# **15-721**DATABASE SYSTEMS

#### Lecture #03 – Concurrency Control Part I

[Source

Andy Pavlo // Carnegie Mellon University // Spring 2016

#### TODAY'S AGENDA

Transaction Models Concurrency Control Overview Many-Core Evaluation

#### TRANSACTION DEFINITION

# TRANSACTION DEFINITION

A txn is a sequence of actions that are executed on a shared database to perform some higherlevel function.

Txns are the basic unit of change in the DBMS. No partial txns are allowed.



# ACTION CLASSIFICATION

#### **Unprotected Actions**

 $\rightarrow$  These lack all of the ACID properties except for consistency. Their effects cannot be depended upon.

#### **Protected Actions**

 $\rightarrow$  These do not externalize their results before they are completely done. Fully ACID.

#### **Real Actions**

 $\rightarrow$  These affect the physical world in a way that is hard or impossible to reverse.

## TRANSACTION MODELS

Flat Txns Flat Txns + Savepoints Chained Txns Nested Txns Saga Txns Compensating Txns

## FLAT TRANSACTIONS

Standard txn model that starts with **BEGIN**, followed by one or more actions, and then completed with either **COMMIT** or **ROLLBACK**.



# LIMITATIONS OF FLAT TRANSACTIONS

The application can only rollback the entire txn (i.e., no partial rollbacks).

All of a txn's work is lost is the DBMS fails before that txn finishes.

Each txn takes place at a single point in time



# LIMITATIONS OF FLAT TRANSACTIONS

#### Multi-Stage Planning

- $\rightarrow$  An application needs to make multiple reservations.
- $\rightarrow$  All the reservations need to occur or none of them.

#### Bulk Updates

- $\rightarrow$  An application needs to update one billion records.
- → This txn could take hours to complete and therefore the DBMS is exposed to losing all of its work for any failure or conflict.



Save the current state of processing for the txn and provide a handle for the application to refer to that savepoint.

The application can control the state of the txn through these checkpoints:

- $\rightarrow$  **ROLLBACK** Revert all changes back to the state of the DB at the savepoint.
- $\rightarrow$  **RELEASE** Destroys a savepoint previously defined in the txn.







CMU 15-721 (Spring 2016)

















Savepoint#2



































Savepoint#2




























Multiple txns executed one after another.

- Combined **COMMIT / BEGIN** operation is atomic.
- → No other txn can change the state of the database as seen by the second txn from the time that the first txn commits and the second txn begins.

Differences with savepoints:

- **COMMIT** allows the DBMS to free locks.
- Cannot rollback previous txns in chain.









Savepoints organize a transaction as a <u>sequence</u> of actions that can be rolled back individually.

- Nested txns form a **<u>hierarchy</u>** of work.
- $\rightarrow$  The outcome of a child txn depends on the outcome of its parent txn.



15





15









#### *Sub-Txn #1.1*























CARNEGIE MELLON DATABASE GROUP

CMU 15-721 (Spring 2016)

## BULK UPDATE PROBLEM

These other txn models are nice, but they still do not solve our bulk update problem.

Chained txns seems like the right idea but they require the application to handle failures and maintain it's own state.

 $\rightarrow$  Has to be able to reverse changes when things fail.

# COMPENSATING TRANSACTIONS

A special type of txn that is designed to semantically reverse the effects of another already committed txn.

Reversal has to be <u>logical</u> instead of physical.  $\rightarrow$  Example: Decrement a counter by one instead of reverting to the original value.



- A sequence of chained txns  $T_1, ..., T_n$  and compensating txns  $C_1, ..., C_{n-1}$  where one of the following is guaranteed:
- $\rightarrow$  The txns will commit in the order **T1**,...,**Tn**
- $\rightarrow$  The txns will commit in the order  $T_1, \dots, T_j, C_j, \dots, C_1$  (where j < n)









CARNEGIE MELLON DATABASE GROUP








# CONCURRENCY CONTROL

The protocol to allow txns to access a database in a multi-programmed fashion while preserving the illusion that each of them is executing alone on a dedicated system.
→ The goal is to have the effect of a group of txns on the database's state is equivalent to any serial execution of all txns.

#### Provides <u>A</u>tomicity + <u>I</u>solation in ACID

# TXN INTERNAL STATE

#### **Undo Log Entries**

- $\rightarrow$  Stored in an in-memory data structure.
- $\rightarrow$  Dropped on commit.

## **Redo Log Entries**

- $\rightarrow$  Append to the in-memory tail of WAL.
- $\rightarrow$  Flushed to disk on commit.

#### **Read/Write Set**

 $\rightarrow$  Depends on the concurrency control scheme.

# CONCURRENCY CONTROL SCHEMES

#### Two-Phase Locking (2PL)

 $\rightarrow$  Assume txns will conflict so they must acquire locks on elements before they are allowed to access them.

#### Timestamp Ordering (T/O)

→ Assume that conflicts are rare so txns do not need to acquire locks and instead check for conflicts at commit time.























