Carnegie Mellon University

ADVANCED DATABASE SYSTEMS

Checkpoint Protocols @Andy_Pavlo // 15-721 // Spring 2018

 \bigcirc

#

Cture

 $\overline{\mathbb{O}}$

COURSE ANNOUNCEMENTS

Mid-Term: Wednesday March 7th @ 3:00pm

Project #2: Monday March 12th @ 11:59pm

Project #3 Proposal: Monday March 19th



TODAY'S AGENDA

In-Memory Checkpoints Shared Memory Restarts





OBSERVATION

Logging allows the DBMS to recover the database after a crash/restart. But this system will have to replay the entire log each time.

Checkpoints allows the systems to ignore large segments of the log to reduce recovery time.



IN-MEMORY CHECKPOINTS

There are different approaches for how the DBMS can create a new checkpoint for an in-memory database.

The choice of approach in a DBMS is tightly coupled with its concurrency control scheme.

The checkpoint thread(s) scans each table and writes out data asynchronously to disk.



IDEAL CHECKPOINT PROPERTIES

Do <u>not</u> slow down regular txn processing.

Do **<u>not</u>** introduce unacceptable latency spikes.

Do **not** require excessive memory overhead.



LOW-OVERHEAD ASYNCHRONOUS CHECKPOINTING IN MAIN-MEMORY DATABASE SYSTEMS SIGMOD 2016



CONSISTENT VS. FUZZY CHECKPOINTS

Approach #1: Consistent Checkpoints

- \rightarrow Represents a consistent snapshot of the database at some point in time. No uncommitted changes.
- \rightarrow No additional processing during recovery.

Approach #2: Fuzzy Checkpoints

- \rightarrow The snapshot could contain records updated from transactions that have not finished yet.
- \rightarrow Must do additional processing to remove those changes.



CHECKPOINT CONTENTS

Approach #1: Complete Checkpoint

→ Write out every tuple in every table regardless of whether were modified since the last checkpoint.

Approach #2: Delta Checkpoint

- → Write out only the tuples that were modified since the last checkpoint.
- \rightarrow Can merge checkpoints together in the background.



FREQUENCY

Taking checkpoints too often causes the runtime performance to degrade.

But waiting a long time between checkpoints is just as bad.

Approach #1: Time-based Approach #2: Log File Size Threshold Approach #3: On Shutdown (always!)



CHECKPOINT IMPLEMENTATIONS

	Туре	Contents	Frequency
MemSQL	Consistent	Complete	Log Size
VoltDB	Consistent	Complete	Time-Based
Altibase	Fuzzy	Complete	Manual?
TimesTen	Consistent (Blocking) Fuzzy (Non-Blocking)	Complete Complete	On Shutdown Time-Based
Hekaton	Consistent	Delta	Log Size
SAP HANA	Fuzzy	Complete	Time-Based



IN-MEMORY CHECKPOINTS

Approach #1: Naïve Snapshots

Approach #2: Copy-on-Update Snapshots

Approach #3: Wait-Free ZigZag

Approach #4: Wait-Free PingPong





NAÏVE SNAPSHOT

Create a consistent copy of the entire database in a new location in memory and then write the contents to disk.

Two approaches to copying database:

- \rightarrow Do it yourself (tuple data only).
- \rightarrow Let the OS do it for you (everything).



HYPER - FORK SNAPSHOTS

Create a snapshot of the database by forking the DBMS process.

- \rightarrow Child process contains a consistent checkpoint if there are not active txns.
- → Otherwise, use the in-memory undo log to roll back txns in the child process.

Continue processing txns in the parent process.





H-STORE - FORK SNAPSHOTS

Workload: TPC-C (8 Warehouses) + OLAP Query





14

COPY-ON-UPDATE SNAPSHOT

During the checkpoint, txns create new copies of data instead of overwriting it. \rightarrow Copies can be at different granularities (block, tuple)

The checkpoint thread then skips anything that was created after it started.





VOLTDB - CONSISTENT CHECKPOINTS

A special txn starts a checkpoint and switches the DBMS into copy-on-write mode.

- \rightarrow Changes are no longer made in-place to tables.
- → The DBMS tracks whether a tuple has been inserted, deleted, or modified since the checkpoint started.

A separate thread scans the tables and writes tuples out to the snapshot on disk.

