



Database Systems

Lecture23 – Chapter 19: Recovery System



Beomseok Nam (남범석)
bnam@skku.edu



Recovery Algorithm



Recovery Algorithm

- **So far:** we covered key concepts
- **Now:** we present the components of the basic recovery algorithm
- **Later:** we present extensions to allow more concurrency



Recovery Algorithm

- **Logging** (during normal operation):
 - $\langle T_i \text{ start} \rangle$ at transaction start
 - $\langle T_i, X_j, V_1, V_2 \rangle$ for each update, and
 - $\langle T_i \text{ commit} \rangle$ at transaction end
- **Transaction rollback (during normal operation)**
 - Let T_i be the transaction to be rolled back
 - Scan log backwards from the end, and for each log record of T_i of the form $\langle T_i, X_j, V_1, V_2 \rangle$
 - Perform the undo by writing V_1 to X_j ,
 - Write a log record $\langle T_i, X_j, V_1 \rangle$
 - such log records are called **compensation log records**
 - Once the record $\langle T_i \text{ start} \rangle$ is found stop the scan and write the log record $\langle T_i \text{ abort} \rangle$



Recovery Algorithm (Cont.)

- **Recovery from failure:** Two phases
 - **Redo phase:** replay updates of **all** transactions, whether they committed, aborted, or are incomplete
 - **Undo phase:** undo all incomplete transactions
- **Redo phase:**
 1. Find last **<checkpoint L>** record, and set undo-list to L .
 2. Scan forward from above **<checkpoint L>** record
 1. Whenever a record $\langle T_i, X_j, V_1, V_2 \rangle$ or $\langle T_i, X_j, V_2 \rangle$ is found, redo it by writing V_2 to X_j
 2. Whenever a log record $\langle T_i \text{ start} \rangle$ is found, add T_i to undo-list
 3. Whenever a log record $\langle T_i \text{ commit} \rangle$ or $\langle T_i \text{ abort} \rangle$ is found, remove T_i from undo-list



Recovery Algorithm (Cont.)

- **Undo phase:**

- 1. Scan log backwards from end

- 1. Whenever a log record $\langle T_i, X_j, V_1, V_2 \rangle$ is found where T_i is in undo-list perform same actions as for transaction rollback:

- 1. perform undo by writing V_1 to X_j .

- 2. write a log record $\langle T_i, X_j, V_1 \rangle$

- 2. Whenever a log record $\langle T_i \text{ start} \rangle$ is found where T_i is in undo-list,