#### **Txn #1**







#### **Txn #2**



















#### **Deadlock Detection**

- $\rightarrow$  Each txn maintains a queue of the txns that hold the locks that it waiting for.
- $\rightarrow$  A separate thread checks these queues for deadlocks.
- → If deadlock found, use a heuristic to decide what txn to kill in order to break deadlock.

#### **Deadlock Prevention**

- $\rightarrow$  Check whether another txn already holds a lock when another txn requests it.
- $\rightarrow$  If lock is not available, the txn will either (1) wait, (2) commit suicide, or (3) kill the other txn.









Record	Read Timestamp	Write Timestamp	
А	10000	10000	
В	10000	10000	

CMU 15-721 (Spring 2016)

.







#### 25







Record	Read Timestamp	Write Timestamp		
А	10001	10000	] 🖌	
В	10000	10001		~

25





Record	Read Timestamp	Write Timestamp	
А	10001	10000	
В	10000	10001	

CMU 15-721 (Spring 2016)

•

CARNEGIE MELLON DATABASE GROUP



Record	Read Timestamp	Write Timestamp	
А	10001	10005	
В	10000	10001	

CMU 15-721 (Spring 2016)

•






## TIMESTAMP ORDERING

#### Basic T/O

- $\rightarrow$  Check for conflicts on each read/write.
- $\rightarrow$  Copy tuples on each access to ensure repeatable reads.

#### Multi-Version Concurrency Control (MVCC)

- $\rightarrow$  Create a new version of a tuple whenever a txn modifies it. Use timestamps as version id.
- $\rightarrow$  Check visibility on every read/write.

#### **Optimistic Currency Control (OCC)**

- $\rightarrow$  Store all changes in private workspace.
- $\rightarrow$  Check for conflicts at commit time and then merge.

DL\_DETECT NO\_WAIT WAIT\_DIE 2PL w/ Deadlock Detection 2PL w/ Non-waiting Prevention 2PL w/ Wait-and-Die Prevention

TIMESTAMP MVCC OCC Basic T/O Algorithm Multi-Version T/O Optimistic Concurrency Control

DL\_DETECT NO\_WAIT WAIT\_DIE 2PL w/ Deadlock Detection 2PL w/ Non-waiting Prevention 2PL w/ Wait-and-Die Prevention

TIMESTAMP MVCC OCC Basic T/O Algorithm Multi-Version T/O Optimistic Concurrency Control

DL\_DETECT NO\_WAIT WAIT\_DIE 2PL w/ Deadlock Detection 2PL w/ Non-waiting Prevention 2PL w/ Wait-and-Die Prevention

TIMESTAMP MVCC OCC Basic T/O Algorithm Multi-Version T/O Optimistic Concurrency Control



CMU 15-721 (Spring 2016)

DL\_DETECT NO\_WAIT WAIT\_DIE 2PL w/ Deadlock Detection 2PL w/ Non-waiting Prevention 2PL w/ Wait-and-Die Prevention

TIMESTAMP MVCC OCC Basic T/O Algorithm Multi-Version T/O Optimistic Concurrency Control

DL\_DETECT NO\_WAIT WAIT\_DIE

2PL w/ Deadlock Detection 2PL w/ Non-waiting Prevention 2PL w/ Wait-and-Die Prevention

TIMESTAMP MVCC OCC Basic T/O Algorithm Multi-Version T/O Optimistic Concurrency Control

DL\_DETECT NO\_WAIT WAIT\_DIE 2PL w/ Deadlock Detection 2PL w/ Non-waiting Prevention 2PL w/ Wait-and-Die Prevention

TIMESTAMP MVCC OCC Basic T/O Algorithm Multi-Version T/O

**Optimistic Concurrency Control** 



CMU 15-721 (Spring 2016)

DL\_DETECT NO\_WAIT WAIT\_DIE 2PL w/ Deadlock Detection 2PL w/ Non-waiting Prevention 2PL w/ Wait-and-Die Prevention

TIMESTAMP MVCC OCC Basic T/O Algorithm Multi-Version T/O Optimistic Concurrency Control

# 1000-CORE CPU SIMULATOR

#### DBx1000 Database System

- $\rightarrow$  In-memory DBMS with pluggable lock manager.
- $\rightarrow$  No network access, logging, or concurrent indexes

#### MIT Graphite CPU Simulator

- $\rightarrow$  Single-socket, tile-based CPU.
- $\rightarrow$  Shared L2 cache for groups of cores.
- $\rightarrow$  Tiles communicate over 2D-mesh network.





28

# TARGET WORKLOAD

Yahoo! Cloud Serving Benchmark (YCSB)

- $\rightarrow$  20 million tuples
- $\rightarrow$  Each tuple is 1KB (total database is ~20GB)
- Each transactions reads/modifies 16 tuples.
- Varying skew in transaction access patterns. Serializable isolation level.

































## BOTTLENECKS

# **Lock Thrashing** $\rightarrow$ DL\_DETECT, WAIT\_DIE

# **Timestamp Allocation** $\rightarrow$ All T/O algorithms + WAIT\_DIE

 $\begin{array}{l} \textbf{Memory Allocations} \\ \rightarrow \text{ OCC + MVCC} \end{array}$ 



## LOCK THRASHING

Each txn waits longer to acquire locks, causing other txn to wait longer to acquire locks.

