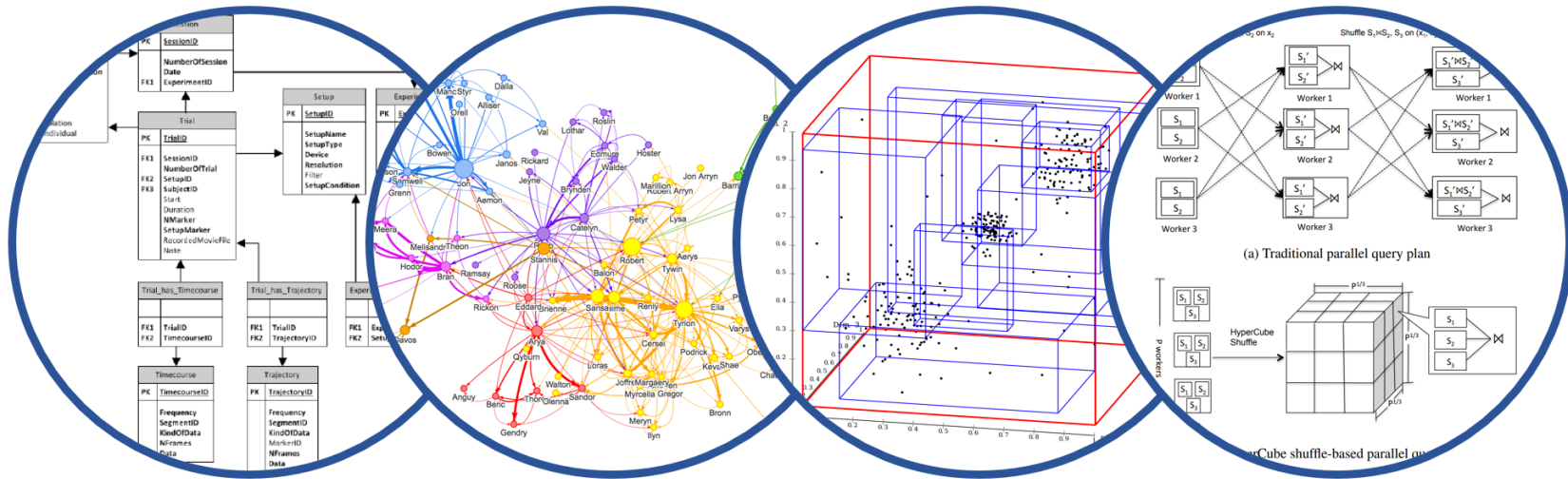


Course evals

Please take a few minutes to fill out the course evaluations:

<https://uw.iasystem.org/survey/258581>

And thank you all for your hard work this quarter!



Database System Internals




Two-tier Replication and NewSQL

Paul G. Allen School of Computer Science and Engineering
 University of Washington, Seattle

Synchronous Replication Properties

- Favours **consistency** over availability
 - Only majority partition can process requests
 - There appears to be a single copy of the db
- **High runtime overhead**
 - Must lock and update at least majority of replicas
 - Two-phase commit
 - Runs at pace of slowest replica in quorum
 - So overall system is now slower
 - Higher deadlock rate (transactions take longer)

