

Lecture 22: Database Recovery

15-445/645 Database Systems (Fall 2017)

Carnegie Mellon University

Prof. Andy Pavlo

ARIES

Algorithms for Recovery and Isolation Exploiting Semantics. Developed at IBM research in early 1990s. Not all systems implement ARIES exactly as defined in the original paper, but they're pretty close.

Main ideas:

1. **Write ahead logging:** Any change is recorded in log on stable storage before the database change is written to disk (STEAL + NO-FORCE).
2. **Repeating history during redo:** On restart, retrace actions and restore database to exact state before crash.
3. **Logging changes during undo:** Record undo actions to log to ensure action is not repeated in the event of repeated failures.

WAL Records

We need to extend log record format to include additional info. Every log record now includes a globally unique **log sequence number** (LSN). Various components in the system keep track of **LSNs** that pertain to them:

1. Each data page contains a **pageLSN**: The LSN of the most recent update to that page.
2. System keeps track of **flushedLSN**: The max **LSN** flushed so far
3. Before page i can be written to disk, we must flush log at least to the point where $pageLSN_i \leq flushedLSN$

Normal Execution

1. **Transaction Commit**
 - (a) Write **COMMIT** record to log.
 - (b) All log records to to transaction's **COMMIT** record are flushed to disk. Note that log flushes are sequential, synchronous writes to disk. There can also be multiple log records per log page.
 - (c) When the commit succeeds, write a special **TXN-END** record to log. This is used for internal bookkeeping and doesn't need to be flushed immediately.

2. Transaction Abort

- (a) Aborting a transaction is actually a special case of the ARIES undo operation applied to only one transaction.
- (b) We need to add another field to our log records:
 - i. **prevLSN**: The previous LSN for the transaction.
 - ii. This maintains a linked-list for each transaction that makes it easy to walk through its records.
- (c) **Compensation Log Record**
 - i. A **CLR** describes the actions taken to undo the actions of a previous update record. It has all the fields of an update log record plus the undoNext pointer (i.e., the next-to-be-undone LSN).
 - ii. **CLRs** are added to the log like any other record but they never need to be undone.
- (d) Algorithm:
 - i. First write **ABORT** record to log.
 - ii. Then play back updates in reverse order to remove their effects. For each update, write a **CLR** entry and restore old value.
 - iii. At end, write a **TXN-END** log record.

Checkpointing

Blocking Checkpoints

The DBMS halts everything when it takes a checkpoint to ensure that it writes a consistent snapshot of the database to disk. This is the same approach discussed in previous lecture:

1. Halt the start of any new transactions.
2. Wait until all active transactions finish executing.
3. Flush dirty pages on disk.

Slightly Better Blocking Checkpoints

Like previous checkpoint scheme except that you the DBMS does not have to wait for active transactions to finish executing. We have to now record internal system state as of the beginning of the checkpoint.

1. Halt the start of any new transactions.
2. Pause transactions while the DBMS takes the checkpoint.
3. Internal State Data Structures:

(a) Active Transaction Table (ATT)

- i. One entry per active transaction
- ii. transactionId: Unique transaction identifier
- iii. status: the current “mode” of the transaction (**R**unning, **C**ommitting, **U**ndo candidate)

- iv. lastLSN: Most recent LSN written by transaction
 - v. Entry is removed when transaction commits or aborts
- (b) **Dirty Page Table (DPT)**
- i. Keep track of pages in the buffer pool contain changes from uncommitted transactions
 - ii. One entry per dirty page containing the recLSN, the LSN of the log record that first caused the page to be dirty.

Fuzzy Checkpoints

A **fuzzy checkpoint** is where the DBMS allows other transactions to continue to run. This is what ARIES uses in its protocol.

1. Add new log records to track checkpoint boundaries
 - (a) **CHECKPOINT-BEGIN**: Indicates the start of the checkpoint
 - (b) **CHECKPOINT-END**: Contains the ATT + DPT

ARIES Recovery

The ARIES protocol is comprised of three phases. Upon start-up after a crash, the DBMS will execute the following phases:

Analysis: Read the WAL to identify dirty pages in the buffer pool and active transactions at the time of the crash.

Redo: Repeat all actions starting from an appropriate point in the log.

Undo: Reverse the actions of transactions that did not commit before the crash.

Analysis Phase

Start from last checkpoint found via MasterRecord

1. Scan log forward from the checkpoint.
2. If you find a **TXN-END** record, remove its transaction from ATT.
3. All other records, add transaction to ATT with status **UNDO**, and on commit, change transaction status to **COMMIT**.
4. For UPDATE records, if page P not in **DPT**, add P to **DPT** and set its recLSN=LSN

Redo Phase

The goal is to repeat history to reconstruct state at the moment of the crash. Reapply all updates (even aborted transactions) and redo **CLRs**:

1. Scan forward from log record containing smallest recLSN in **PDT**.
2. For each update log record or CLR with a given LSn, redo the action unless:
 - (a) Affected page is not in the **DPT**, or

- (b) Affected page is in **DPT** but that record's LSN is greater than smallest recLSN, or
 - (c) Affected pageLSN (on disk) \geq LSN.
3. To redo an action:
- (a) Reapply logged action.
 - (b) Set pageLSN to log records **LSN**.
 - (c) No additional logging, no forcing.
 - (d) At the end of the redo phase, write **TXN-END** log records for all transactions with status "C" and remove them from the **ATT**.

Undo Phase

1. Undo All transactions active at the time of crash
2. These are all transactions with "U" status in the **ATT** after the Analysis phase
3. Process them in reverse **LSN** order using the lastLSN to speed up traversal
4. Write a CLR for every modification

Conclusion

1. Main ideas of ARIES
 - (a) WAL, STEAL + NO-FORCE.
 - (b) Fuzzy checkpoints (snapshot of dirty page ids).
 - (c) Redo everything since the earliest dirty page; undo "loser" transactions.
 - (d) Write CLR's when undoing to survive failures during restarts.
2. Log Sequence Numbers (LSN)
 - (a) **LSNs** identify log records, linked into backwards chains per transactions via prevLSN.
 - (b) pageLSN allows comparison of data page and log records.