- 1. Write a log record $\langle T_i \text{ abort} \rangle$

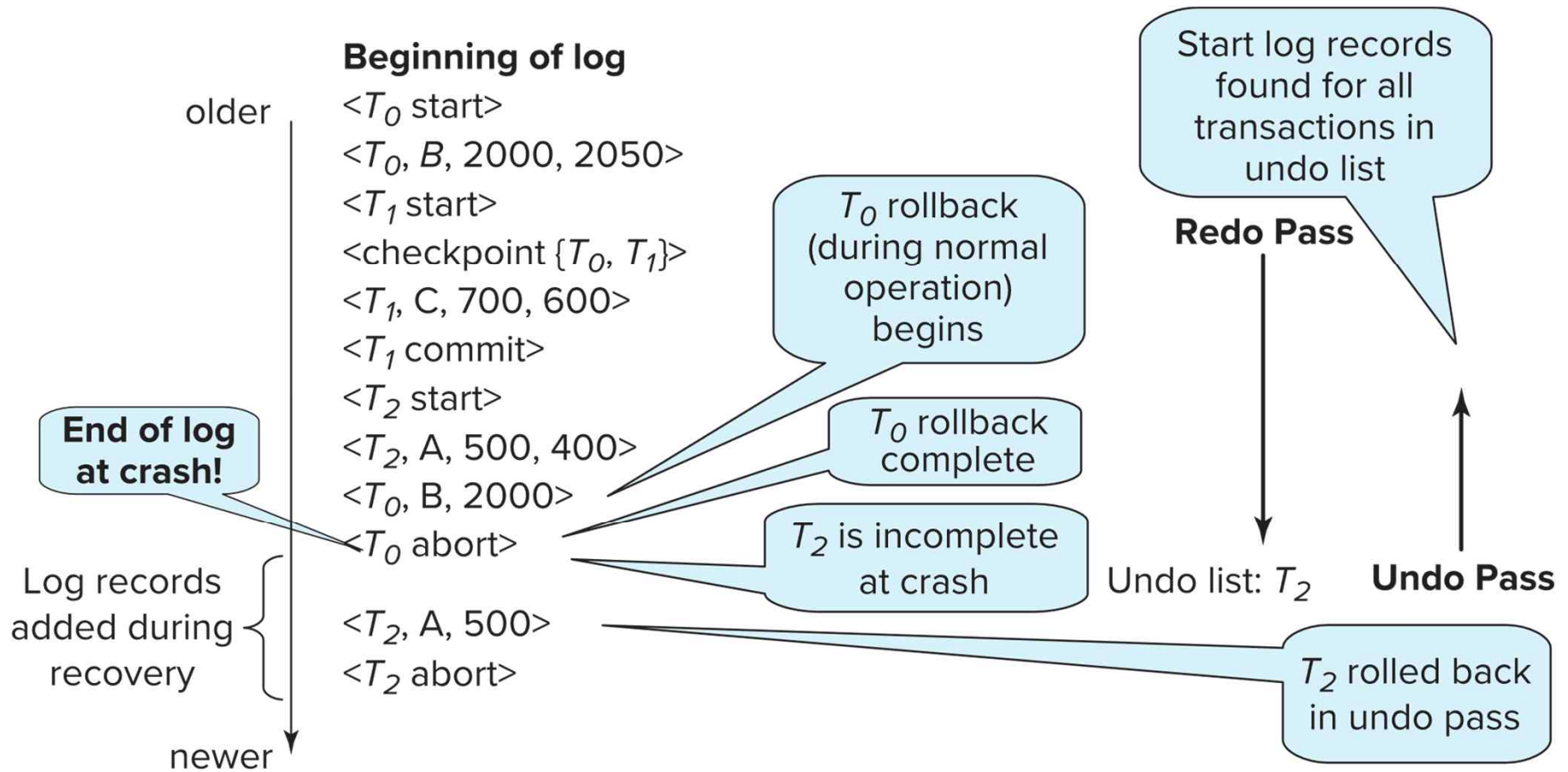
- 2. Remove T_i from undo-list

- 3. Stop when undo-list is empty

- 1. i.e., $\langle T_i \text{ start} \rangle$ has been found for every transaction in undo-list

- After undo phase completes, normal transaction processing can commence

Example of Recovery





Log Record Buffering

- **Log record buffering:** log records are buffered in main memory, instead of being output directly to stable storage.
 - Log records are output to stable storage when a block of log records in the buffer is full, or a **log force** operation is executed.
- Log force is performed to commit a transaction by forcing all its log records (including the commit record) to stable storage.
- Several log records can thus be output using a single output operation, reducing the I/O cost.



Log Record Buffering (Cont.)

- The rules below must be followed if log records are buffered:
 - Log records are output to stable storage in the order in which they are created.
 - Transaction T_i enters the commit state only when the log record $\langle T_i, \text{commit} \rangle$ has been output to stable storage.
 - Before a block of data in main memory is output to the database, all log records pertaining to data in that block must have been output to stable storage.
 - This rule is called the **write-ahead logging** or **WAL** rule
 - Strictly speaking, WAL only requires undo information to be output



Database Buffering

- Database maintains an in-memory buffer of data blocks
 - When a new block is needed, if buffer is full an existing block needs to be removed from buffer
 - If the block chosen for removal has been updated, it must be output to disk
- The recovery algorithm supports the **no-force policy**: i.e., updated blocks need not be written to disk when transaction commits
 - **force policy**: requires updated blocks to be written at commit
 - More expensive commit
- The recovery algorithm supports the **steal policy**: i.e., blocks containing updates of uncommitted transactions can be written to disk, even before the transaction commits



Database Buffering (Cont.)

- If a block with uncommitted updates is output to disk, log records with undo information for the updates are output to the log on stable storage first
 - (Write ahead logging)
- No updates should be in progress on a block when it is output to disk. Can be ensured as follows.
 - Before writing a data item, transaction acquires exclusive lock on block containing the data item
 - Lock can be released once the write is completed.
 - Such locks held for short duration are called **latches**.
- **To output a block to disk**
 1. First acquire an exclusive latch on the block
 - Ensures no update can be in progress on the block
 2. Then perform a **log flush**
 3. Then output the block to disk
 4. Finally release the latch on the block



Buffer Management (Cont.)