- \rightarrow Ignore anything changed after checkpoint.
- \rightarrow Clean up old versions as it goes along.



OBSERVATION

Txns have to wait for the checkpoint thread when using naïve snapshots.

Txns may have to wait to acquire latches held by the checkpoint thread under copy-on-update if not using MVCC.



WAIT-FREE ZIGZAG

Maintain two copies of the entire database \rightarrow Each txn write only updates one copy.

Use two BitMaps to keep track of what copy a txn should read/write from per tuple.

 \rightarrow Avoid the overhead of having to create copies on the fly as in the copy-on-update approach.



WAIT-FREE ZIGZAG





Trade extra memory + CPU to avoid pauses at the end of the checkpoint.

- Maintain two copies of the entire database at all times plus a third "base" copy.
- \rightarrow Pointer indicates which copy is the current master.
- \rightarrow At the end of the checkpoint, swap these pointers.



-

Base Copy



Co	ру)#1
	0	-
	0	_

0

0

0

0

Copy #2



Master: Copy #1 Shadow: Copy #2

CARNEGIE MELLON DATABASE GROUP





Copy #1





Master: Copy #1



CARNEGIE MELLON DATABASE GROUP



Copy #2





Master: Copy #1

CARNEGIE MELLON DATABASE GROUP



Copy #1

Copy #2



Base Copy



Copy #1					
	1	6			
	0	-			
	1	1			
	1	9			
	0	-			
	0	-			

Copy #2



Master:Copy #2Shadow:Copy #1

CARNEGIE MELLON DATABASE GROUP







CHECKPOINT IMPLEMENTATIONS

Bulk State Copying

 \rightarrow Pause txn execution to take a snapshot.

Locking / Latching

 \rightarrow Use latches to isolate the checkpoint thread from the worker threads if they operate on shared regions.

Bulk Bit-Map Reset:

 \rightarrow If DBMS uses BitMap to track dirty regions, it must perform a bulk reset at the start of a new checkpoint.

Memory Usage:

 \rightarrow To avoid synchronous writes, the method may need to allocate additional memory for data copies.



IN-MEMORY CHECKPOINTS

	Bulk		Bulk Bit-	Memory
	Copying	Locking	Map Reset	Usage
Naïve Snapshot	Yes	No	No	2x
Copy-on-Update	No	Yes	Yes	2x
Wait-Free ZigZag	No	No	Yes	2x
Wait-Free Ping-Pong	No	No	No	3x



OBSERVATION

Not all DBMS restarts are due to crashes.

- \rightarrow Updating OS libraries
- \rightarrow Hardware upgrades/fixes

 \rightarrow Updating DBMS software

Need a way to be able to quickly restart the DBMS without having to re-read the entire database from disk again.



FACEBOOK SCUBA - FAST RESTARTS

Decouple the in-memory database lifetime from the process lifetime.

By storing the database shared memory, the DBMS process can restart and the memory contents will survive.





CMU 15-721 (Spring 2018)

33

FACEBOOK SCUBA

Distributed, in-memory DBMS for time-series event analysis and anomaly detection.

Heterogeneous architecture

- → Leaf Nodes: Execute scans/filters on in-memory data
- \rightarrow **Aggregator Nodes:** Combine results from leaf nodes



FACEBOOK SCUBA - ARCHITECTURE



SHARED MEMORY RESTARTS

Approach #1: Shared Memory Heaps

- \rightarrow All data is allocated in SM during normal operations.
- \rightarrow Have to use a custom allocator to subdivide memory segments for thread safety and scalability.
- \rightarrow Cannot use lazy allocation of backing pages with SM.

Approach #2: Copy on Shutdown

- \rightarrow All data is allocated in local memory during normal operations.
- \rightarrow On shutdown, copy data from heap to SM.



FACEBOOK SCUBA - FAST RESTARTS

When the admin initiates restart command, the node halts ingesting updates.

DBMS starts copying data from heap memory to shared memory.

 \rightarrow Delete blocks in heap once they are in SM.

Once snapshot finishes, the DBMS restarts.

- \rightarrow On start up, check to see whether the there is a valid database in SM to copy into its heap.
- \rightarrow Otherwise, the DBMS restarts from disk.



PARTING THOUGHTS

I think that copy-on-update checkpoints are the way to go especially if you are using MVCC

Shared memory does have some use after all...



NEXT CLASS

Networking Protocols