Types of Replication

	Master	Group
Synchronous		
Asynchronous		

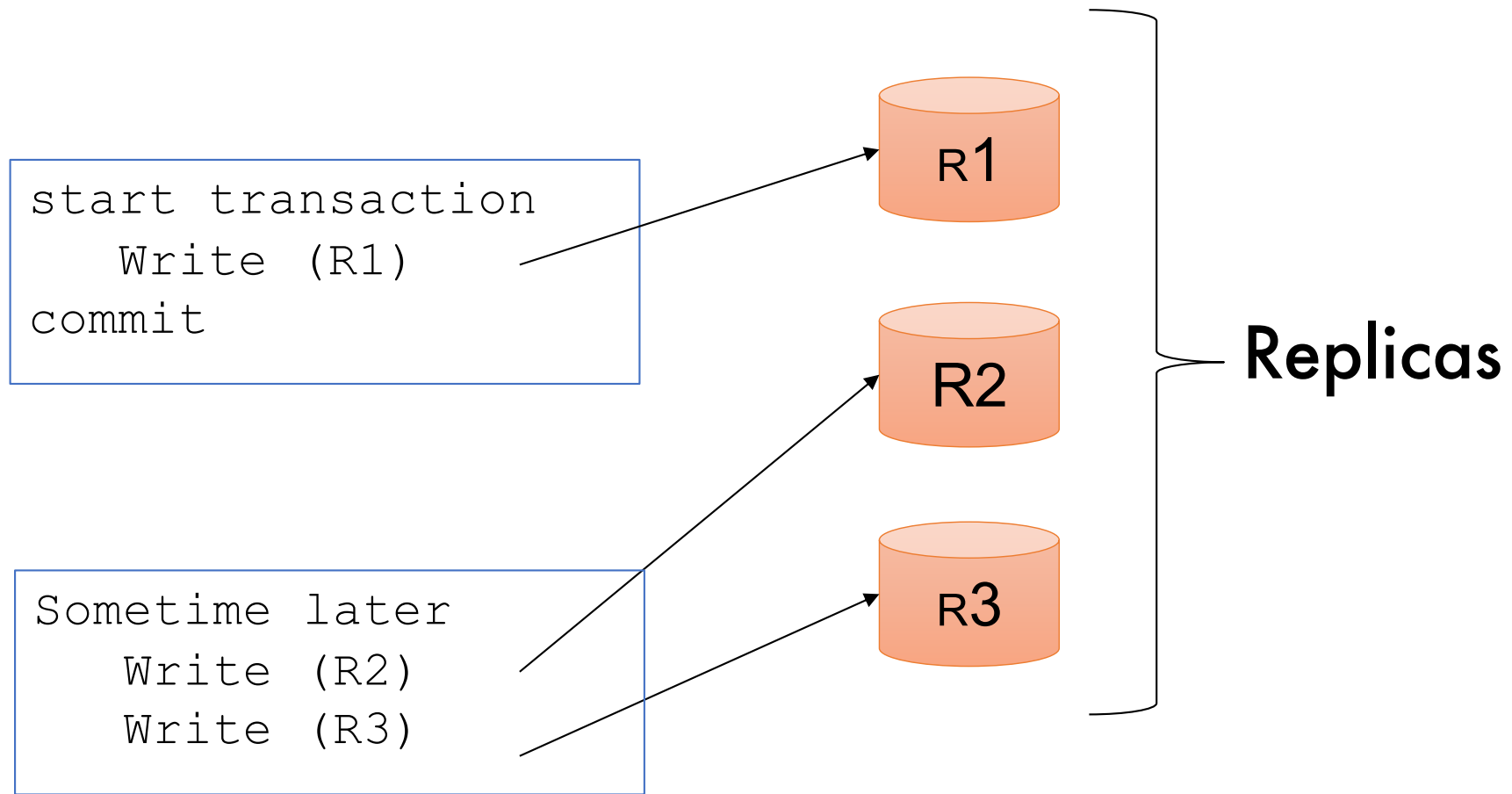
Asynchronous Replication

- Also called **lazy replication**
- Also called **optimistic replication**

- Main goals: availability and performance

- Approach
 - One replica updated by original transaction
 - Updates propagate asynchronously to other replicas

Asynchronous Replication



Asynchronous Master Replication

One master holds primary copy

- Transactions update primary copy
- Master asynchronously propagates updates to replicas, which process them in same order
E.g. through **log shipping**
- Ensures single-copy serializability

What happens when master/primary fails?