- Database buffer can be implemented either
 - In an area of real main-memory reserved for the database, or
 - In virtual memory
- Implementing buffer in reserved main-memory has drawbacks:
 - Memory is partitioned before-hand between database buffer and applications, limiting flexibility.
 - Needs may change, and although operating system knows best how memory should be divided up at any time, it cannot change the partitioning of memory.



Buffer Management (Cont.)

- Database buffers are generally implemented in virtual memory in spite of some drawbacks:
 - When operating system needs to evict a page that has been modified, the page is written to swap space on disk.
 - When database decides to write buffer page to disk, buffer page may be in swap space, and may have to be read from swap space on disk and output to the database on disk, resulting in extra I/O!
 - Known as **dual paging** problem.
 - Ideally when OS needs to evict a page from the buffer, it should pass control to database, which in turn should
 1. Output the page to database instead of to swap space (making sure to output log records first), if it is modified
 2. Release the page from the buffer, for the OS to use
 - Dual paging can thus be avoided, but common operating systems do not support such functionality.

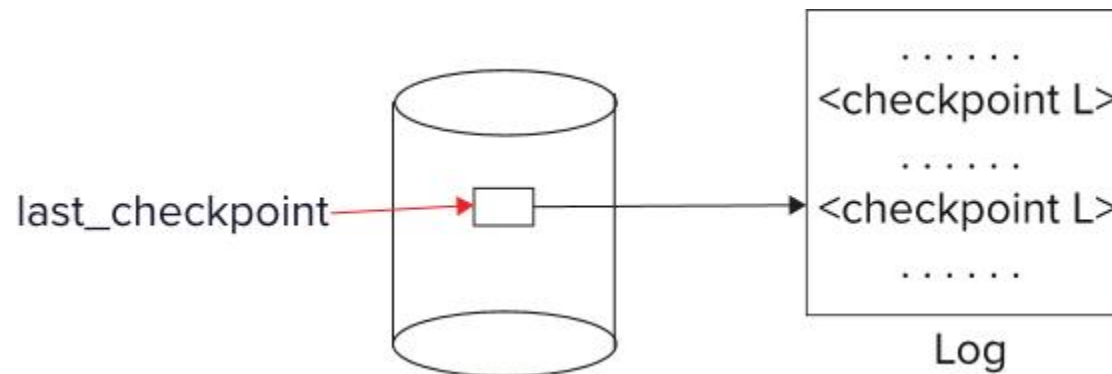


Fuzzy Checkpointing

- To avoid long interruption of normal processing during checkpointing, allow updates to happen during checkpointing
- **Fuzzy checkpointing** is done as follows:
 1. Temporarily stop all updates by transactions
 2. Write a **<checkpoint L>** log record and force log to stable storage
 3. Note list *M* of modified buffer blocks
 4. Now permit transactions to proceed with their actions
 5. Output to disk all modified buffer blocks in list *M*
 - blocks should not be updated while being output
 - Follow WAL: all log records pertaining to a block must be output before the block is output
 6. Store a pointer to the **checkpoint** record in a fixed position **last_checkpoint** on disk

Fuzzy Checkpointing (Cont.)

- When recovering using a fuzzy checkpoint, start scan from the **checkpoint** record pointed to by **last_checkpoint**
 - Log records before **last_checkpoint** have their updates reflected in database on disk, and need not be redone.
 - Incomplete checkpoints, where system had crashed while performing checkpoint, are handled safely





