Lazy Evaluation for Concurrent OLTP and Bulk Transactions

Lesley Wevers Marieke Huisman Maurice van Keulen

March 4, 2014

Bitcoin bank Flexcoin shuts down after massive theft

Some \$600,000 worth of bitcoins were stolen from the site

March 5, 2014

Yet another exchange hacked: Poloniex loses around \$50,000 in bitcoin

Firm can't cover losses: "All balances will temporarily be deducted by 12.3%."

On March 4th, 2014, about 12.3% of the BTC on Poloniex was stolen.

How Did It Happen?

The hacker found a vulnerability in the code that takes withdrawals. Here's what happens when you place a withdrawal:

- 1. Input validation.
- 2. Your balance is checked to see if you have enough funds.
- 3. If you do, your balance is deducted.
- 4. The withdrawal is inserted into the database.
- 5. The confirmation email is sent.

6. After you confirm the withdrawal, the withdrawal daemon picks it up and processes the withdrawal.

The hacker discovered that if you place several withdrawals all in practically the same instant, they will get processed at more or less the same time. This will result in a negative balance, but valid insertions into the database, which then get picked up by the withdrawal daemon.

	balance
Alice	10







	balance	
Alice	-50	



Transactions & Correctness

Databases provide **transactions** to ensure correct execution of database operations.

Correctness guarantees:

def withdraw(account, amount):
 atomic:
 if(balance[account] < amount)</pre>

abort

balance[account] -= amount
withdrawals += (account, amount)

Atomicity Consistency Isolation Durability

Weak Isolation

High volume database systems often trade correctness for performance.



Weak Isolation

High volume database systems often trade correctness for performance.



Can we make transactions faster without sacrificing correctness?

Transactional Workloads in Database Systems

	Variety	Footprint	Access pattern
OLTP Online transaction processing	Regular	Small	Read / write
OLAP Online analytical processing	Ad-hoc	Large / Everything	Read-only

Transactional Workloads in Database Systems

	Variety	Footprint	Access pattern
OLTP Online transaction processing	Regular	Small	Read / write
OLAP Online analytical processing	Ad-hoc	Large / Everything	Read-only
OLBP? Online bulk processing	Ad-hoc	Large / Everything	Read / write

Concurrency Control

Two-phase locking

- Lock data elements while executing the transaction.
- Do not release locks until the transaction has finished.

Optimistic concurrency control

- Check validity of a transaction after executing.
- Abort and re-execute if inconsistent.

Multi-version concurrency control

• Perform read-only transactions on a snapshot.

Multi-Core Scalability



The YCSB benchmark under high contention.

Xiangyao Yu et al. Staring Into The Abyss: An Evaluation Of Concurrency Control With One Thousand Cores. VLDB 2014.

Runtime Breakdown



Xiangyao Yu et al. Staring Into The Abyss: An Evaluation Of Concurrency Control With One Thousand Cores. VLDB 2014.

Contention

Sources of Contention in OLTP transactions

- Counters
- Summary data
- Popular items
- Newest / oldest items
- Correlated access
- Zipfian access skew

Bulk Updates

Bulk updates content with most or all concurrent transactions

• Bulk operations are often required in schema transformations

WikiMedia schema revisions:



- 90% require a write lock.
- Largest took 22 hours to complete for wikipedia.

Curino et al. http://yellowstone.cs.ucla.edu/schema-evolution/index.php/Schema_Evolution_Benchmark

Bulk Updates



Bulk updates block all concurrent transactions, or are executed incorrectly.

L. Wevers et al. A Benchmark for Online Non-Blocking Schema Transformations. DATA 2015. L. Wevers et al. Analysis of the Blocking Behaviour of Schema Transformations in Relational Database Systems. ADBIS 2015.

Contributions

Existing concurrency control mechanisms:

- Cannot execute contended workloads concurrently
- Block concurrent transactions during bulk operations

We have developed a concurrency control mechanism that can:

- Significantly improve parallelism in contended workloads
- Execute bulk operations without blocking

Lazy Transactions

Lazy Transactions

Main Ideas

- Re-order operations within a transaction based on demand from readers
- Use results from partially executed transactions

Mechanism: move responsibility of evaluation from writers to readers

- Transactions commit before executing their operations
- Readers can prioritize evaluation of those operations that they need



Split a transaction into update functions on variables.