- Can lose most recent transactions when primary fails!
- After electing a new primary, secondaries must agree who is most up-to-date

Discussion: Log Shipping

A general problem:

- A master operates on a database
- The DB needs to be replicated to one or several replicas (e.g. hot stand-by databases)

Discussion: Log Shipping

A general problem:

- A master operates on a database
- The DB needs to be replicated to one or several replicas (e.g. hot stand-by databases)
- Log Shipping Technique

Discussion: Log Shipping

A general problem:





- A master operates on a database
- The DB needs to be replicated to one or several replicas (e.g. hot stand-by databases)
- Log Shipping Technique:
 - Master node ships the tail of the log to the replicas
E.g. when it flushes the log tail to disk
 - Replicas REDO the log; this is very efficient
 - Need very little systems development: we create the log anyway, and we have the REDO function anyway

Discussion: Log Shipping

A general problem:

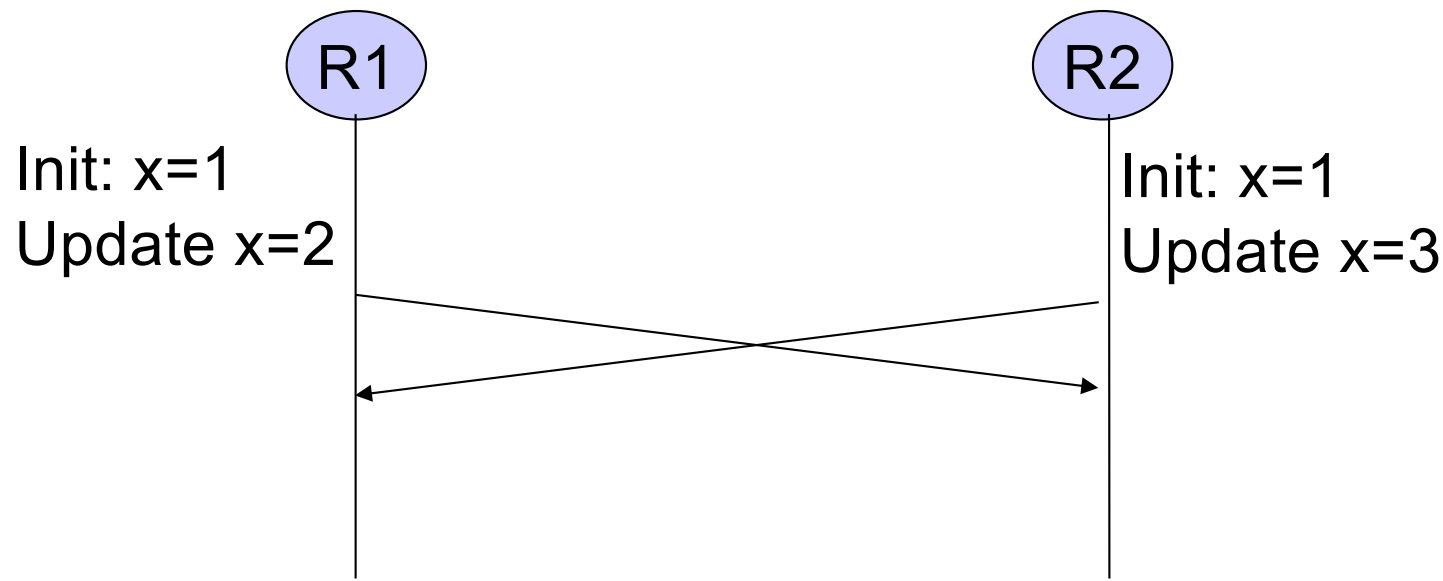
- A master operates on a database
- The DB needs to be replicated to one or several replicas (e.g. hot stand-by databases)
- Log Shipping Technique:
 - Master node ships the tail of the log to the replicas
E.g. when it flushes the log tail to disk
 - Replicas REDO the log; this is very efficient
 - Need very little systems development: we create the log anyway, and we have the REDO function anyway
 - Complications due to the need to "remove" updates of active transactions (they may later abort)

Types of Replication

	Master	Group
Synchronous		
Asynchronous		

Asynchronous Group Replication

- Also called **multi-master**
- Best scheme for availability
- **Cannot guarantee one-copy serializability!**