Logical Undo Logging

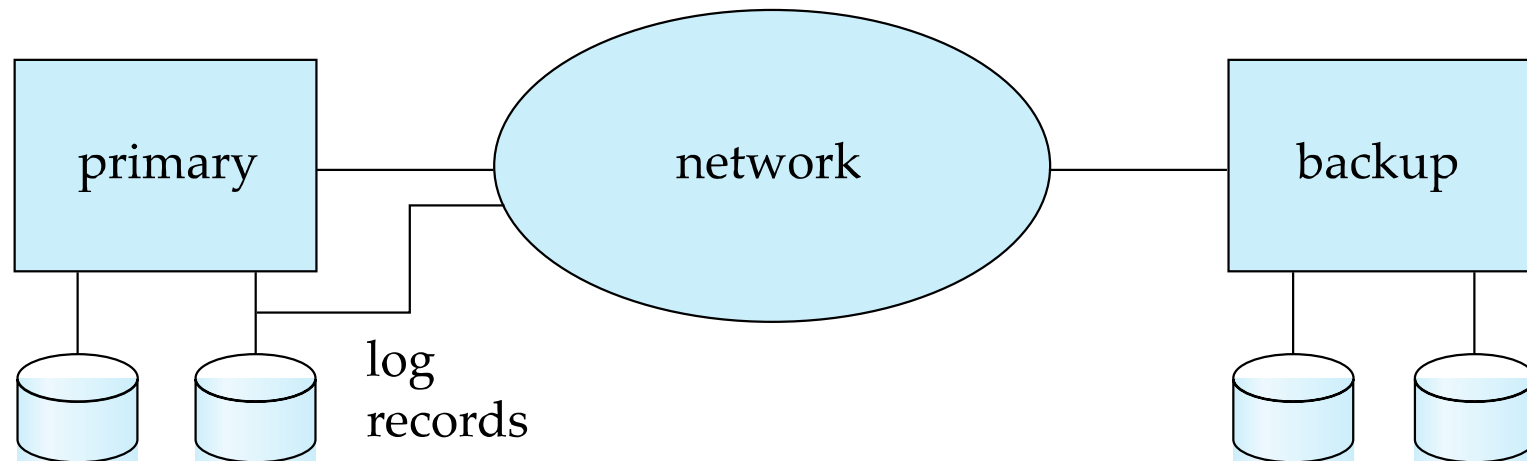
- Operations like B⁺-tree insertions and deletions release locks early.
 - They cannot be undone by restoring old values (**physical undo**), since once a lock is released, other transactions may have updated the B⁺-tree.
 - Instead, insertions (resp. deletions) are undone by executing a deletion (resp. insertion) operation (known as **logical undo**).
- For such operations, undo log records should contain the undo operation to be executed
 - Such logging is called **logical undo logging**, in contrast to **physical undo logging**
 - Operations are called **logical operations**
 - Other examples:
 - delete of tuple, to undo insert of tuple
 - allows early lock release on space allocation information
 - subtract amount deposited, to undo deposit
 - allows early lock release on bank balance



Remote Backup Systems

Remote Backup Systems

- Remote backup systems provide high availability by allowing transaction processing to continue even if the primary site is destroyed.





Remote Backup Systems (Cont.)

- **Detection of failure:** Backup site must detect when primary site has failed
 - to distinguish primary site failure from link failure maintain several communication links between the primary and the remote backup.
 - Heart-beat messages
- **Transfer of control:**
 - To take over control backup site first perform recovery using its copy of the database and all the log records it has received from the primary.
 - Thus, completed transactions are redone and incomplete transactions are rolled back.
 - When the backup site takes over processing it becomes the new primary
 - To transfer control back to old primary when it recovers, old primary must receive redo logs from the old backup and apply all updates locally.



Remote Backup Systems (Cont.)

- **Time to recover:** To reduce delay in takeover, backup site periodically process the redo log records (in effect, performing recovery from previous database state), performs a checkpoint, and can then delete earlier parts of the log.
- **Hot-Spare** configuration permits very fast takeover:
 - Backup continually processes redo log record as they arrive, applying the updates locally.
 - When failure of the primary is detected the backup rolls back incomplete transactions, and is ready to process new transactions.
- Alternative to remote backup: distributed database with replicated data
 - Remote backup is faster and cheaper, but less tolerant to failure



Remote Backup Systems (Cont.)

- Ensure durability of updates by delaying transaction commit until update is logged at backup; avoid this delay by permitting lower degrees of durability.
- **One-safe**: commit as soon as transaction's commit log record is written at primary
 - Problem: updates may not arrive at backup before it takes over.
- **Two-very-safe**: commit when transaction's commit log record is written at primary and backup
 - Reduces availability since transactions cannot commit if either site fails.
- **Two-safe**: proceed as in two-very-safe if both primary and backup are active. If only the primary is active, the transaction commits as soon as its commit log record is written at the primary.
 - Better availability than two-very-safe
 - Avoids problem of lost transactions in one-safe.



ARIES Recovery Algorithm



ARIES

- ARIES is a state of the art recovery method
 - Incorporates numerous optimizations to reduce overheads during normal processing and to speed up recovery
 - The recovery algorithm we studied earlier is modeled after ARIES, but greatly simplified by removing optimizations
- Unlike the recovery algorithm described earlier, ARIES
 1. Uses **log sequence number (LSN)** to identify log records
 - Stores LSNs in pages to identify what updates have already been applied to a database page
 2. Physiological redo
 3. Dirty page table to avoid unnecessary redos during recovery
 4. Fuzzy checkpointing that only records information about dirty pages,
and does not require dirty pages to be written out at checkpoint time
 - More coming up on each of the above ...



