

Lecture 9: ARIES from First Principles

CREATING THE NEXT®

Recap

Definitions

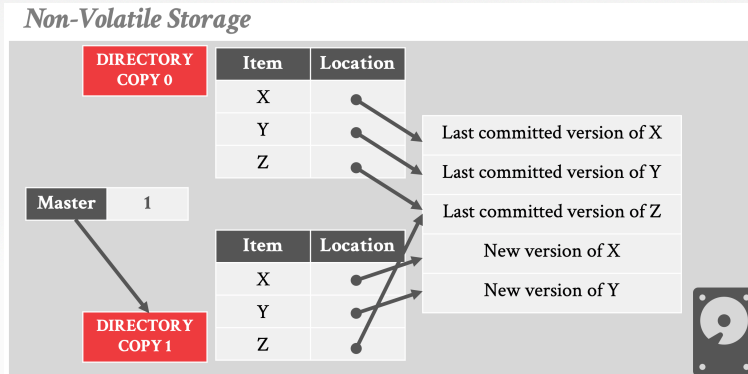
Policy vs Mechanism

- Policy
 - ▶ Specifies the desired behavior of the system (**what**).
 - ▶ Example: Buffer manager may adopt the LRU policy for evicting pages from the buffer.
- Mechanism
 - ▶ Specifies how that behavior must be realized (**how**)
 - ▶ Example: We may implement the policy using: (1) uni-directional map + linked list, or (2) bi-directional map. Optimize the code for specific hardware technology.

Deriving ARIES

- DRAM is volatile

V1: SHADOW PAGING



- DRAM is volatile
- Avoid random writes to database on disk (NO FORCE)

Constraints

- DRAM is volatile
- Avoid random writes to database on disk (NO FORCE)
- Support transactions with change sets > DRAM (STEAL)

V3: WAL

- **Phase 1 – Analysis**

- ▶ Read the WAL to identify dirty pages in the buffer pool and active txns at the time of the crash.

- **Phase 2 – Redo**

- ▶ Repeat all actions starting from an appropriate point in the log.

- **Phase 3 – Undo**

- ▶ Reverse the actions of txns that did not commit before the crash.

V3: WAL

<u>LSN Type</u>	<u>Where</u>	<u>Definition</u>
flushedLSN	Memory	Last LSN in log on disk
pageLSN	$page_x$	Newest update to $page_x$
prevLSN	log record	LSN of prior log record by same txn
recLSN	DPT	Oldest update to $page_x$ since it was last flushed
lastLSN	ATT	Latest action of txn T_i

V3: WAL

- RecLSN (in memory – Dirty Page Table)
 - ▶ Determine whether page state has not made it to disk.
 - ▶ If there is a suspicion, then page has to be accessed.
 - ▶ Serves to limit the number of pages whose PageLSN has to be examined
 - ▶ If a file sync operation is found in the log, all the pages in the file are removed from the dirty page table
- LastLSN (in memory – Active Transaction Table)
 - ▶ Determine log records which have to be rolled back for the yet-to-be-completely-undone uncommitted transactions

V3: WAL

- Advantages
 - ▶ Maximum flexibility for buffer manager
- Disadvantages
 - ▶ Log will keep growing over time thereby slowing down recovery and taking up more storage space.

Constraints

- DRAM is volatile
- Avoid random writes to database on disk (NO FORCE)
- Support transactions with change sets > DRAM (STEAL)
- Recovery time must be bounded.

V4: COMMIT-CONSISTENT CHECKPOINTS

<u>LSN Type</u>	<u>Where</u>	<u>Definition</u>
flushedLSN	Memory	Last LSN in log on disk
pageLSN	$page_x$	Newest update to $page_x$
prevLSN	log record	LSN of prior log record by same txn
recLSN	DPT	Oldest update to $page_x$ since it was last flushed
lastLSN	ATT	Latest action of txn T_i
MasterRecord	Disk	LSN of latest checkpoint

V4: COMMIT-CONSISTENT CHECKPOINTS

- **Phase 1 – Analysis**
 - ▶ Read the WAL starting from the latest checkpoint.
- **Phase 2 – Redo**
 - ▶ Repeat all actions starting from an appropriate point in the log.
- **Phase 3 – Undo**
 - ▶ Reverse the actions of txns that did not commit before the crash.

V4: COMMIT-CONSISTENT CHECKPOINTS

- Advantages
 - ▶ Recovery time is bounded due to checkpoints.
- Disadvantages
 - ▶ With commit consistent checkpointing, DBMS must stop processing transactions while taking checkpoint
 - ▶ Users will suffer long delays due to checkpointing

Constraints

- DRAM is volatile
- Avoid random writes to database on disk (NO FORCE)
- Support transactions with change sets > DRAM (STEAL)
- Recovery time must be bounded.
- Users must not suffer long delays due to checkpointing.

V5: FUZZY CHECKPOINTS

- Instead of flushing **all** dirty pages, only flush those dirty pages that have not been flushed since before the **previous checkpoint**.
- This guarantees that, at any time, all updates of committed transactions that occurred before the **penultimate** (*i.e.*, second to last) checkpoint have been applied to database on disk - during the last checkpoint, if not earlier.