Asynchronous Group Replication

- **Cannot guarantee one-copy serializability!**
- **Instead guarantee convergence**
 - Db state does not reflect any serial execution
 - But all replicas have the same state
- **Called “Eventual Consistency” = if the DB stops operations, then eventually all copies are equal**

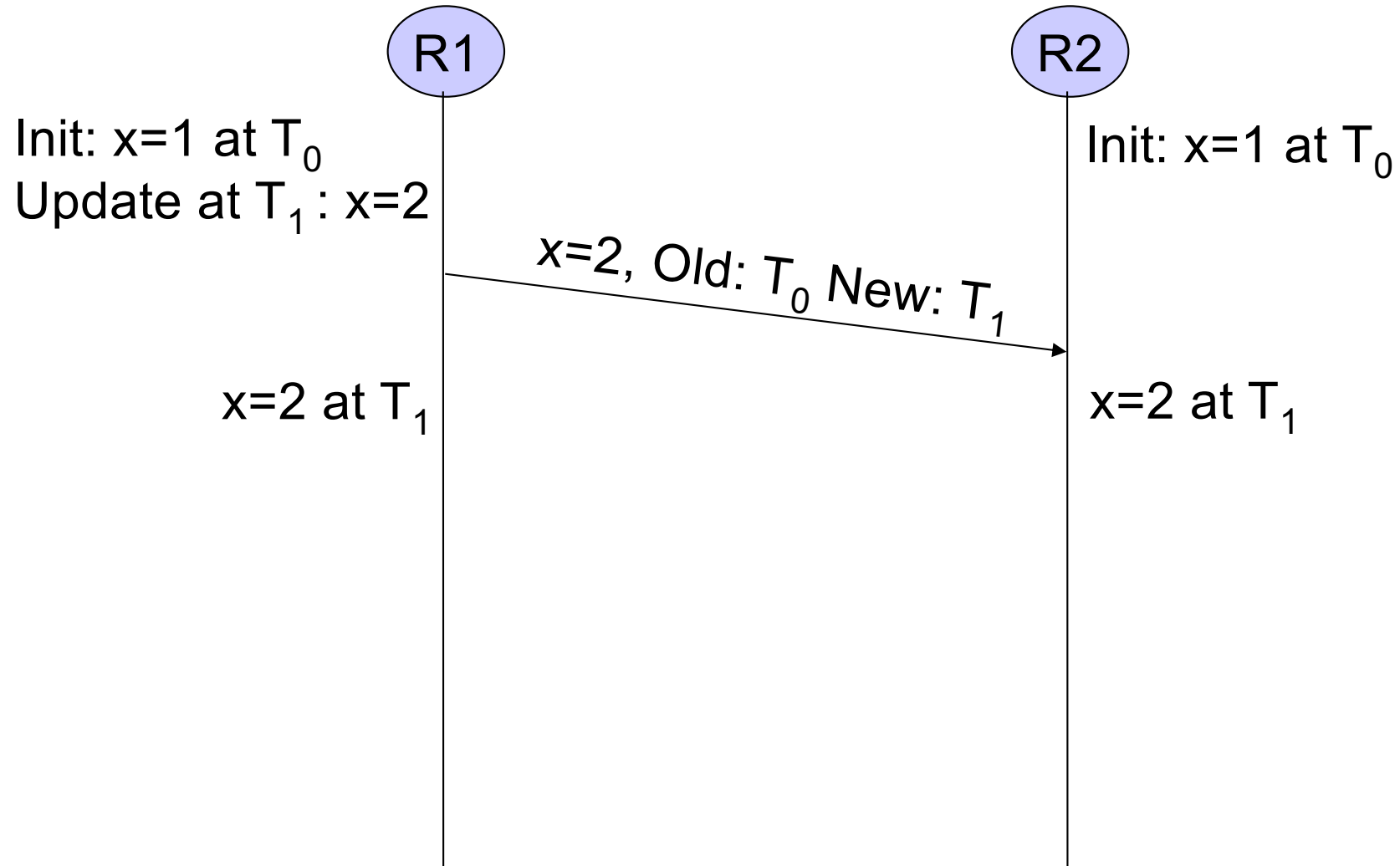
Asynchronous Group Replication

- **Cannot guarantee one-copy serializability!**
- **Instead guarantee convergence**
 - Db state does not reflect any serial execution
 - But all replicas have the same state
- **Called “Eventual Consistency” = if the DB stops operations, then eventually all copies are equal**
- **Detect conflicts and reconcile replica states**