ARIES Optimizations

■ **Physiological redo**

- Affected page is physically identified, action within page can be logical
 - Used to reduce logging overheads
 - e.g. when a record is deleted and all other records have to be moved to fill hole
 - Physiological redo can log just the record deletion
 - Physical redo would require logging of old and new values for much of the page
 - Requires page to be output to disk atomically
 - Easy to achieve with hardware RAID, also supported by some disk systems
 - Incomplete page output can be detected by checksum techniques,
 - But extra actions are required for recovery
 - Treated as a media failure



ARIES Data Structures

- ARIES uses several data structures
 - Log sequence number (LSN) identifies each log record
 - Must be sequentially increasing
 - Typically an offset from beginning of log file to allow fast access
 - Easily extended to handle multiple log files
 - Page LSN
 - Log records of several different types
 - Dirty page table



ARIES Data Structures: Page LSN

- Each page contains a **PageLSN** which is the LSN of the last log record whose effects are reflected on the page
 - To update a page:
 - X-latch the page, and write the log record
 - Update the page
 - Record the LSN of the log record in PageLSN
 - Unlock page
 - To flush page to disk, must first S-latch page
 - Thus page state on disk is operation consistent
 - Required to support physiological redo
 - PageLSN is used during recovery to prevent repeated redo
 - Thus ensuring idempotence

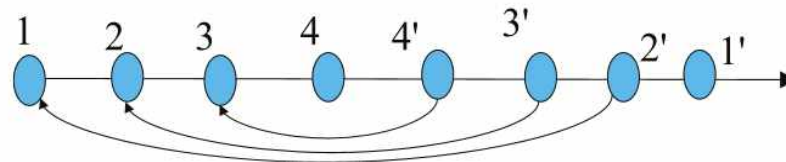
ARIES Data Structures: Log Record

- Each log record contains LSN of previous log record of the same transaction

LSN	TransID	PrevLSN	RedoInfo	UndoInfo
-----	---------	---------	----------	----------

- LSN in log record may be implicit
- Special redo-only log record called **compensation log record (CLR)** used to log actions taken during recovery that never need to be undone
 - Serves the role of operation-abort log records used in earlier recovery algorithm
 - Has a field UndoNextLSN to note next (earlier) record to be undone
 - Records in between would have already been undone
 - Required to avoid repeated undo of already undone actions

LSN	TransID	UndoNextLSN	RedoInfo
-----	---------	-------------	----------



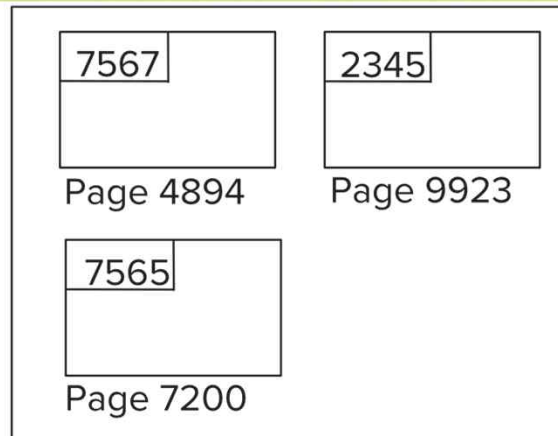


ARIES Data Structures: DirtyPage Table

▪ **DirtyPageTable**

- List of pages in the buffer that have been updated
- Contains, for each such page
 - **PageLSN** of the page
 - **RecLSN** is an LSN such that log records before this LSN have already been applied to the page version on disk
 - Set to current end of log when a page is inserted into dirty page table (just before being updated)
 - Recorded in checkpoints, helps to minimize redo work