Transaction 3







Transaction 2



Transaction 3















Transaction 3





Encoding Transactions

Transaction Execution Phases

1. Preparation Phase

- Prepare update functions
- **Optimistic reads** on a snapshot

2. Commit Phase

- Queue update functions
- Blocking reads on the current state

3. Lazy Phase

- Work performed inside update functions
- Lazy reads on a snapshot

4. Result Phase

• Compute observable results from a snapshot



commit:

update(y, get(x))

commit:

s <- snapshot update(y, v => s.get(x))

Optimistic Reads

commit:

update(get(y), a)

prepare:

x <- get(y)

commit:

s <- snapshot update(x, v => if(x == s.get(y)) a else v

Reading transaction results

commit:

update(x, v) return get(x) + get(y) + get(z)

commit:

update(x, v) s <- snapshot

result:

return s.get(x) + s.get(y) + s.get(z)

Read-only

prepare:

```
s <- snapshot
return s.get(x) + s.get(y) + s.get(z)
```











How do we deal with errors?

Aborts

If a transaction **aborts**: all update functions must return the original value If a transaction **commits**: all update functions must not fail

Check for errors before committing

commit:

```
s <- snapshot
lazy check = s.get(x) > 10
update(x, v => if(check) f(v) else v)
update(y, v => if(check) g(v) else v)
update(z, v => if(check) h(v) else v)
```

Lazy Indexes & Bulk Operations

Lazy Indexes

Indexes are the basic building block of database management systems

- Access by primary key
- Secondary indexes

Goal: develop an index structure for lazy transactions

- Single-key and bulk operations
- Sequential access by writers to queue lazy operations
- Concurrent access by readers, which evaluate lazy operations
- Snapshots, for lazy reads

Immutable Data Structures

We use immutable data structures to provide snapshots

- Instances cannot be mutated: any instance is a snapshot
- Updates create a new version of the data structure
- Unmodified parts are shared (copy-on-write)



How can we atomically write many update functions without blocking?













Tries

Trie[K,V]

Empty Leaf(k : K, v : V) Branch(children : Array[Trie[K,V]])

Benefits

- Static balancing scheme
- Wide branches
- Ordered
- Snapshots: copy on write



Lazy Tries

Trie[K,V]
Empty
Leaf(k : K, v : Lazy[V])
Branch(
children : Array[Lazy[Trie[K,V]]])



Single key updates



Bulk range updates



Bulk random-access updates



Scheduling Lazy Operations

Immediate Evaluation

commit:

update(x, $v \Rightarrow ...$) update(y, $v \Rightarrow ...$) update(z, $v \Rightarrow ...$) s <- snapshot

result:

s.get(x) s.get(y) s.get(z)

Scheduling Opportunities

Delay evaluation of operations that are not needed immediately

- + Better temporal locality
- + Lower write latency
- + Can avoid work
- Higher read latency
- Higher memory consumption

Execute lazy operations in batches

- + Better spatial locality
- + Amortization of scheduling overhead
- Higher read latency

Experimental Results

Experimental Setup

We have implemented a lazy trie in Scala (JVM)

How does our system handle:

- Highly contended workloads
- Bulk operations: effect on concurrent transactions
- Realistic workloads: TPC-C

We compare against ScalaSTM

- Software transactional memory for Scala
- Transactional maps
- Optimistic concurrency control

Workload:

- update 2 hot records
- update 8 normal records





Random Access

Sequential Access



Lazy Reads

Optimistic Reads

Effect of bulk updates on concurrent OLTP workload







Conclusions

Conclusions

Lazy transactions execute contended transactions concurrently

- Non-blocking bulk updates
- Nearly ideal scaling on single-processor multi-core systems

Current concurrency control mechanisms cannot do this

• This is surprising, considering there has been 40 years of DB research

Requires core changes to database architecture

• Transactions are be submitted as programs

Future Work

- Scalability on NUMA systems: solving root contention
- Lazy relational schema transformations
- Lazy transactions for on-disk databases
- Language support for lazy transactions

Further Information

Read the paper

Lesley Wevers, Marieke Huisman and Maurice van Keulen. Lazy Evaluation for Concurrent OLTP and Bulk Transactions.

Implementation and benchmarks

https://github.com/utwente-fmt/lazy-transactions-IDEAS16