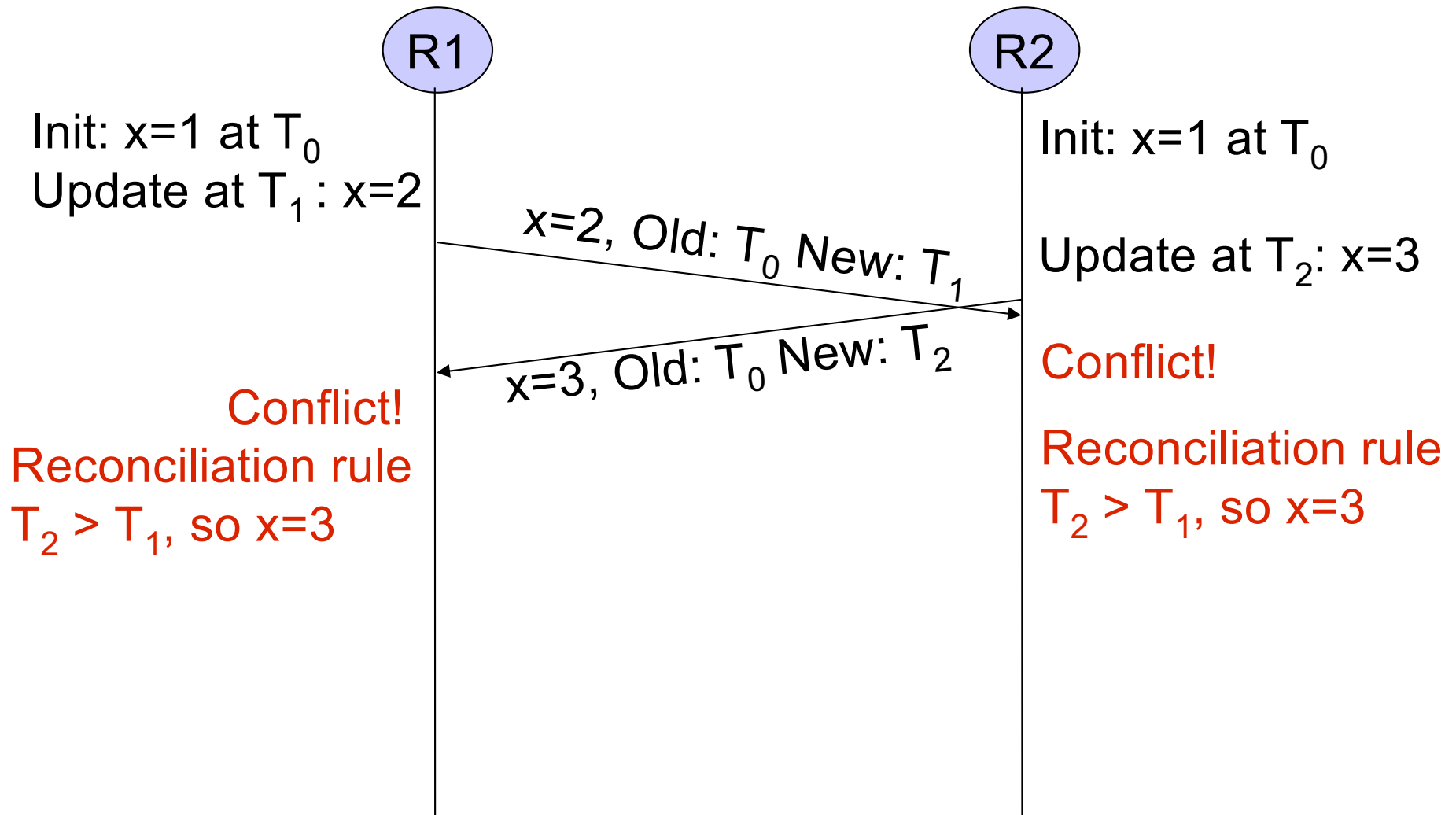
Asynchronous Group Replication

- **Cannot guarantee one-copy serializability!**
- **Instead guarantee convergence**
 - Db state does not reflect any serial execution
 - But all replicas have the same state
- **Called “Eventual Consistency” = if the DB stops operations, then eventually all copies are equal**
- **Detect conflicts and reconcile replica states**
- **Reconciliation techniques:**
 - Most recent timestamp wins
 - Site A wins over site B
 - But also: user-defined rules, or even manual

Detecting Conflicts Using Timestamps



Detecting Conflicts Using Timestamps



Asynchronous Group Replication Properties

- Favours **availability** over consistency
 - Can read and update any replica
 - High runtime performance
- **Weak consistency**
 - Conflicts and reconciliation

Outline

- Goals of replication
- Three types of replication
 - Synchronous (aka eager) replication
 - Asynchronous (aka lazy) replication
 - Two-tier replication

Two-Tier Replication

- **Benefits of lazy master and lazy group**
- Each object has a master with primary copy
- When disconnected from master
 - Secondary can only run **tentative transactions**
- When reconnects to master
 - Master reprocesses all tentative transactions
 - Checks an acceptance criterion
 - If passes, we now have **final commit order**
 - Secondary **undoes tentative and redoes committed**

Conclusion

- **Replication is a very important problem**
 - Fault-tolerance (various forms of replication)
 - Caching (lazy master)
 - Warehousing (lazy master)
 - Mobility (two-tier techniques)
- **Replication is complex, but basic techniques and trade-offs are **very well known****
 - Synchronous or asynchronous replication
 - Master or quorum

SCALABILITY

HIGH
(Many Nodes)

NOSQL