V5: FUZZY CHECKPOINTS

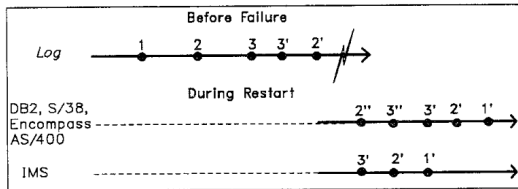
- Advantages
 - ▶ With fuzzy checkpointing, DBMS can concurrently process transactions while taking checkpoints.
- Problem
 - ▶ Repeated failures during recovery can lead to unbounded amount of logging during recovery

Constraints

- DRAM is volatile
- Avoid random writes to database on disk (NO FORCE)
- Support transactions with change sets > DRAM (STEAL)
- Recovery time must be bounded.
- Users must not suffer long delays due to checkpointing.
- Cope with failures during recovery.

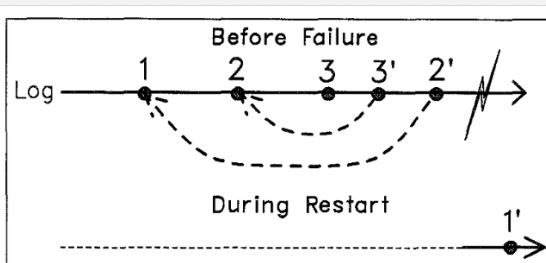
V6: COMPENSATION LOG RECORDS

- Problems: (1) compensating compensations and (2) duplicate compensations



I' is the CLR for I and I'' is the CLR for I'

V6: COMPENSATION LOG RECORDS



I' is the Compensation Log Record for I
 I' points to the predecessor, if any, of I

V6: COMPENSATION LOG RECORDS

<u>LSN Type</u>	<u>Where</u>	<u>Definition</u>
flushedLSN	Memory	Last LSN in log on disk
pageLSN	$page_x$	Newest update to $page_x$
prevLSN	log record	LSN of prior log record by same txn
recLSN	DPT	Oldest update to $page_x$ since it was last flushed
lastLSN	ATT	Latest action of txn T_i
MasterRecord	Disk	LSN of latest checkpoint
undoNextLSN	log record	LSN of prior to-be-undone record

Constraints

- DRAM is volatile
- Avoid random writes to database on disk (NO FORCE)
- Support transactions with change sets > DRAM (STEAL)
- Recovery time must be bounded.
- Users must not suffer long delays due to checkpointing.
- Cope with repeated failures during recovery.
- Increase concurrency of undo.

V7: LOGICAL UNDO

- Record logical operations to be undone instead of physical offsets
 - ▶ Undo action need not be exact physical inverse of original action (*i.e.*, page offsets need not be recorded)
 - ▶ Example: Insert key X in B+tree
 - ▶ X can be initially inserted in Page 10 by T_1
 - ▶ X may be moved to Page 20 by another txn T_2 before T_1 commits
 - ▶ Later, if T_1 is aborted, logical undo (Delete key X in B+tree) will automatically remove it from Page 20

V7: LOGICAL UNDO

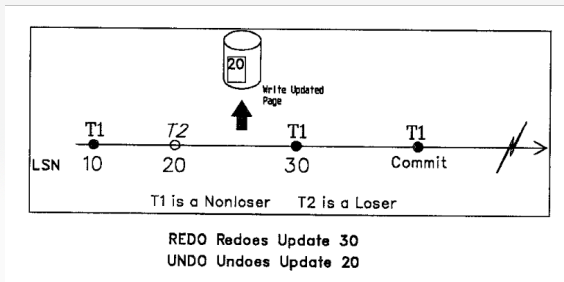
- Logical undo enables:
 - ▶ Highly-parallel transaction-oriented logical undo
 - ▶ Works with fast page-oriented physical redo
 - ▶ Hence, this protocol performs physiological logging
- Record logical ops for index and space management (*i.e.*, garbage collection)
 - ▶ Avoid rebuilding indexes from scratch during recovery
 - ▶ Reclaim storage space of deleted records
 - ▶ Example: Put in slot 5 (instead of Put at offset 30)

Constraints

- DRAM is volatile
- Avoid random writes to database on disk (NO FORCE)
- Support transactions with change sets > DRAM (STEAL)
- Recovery time must be bounded.
- Users must not suffer long delays due to checkpointing.
- Cope with repeated failures during recovery.
- Increase concurrency of undo (logical undo).
- Support record-level locking

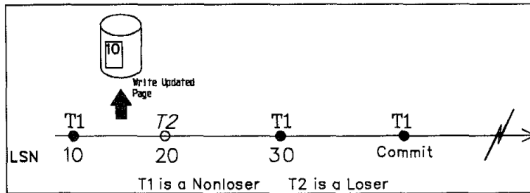
V8: AVOID SELECTIVE REDO

- Problem-free scenario



V8: AVOID SELECTIVE REDO

- Problematic scenario: **UNDOing non-existent changes**



REDO Redoes Update 30

UNDO Will Try to Undo 20 Even
Though Update Is NOT on Page

ERROR?!

- Solution:
 - ▶ Replay history of both committed and uncommitted transactions
 - ▶ Rather than selectively redo-ing committed transactions.
 - ▶ Then state of database guaranteed to be equivalent to that at the time of failure

Conclusion

- Protocols evolve over time to better handle user, workload, and hardware constraints.
- Deconstructing protocols will help you better appreciate the internals of complex software systems and learn the art of designing protocols.