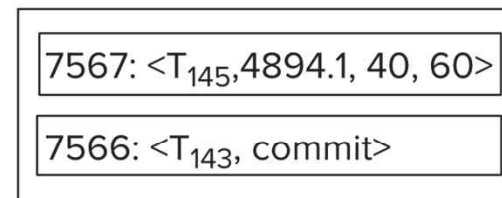
ARIES Data Structures



Database Buffer

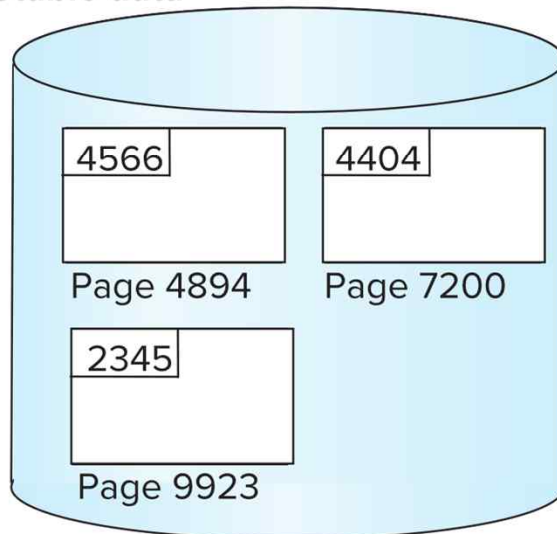
PageID	PageLSN	RecLSN
4894	7567	7564
7200	7565	7565

Dirty Page Table

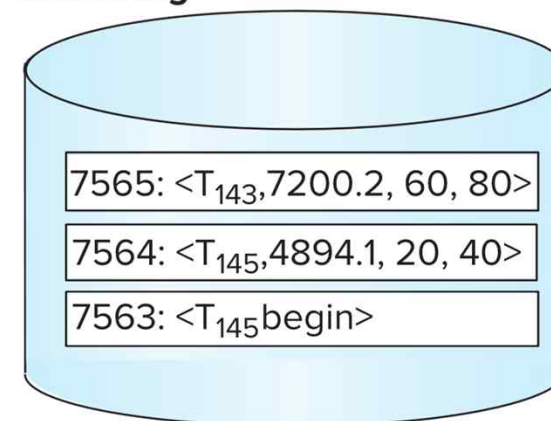


Log Buffer (PrevLSN and UndoNextLSN fields not shown)

Stable data



Stable log





ARIES Data Structures: Checkpoint Log

- **Checkpoint log record**

- Contains:
 - DirtyPageTable and list of active transactions
 - For each active transaction, LastLSN, the LSN of the last log record written by the transaction
- Fixed position on disk notes LSN of last completed checkpoint log record

- Dirty pages are not written out at checkpoint time
 - Instead, they are flushed out continuously, in the background
- Checkpoint is thus very low overhead
 - can be done frequently



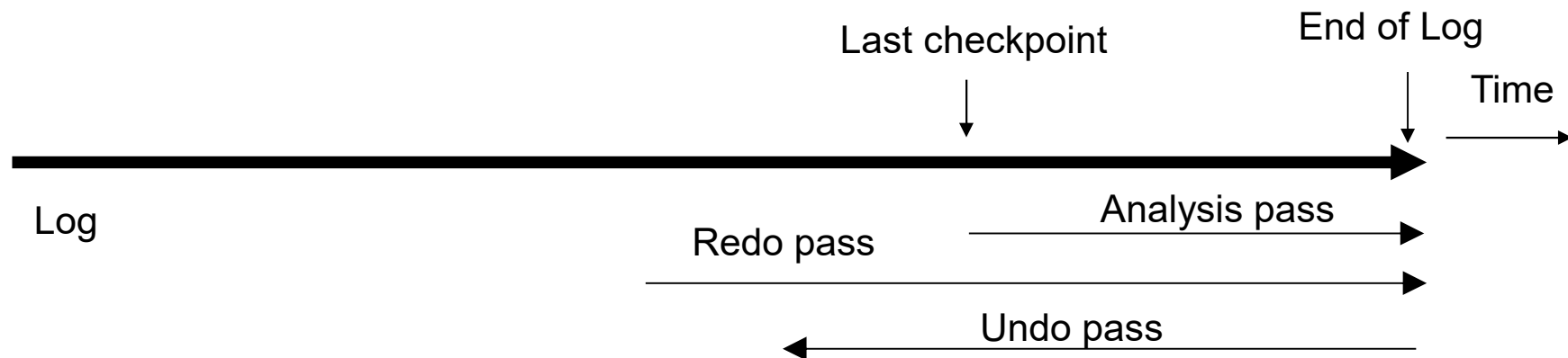
ARIES Recovery Algorithm

ARIES recovery involves three passes

- **Analysis pass:** Determines
 - Which transactions to undo
 - Which pages were dirty (disk version not up to date) at time of crash
 - **RedoLSN:** LSN from which redo should start
- **Redo pass:**
 - Repeats history, redoing all actions from RedoLSN
 - RecLSN and PageLSNs are used to avoid redoing actions already reflected on page
- **Undo pass:**
 - Rolls back all incomplete transactions
 - Transactions whose abort was complete earlier are not undone
 - Key idea: no need to undo these transactions: earlier undo actions were logged, and are redone as required

