

Lecture #20: Logging Schemes

15-445/645 Database Systems (Fall 2018)

<https://15445.courses.cs.cmu.edu/fall2018/>

Carnegie Mellon University

Prof. Andy Pavlo

1 Crash Recovery

Recovery algorithms are techniques to ensure database consistency, transaction atomicity, and durability despite failures. DBMS is divided into different components based on the underlying storage device. We must also classify the different types of failures that the DBMS needs to handle.

Every recovery algorithm has two parts:

- Actions during normal transaction processing to ensure that the DBMS can recover from a failure.
- Actions after a failure to recover the database to a state that ensures atomicity, consistency, and durability.

The key primitives that we are going to use in a recovery algorithm are UNDO and REDO. Not all algorithms use both of these:

- **UNDO**: The process of removing the effects of an incomplete or aborted transaction.
- **REDO**: The process of re-instating the effects of a committed transaction for durability.

2 Storage Types

- **Volatile Storage**
 - Data does not persist after power is cut.
 - Examples: DRAM, SRAM,.
- **Non-Volatile Storage**
 - Data persists after losing power.
 - Examples: HDD, SDD.
- **Stable Storage**
 - A non-existent form of non-volatile storage that survives all possible failures scenarios.
 - Use multiple storage devices to approximate.

3 Failure Classification

- **Type #1: Transaction Failures**
 - **Logical Errors**: A transaction cannot complete due to some internal error condition (e.g., integrity, constraint violation).
 - **Internal State Errors**: The DBMS must terminate an active transaction due to an error condition (e.g. deadlock)
- **Type #2: System Failures**
 - **Software Failure**: There is a problem with the DBMS implementation (e.g. uncaught divide-by-zero exception) and the system has to halt.

- **Hardware Failure:** The computer hosting the DBMS crashes. We assume that non-volatile storage contents are not corrupted by system crash.
- **Type #3: Storage Media Failure**
 - **Non-Repairable Hardware Failure:** A head crash or similar disk failure destroys all or parts of non-volatile storage. Destruction is assumed to be detectable. No DBMS can recover from this. Database must be restored from archived version

4 Buffer Pool Management Policies

Steal Policy: Whether the DBMS allows an uncommitted transaction to overwrite the most recent committed value of an object in non-volatile storage (can a transaction write uncommitted changes to disk).

- **STEAL:** is allowed
- **NO-STEAL:** is not allowed.

Force Policy: Whether the dbms ensures that all updates made by a transaction are reflected on non-volatile storage before the transaction is allowed to commit

- **FORCE:** Is enforced
- **NO-FORCE:** Is not enforced

Force writes makes it easier to recover but results in poor runtime performance.

Easiest System to implement: NO-STEAL + FORCE

- The DBMS never has to undo changes of an aborted transaction because the changes were not written to disk.
- It also never has to redo changes of a committed transaction because all the changes are guaranteed to be written to disk at committed.
- **Limitation:** If all of the data that a transaction needs to modify does not fit on memory, then that transaction cannot execute because the DBMS is not allowed to write out dirty pages to disk before the transaction commits.

5 Shadow Paging

The DBMS maintains two separate copies of the database (**master, shadow**). Updates are only made in the shadow copy. When a transaction commits, atomically switch the shadow to become the new master. This is an example of a NO-STEAL + FORCE system.

Disadvantages: Copying the entire page table is expensive and the commit overhead is high.

Implementation:

- Organize the database pages in a tree structure where the root is a single disk page.
- There are two copies of the tree, the master and the shadow:
 - The root points to the master copy
 - Updates are applied to the shadow copy
- To install updates, overwrite the root so it points to the shadow, thereby swapping the master and shadow.
 - Before overwriting the root, none of the transactions updates are part of the disk-resident database.
 - After overwriting the root, all of the transactions updates are part of the disk resident database.
- **UNDO:** Remove the shadow pages. Leave master and the DB root pointer alone
- **REDO:** Not needed at all

6 Write-Ahead Logging

The DBMS records all the changes made to the database in a log file (on stable storage) before the change is made to a disk page. The log contains sufficient information to perform the necessary undo and redo actions to restore the database after a crash. This is an example of a STEAL + NO-FORCE system.

Almost every DBMS uses write-ahead logging (WAL) because it has the fastest runtime performance. But its recovery time is slower because it has to replay the log.

Implementation:

- All log records pertaining to an updated page are written to non-volatile storage before the page itself is allowed to be overwritten in non-volatile storage.
- A transaction is not considered committed until all its log records have been written to stable storage.
- When the transaction starts, write a <BEGIN> record to the log for each transaction to mark its starting point.
- When a transaction finishes, write a <COMMIT> record to the log and make sure all log records are flushed before it returns an acknowledgement to the application.
- Each log entry contains information about the change to a single object:
 - Transaction ID.
 - Object ID.
 - Before Value (used for UNDO).
 - After Value (used for REDO).
- Log entries to disk should be done when transaction commits. You can use group commit to batch multiple log flushes together to amortize overhead.
- **Deferred Updates**
 - If we prevent the DBMS from writing dirty records to disk until the transaction commits, then we do not need to store their original values.
 - This will not work if the change set of a transaction is larger than the amount of memory available.
 - The DBMS cannot undo changes for an aborted transaction if it does not have the original values in the log.
 - Thus, this is why the DBMS needs to use the **STEAL** policy.

7 Checkpoints

The main problem with write-ahead logging is that the log file will grow forever. After a crash, the DBMS has to replay the entire log, which can take a long time if the log file is large. Thus, the DBMS can periodically take a *checkpoint* where it flushes all buffers out to disk.

It is not obvious how often the DBMS should take a checkpoint. Checkpointing too often causes the runtime performance to degrade. But waiting a long time is just as bad, as recovery time increases.

Blocking Checkpoint Implementation:

- The DBMS stops accepting new transactions and waits for all active transactions to complete.
- Flush all log records and dirty blocks currently residing in main memory to stable storage.
- Write a <CHECKPOINT> entry to the log and flush to stable storage.

8 Logging Schemes

- **Physical Logging**

- Record the changes made to a specific location in the database
- Example: Position of a record in a page
- **Logical Logging**
 - Record the high level operations executed by transactions. Not necessarily restricted to single page. Requires less data written in each log record than physical logging. Difficult to implement recovery with logical logging if you have concurrent transactions in a non-deterministic concurrency control scheme.
 - Example: The UPDATE, DELETE, and INSERT queries invoked by a transaction.
- **Physiological logging**
 - Hybrid approach where log records target a single page but do not specify data organization of the page.
 - Most commonly used approach.

9 Conclusion

- Write-Ahead Logging is most commonly used approach to handle loss of volatile storage
- Use incremental updates (STEAL + NO-FORCE) with checkpoints
- On recovery: undo uncommitted transactions + redo committed transactions