciliation problem when they're reunited. This is multi-master.
- To avoid this, one partition must know it's the only one that can operate, and can't communicate with other partitions to figure this out.
- Could make a static decision.
The partition that has the primary wins.
- Dynamic solutions are based on Majority Consensus

3/7/01

19

Majority Consensus

- Whenever a set of communicating replicas detects a replica failure or recovery, they test if they have a majority (more than half) of the replicas.
- If so, they can elect a primary
- Only one set of replicas can have a majority.
- Doesn't work well with even number of copies.
 - Useless with 2 copies
- Quorum consensus
 - Give a weight to each replica
 - The replica set that has a majority of the weight wins
 - E.g. 2 replicas, one has weight 1, the other weight 2

3/7/01

20

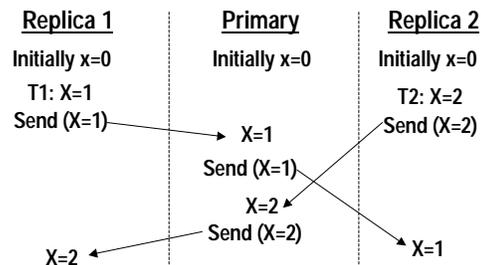
3. Multi-Master Replication

- Some systems must operate when partitioned.
 - Requires many updatable copies, not just one primary
 - Conflicting updates on different copies are detected late
- Classic example - salesperson's disconnected laptop
 - Customer table (rarely updated) Orders table (insert mostly)
 - Customer log table (append only)
 - So conflicting updates from different salespeople are rare
- Use primary-copy algorithm, with multiple masters
 - Each master exchanges updates ("gossips") with other replicas when it reconnects to the network
 - Conflicting updates require reconciliation (i.e. merging)
- In Lotus Notes, Access, SQL Server, Oracle, ...

3/7/01

21

Example of Conflicting Updates A Classic Race Condition



- Replicas end up in different states

3/7/01

22

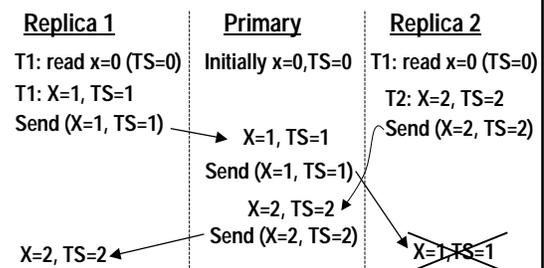
Thomas' Write Rule

- To ensure replicas end up in the same state
 - Tag each data item with a timestamp
 - A transaction updates the value and timestamp of data items (timestamps monotonically increase)
 - An update to a replica is applied only if the update's timestamp is greater than the data item's timestamp
 - You only need to keep timestamps of data items that were recently updated (where an older update could still be floating around the system)
- All multi-master products use some variation of this
- Robert Thomas, *ACM TODS*, June '79
 - Same article that invented majority consensus

3/7/01

23

Thomas Write Rule $\not\Rightarrow$ Serializability



- Replicas end in the same state, but neither T1 nor T2 reads the other's output, so the execution isn't serializable.

3/7/01

24

Multi-Master Performance

- The longer a replica is disconnected and performing updates, the more likely it will need reconciliation
- The amount of propagation activity increases with more replicas
 - If each replica is performing updates, the effect is quadratic

3/7/01

25

Microsoft Access and SQL Server

- Multi-master replication without a primary
- Each row R of a table has 4 additional columns
 - globally unique id (GUID)
 - generation number, to determine which updates from other replicas have been applied
 - version number = the number of updates to R
 - array of [replica, version number] pairs, identifying the largest version number it got for R from every other replica
- Uses Thomas' write rule, based on version numbers
 - Access uses replica id to break ties. SQL Server 7 uses subscriber priority or custom conflict resolution.

3/7/01

26

Generation Numbers (Access/SQL cont'd)

- Each replica has a current generation number
- A replica updates a row's generation number whenever it updates the row
- A replica knows the generation number it had when it last exchanged updates with R', for every replica R'.
- A replica increments its generation number every time it exchanges updates with another replica.
- So, when exchanging updates with R', it should send all rows with a generation number larger than what it had when last exchanging updates with R'.

3/7/01

27

Duplicate Updates (Access/SQL cont'd)

- Some rejected updates are saved for later analysis
- To identify duplicate updates to discard them -
 - When applying an update to x, replace x's array of [replica, version#] pairs by the update's array.
 - To avoid processing the same update via many paths, check version number of arriving update against the array
- Consider a rejected update to x at R from R', where
 - [R', V] describes R' in x's array, and
 - V' is the version number sent by R'.
 - If $V \geq V'$, then R saw R''s updates
 - If $V < V'$, then R didn't see R''s updates, so store it in the conflict table for later reconciliation

3/7/01

28

Oracle 8i (revisited)

- Masters replicate entire tables
 - Updates are pushed from master to masters and to snapshots (synchronous or asynchronous)
 - Updates include before values (you can disable if conflicts are impossible)
 - They recommend masters should always be connected
- Snapshots are updatable \Rightarrow "multi-master"
 - Each propagation transaction updates its queue entry (instead of update-oriented generation numbers)
- Conflict detection
 - Before-value at replica is different than in update
 - Uniqueness constraint is violated
 - Row with the update's key doesn't exist

3/7/01

29

Oracle 8i Conflict Resolution

- Built-in resolution strategies (defined per column-group)
 - Add difference between the old and new values of the originating site to the destination site
 - Average the value of the current site and the originating site
 - Min or max of the two values
 - The one with min or max timestamp
 - The site or value with maximum priority
 - Can apply methods in sequence: e.g., by time, then by priority.
- Can call custom procs to log, notify, or resolve the conflict
 - Parameters - update's before/after value and row's current value
- For a given update, if no built-in or custom conflict resolution applies, then the entire transaction is logged.

3/7/01

30

4. Other Approaches

- Non-transactional replication using timestamped updates and variations of Thomas' write rule
 - directory services are managed this way
- Quorum consensus per-transaction
 - Read and write a quorum of copies
 - Each data item has a version number and timestamp
 - Each read chooses a replica with largest version number
 - Each write increments version number one greater than any one it has seen
 - No special work needed during a failure or recovery

3/7/01

31

Other Approaches (cont'd)

- Read-one replica, write-all-available replicas
 - Requires careful management of failures and recoveries
- E.g., Virtual partition algorithm
 - Each node knows the nodes it can communicate with, called its view
 - Transaction T can execute if its home node has a view including a quorum of T's readset and writeset (i.e. the data it can read or write)
 - If a node fails or recovers, run a view formation protocol (much like an election protocol)
 - For each data item with a read quorum, read the latest version and update the others with smaller version #.

3/7/01

32

Summary

- State-of-the-art products have rich functionality.
 - It's a complicated world for app designers
 - Lots of options to choose from
- Most failover stories are weak
 - Fine for data warehousing
 - For 24x7 TP, need better integration with cluster node failover

3/7/01

33