Aries Recovery: 3 Passes

- Analysis, redo and undo passes
- Analysis determines where redo should start
- Undo has to go back till start of earliest incomplete transaction





ARIES Recovery: Analysis

Analysis pass

- Starts from last complete checkpoint log record
 - Reads DirtyPageTable from log record
 - Sets RedoLSN = min of RecLSNs of all pages in DirtyPageTable
 - In case no pages are dirty, RedoLSN = checkpoint record's LSN
 - Sets undo-list = list of transactions in checkpoint log record
 - Reads LSN of last log record for each transaction in undo-list from checkpoint log record
- Scans forward from checkpoint
- .. Cont. on next page ...



ARIES Recovery: Analysis (Cont.)

Analysis pass (cont.)

- Scans forward from checkpoint
 - If any log record found for transaction not in undo-list, adds transaction to undo-list
 - Whenever an update log record is found
 - If page is not in DirtyPageTable, it is added with RecLSN set to LSN of the update log record
 - If transaction end log record found, delete transaction from undo-list
 - Keeps track of last log record for each transaction in undo-list
 - May be needed for later undo
- At end of analysis pass:
 - RedoLSN determines where to start redo pass
 - RecLSN for each page in DirtyPageTable used to minimize redo work
 - All transactions in undo-list need to be rolled back



ARIES Redo Pass

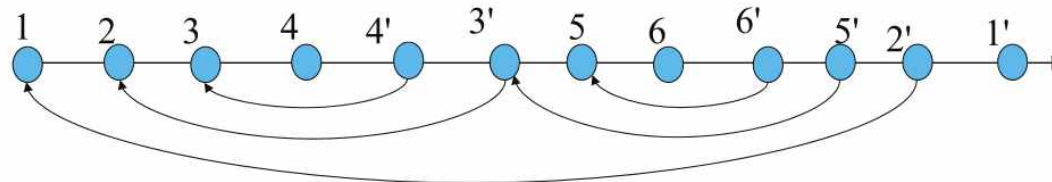
Redo Pass: Repeats history by replaying every action not already reflected in the page on disk, as follows:

- Scans forward from RedoLSN. Whenever an update log record is found:
 1. If the page is not in DirtyPageTable or the LSN of the log record is less than the RecLSN of the page in DirtyPageTable, then skip the log record
 2. Otherwise fetch the page from disk. If the PageLSN of the page fetched from disk is less than the LSN of the log record, redo the log record

NOTE: if either test is negative the effects of the log record have already appeared on the page. First test avoids even fetching the page from disk!

ARIES Undo Actions

- When an undo is performed for an update log record
 - Generate a CLR containing the undo action performed (actions performed during undo are logged physically or physiologically).
 - CLR for record n noted as n' in figure below
 - Set UndoNextLSN of the CLR to the PrevLSN value of the update log record
 - Arrows indicate UndoNextLSN value
- ARIES supports partial rollback
 - Used e.g. to handle deadlocks by rolling back just enough to release reqd. locks
 - Figure indicates forward actions after partial rollbacks
 - records 3 and 4 initially, later 5 and 6 then full rollback