NEWSQL

LOW
(One Node)

TRADITIONAL

WEAK
(None/Limited)

GUARANTEES

STRONG
(ACID)

Slide from Andy Pavlo @ CMU

Some Popular NewSQL Systems

■ H-Store

- Research system from Brown U., MIT, CMU, and Yale
- Commercialized as VoltDB

■ Hekaton

- Microsoft
- Fully integrated into SQL Server

■ Hyper

- Hybrid OLTP/OLAP
- Research system from TU Munich. Bought by Tableau

■ Spanner

- Google

Hekaton

- Focus: DBMS with large main memories and many core CPUs
- Integrated with SQL Server
- Key user-visible features
 - Simply declare a table “memory resident”
 - Hekaton tables are fully durable and transactional, though non-durable tables are also supported
 - Query can touch both Hekaton and regular tables

Hekaton Key Details

- **Idea: To increase transaction throughput must decrease number of instructions / transaction**
- **Main-memory DBMS**
 - Optimize indexes for memory-resident data
 - Durability by logging and checkpointing records to external storage
- **No partitioning**
 - Any thread can touch any row of any table
- **No locking**
 - Uses a new MVCC method for isolation

Hekaton More Details

- **Optimized stored procedures**
 - Compile statements and stored procedures into customized, highly efficient machine code

- Hybrid OLTP and OLAP
- In-memory data management
 - Including optimized indexes for memory-resident data
 - Data compression for cold data
- Data-centric code generation
 - SQL translated to LLVM
- OLAP separated from OLTP using MVCC
- Exploits hardware transactional memory
- Data shuffling and distribution optimizations

Conclusion

- **Many innovations recently in**
 - Big data analytics
 - Transaction processing at very large scale
- **Many more problems remain open**
- **This course teaches foundations**
- **Innovate with an open mind!**