Can measure this phenomenon by removing deadlock detection/prevention overhead.  $\rightarrow$  Force txns to acquire locks in primary key order.

 $\rightarrow$  Deadlocks are not possible.

#### LOCK THRASHING



#### LOCK THRASHING



#### [K THRASHING



By reducing the frequency of lock conversion deadlocks, we have dispensed with deadlock as a major performance consideration, so we are left with blocking situations. Blocking affects performance in a rather dramatic

> Number of Active Transactions

the workload decreases the throughput.

Throughput High

Low

and few transactions can make progress.

SQL Server. SQL Server also allows update locks to be obtained in a SELECT (i.e., read) statement, but in this case, it will not downgrade the update locks to read locks, since it doesn't know when it is safe to do so.

converts the update lock to a write lock. This lock conversion can't lead to a lock conversion deadlock, because at most one transaction can have an update lock on the data item. (Two transactions must try to convert the lock at the same time to create a lock conversion deadlock.) On the other hand, the benefit of this approach is that an update lock does not block other transactions that read without expecting to update later on. The weakness is that the request to convert the update lock to a write lock may be delayed by other read locks. If a large number of data items are read and only a few of them are updated, the tradeoff is worthwhile. This approach is used in Microsoft

Thrashing Region

High

6.4 Performance

155

the workload decreases the throughput. Throughput High Thrashing Region Number of Active Transactions High 107 Lock Thrashing. When the number of active transactions gets too high, many transactions suddenly become blocked, and few transactions can make progress.

which is measured by the number of active transactions. When the system is idle, the first transaction to run cannot block due to locks, because it's the only one requesting locks. As the number of active transactions grows, each successive transaction has a higher probability of becoming blocked due to transactions already running. When the number of active transactions is high enough, the next transaction to be started has virtually no chance of running to completion without blocking for some lock. Worse, it probably will get some locks before encountering one that blocks it, and these locks contribute to the likelihood that other active transactions will become blocked. So, not only does it not contribute to increased throughput, but by getting some locks that block other transactions, it actually reduces throughput. This leads to thrashing, where increasing

One way to understand lock thrashing is to consider the effect of slowly increasing the transaction load,

By reducing the frequency of lock conversion deadlocks, we have dispensed with deadlock as a major performance consideration, so we are left with blocking situations. Blocking affects performance in a rather dramatic way. Until lock usage reaches a saturation point, it introduces only modest delays-significant, but not a serious problem. At some point, when too many transactions request locks, a large number of transactions suddealy become blocked, and few transactions can make progress. Thus, transaction throughput stops growing. Surprisingly, if enough transactions are initiated, throughput actually decreases. This is called lock thrashing (see Figure 6.7). The main issue in locking performance is to maximize throughput without reaching the point

case, it will not downgrade the update locks to read locks, since it doesn't know when it is safe to do so.

converts the update lock to a write lock. This lock conversion can't lead to a lock conversion deadlock, because at most one transaction can have an update lock on the data item. (Two transactions must try to convert the lock at the same time to create a lock conversion deadlock.) On the other hand, the benefit of this approach is that an update lock does not block other transactions that read without expecting to update later on. The weakness is that the request to convert the update lock to a write lock may be delayed by other read locks. If a large number of data items are read and only a few of them are updated, the tradeoff is worthwhile. This approach is used in Microsoft SQL Server. SQL Server also allows update locks to be obtained in a SELECT (i.e., read) statement, but in this

155

6.4 Performance





## TIMESTAMP ALLOCATION

#### Mutex

 $\rightarrow$  Worst option.

#### **Atomic Addition**

 $\rightarrow$  Requires cache invalidation on write.

#### **Batched Atomic Addition**

 $\rightarrow$  Needs a back-off mechanism to prevent fast burn.

#### Hardware Clock

 $\rightarrow$  Not sure if it will exist in future CPUs.

#### Hardware Counter

 $\rightarrow$  Not implemented in existing CPUs.

#### TIMESTAMP ALLOCATION



## MEMORY ALLOCATIONS

Copying data on every read/write access slows down the DBMS because of contention on the memory controller.

 $\rightarrow$  In-place updates and non-copying reads are not affected as much.

Default libc malloc is slow. Never use it.

# PARTITION-LEVEL LOCKING

The database is split up into horizontal partitions:

- $\rightarrow$  Each partition is assigned a single-threaded execution engine that has exclusive access to its data.
- $\rightarrow$  In-place updates.

Only one txn can execute at a time per partition.

- $\rightarrow$  Order txns based on when they arrive at the DBMS.
- $\rightarrow$  A txn acquires the lock for a partition when it has the lowest timestamp.
- $\rightarrow$  It is not allowed to access any partition that it does not hold the lock for.
## READ-ONLY WORKLOAD



## MULTI-PARTITION WORKLOADS



## PARTING THOUGHTS

Concurrency control is hard to get correct and perform well.

Evaluation did not consider HTAP workloads.



CMU 15-721 (Spring 2016)

## NEXT CLASS

Isolation Levels Modern MVCC



43

CMU 15-721 (Spring 2016)