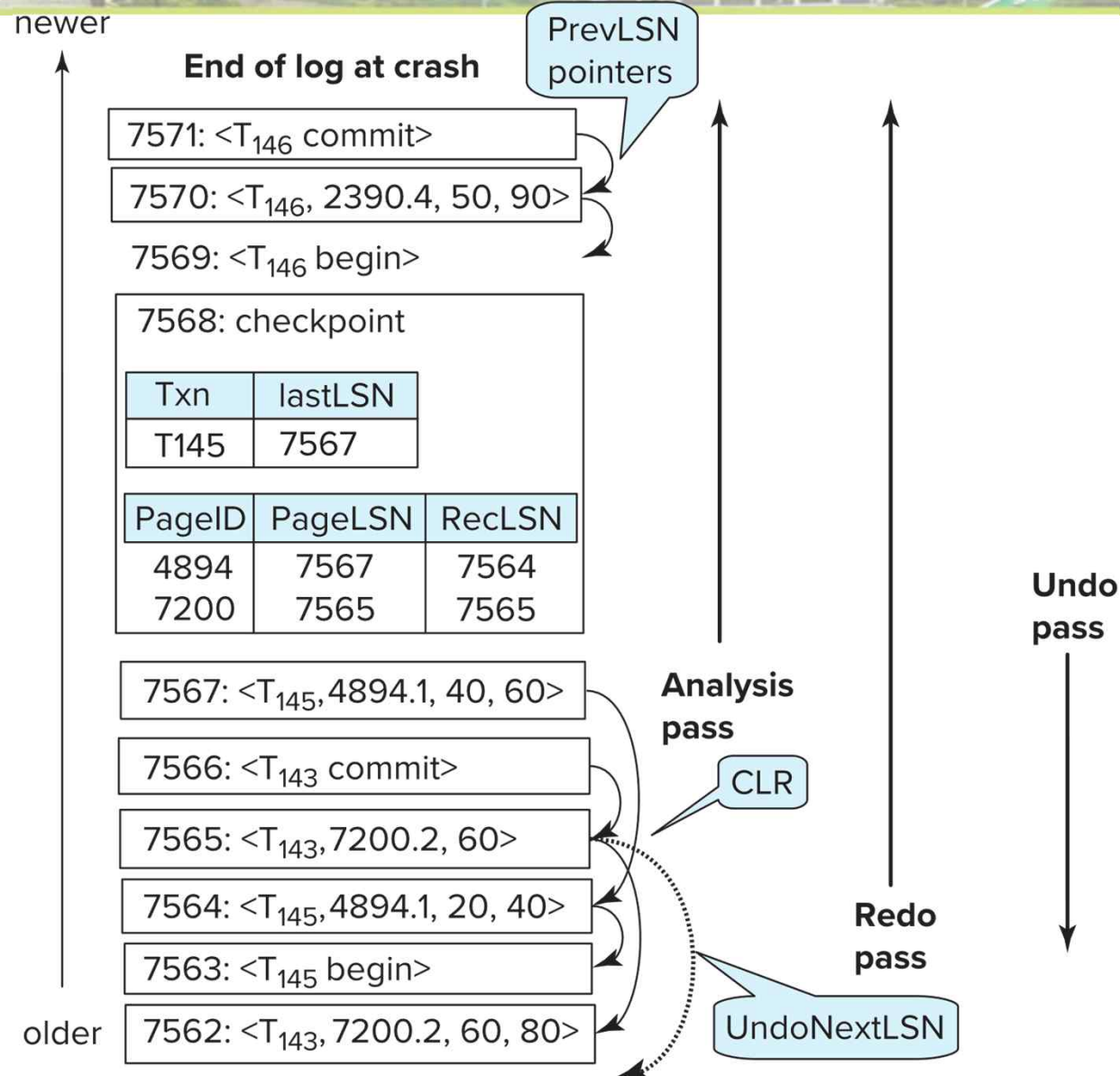


ARIES: Undo Pass

Undo pass:

- Performs backward scan on log undoing all transaction in undo-list
 - Backward scan optimized by skipping unneeded log records as follows:
 - Next LSN to be undone for each transaction set to LSN of last log record for transaction found by analysis pass.
 - At each step pick largest of these LSNs to undo, skip back to it and undo it
 - After undoing a log record
 - For ordinary log records, set next LSN to be undone for transaction to PrevLSN noted in the log record
 - For compensation log records (CLRs) set next LSN to be undo to UndoNextLSN noted in the log record
 - All intervening records are skipped since they would have been undone already
- Undos performed as described earlier

Recovery Actions in ARIES





Other ARIES Features

- Recovery Independence
 - Pages can be recovered independently of others
 - E.g. if some disk pages fail they can be recovered from a backup while other pages are being used
- Savepoints:
 - Transactions can record savepoints and roll back to a savepoint
 - Useful for complex transactions
 - Also used to rollback just enough to release locks on deadlock



Other ARIES Features (Cont.)

- Fine-grained locking:
 - Index concurrency algorithms that permit tuple level locking on indices can be used
 - These require logical undo, rather than physical undo, as in earlier recovery algorithm
- Recovery optimizations: For example:
 - Dirty page table can be used to **prefetch** pages during redo
 - Out of order redo is possible:
 - redo can be postponed on a page being fetched from disk, and performed when page is fetched.
 - Meanwhile other log records can continue to be processed



Recovery in Main Memory Databases

- Normal recovery algorithms can be used with main-memory databases
- But optimizations are possible
 - No redo-logging for indices: indices can be rebuilt quickly in-memory on recovery
 - No undo logging if only committed data is written to disk
 - Parallel recovery is important to load large amounts of data in-memory and perform recover with minimum delay