

Carnegie Mellon University

# DATABASE SYSTEMS

## Database Recovery

LECTURE #22 » 15-445/645 FALL 2025 » PROF. ANDY PAVLO



# ADMINISTRIVIA

---



**Homework #5** is due Sunday Nov 23<sup>rd</sup> @ 11:59pm

**Project #4** is due Sunday Dec 7<sup>th</sup> @ 11:59pm

→ Recitation Slides + Video ([@280](#))

**Final Exam** is on Thursday Dec 11<sup>th</sup> @ 1:00pm

→ Do not make travel plans before this date!

# CRASH RECOVERY

---



Recovery algorithms are techniques to ensure database consistency, transaction atomicity, and durability despite failures.

Recovery algorithms have two parts:

- Actions during normal txn 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.

*Today*

# CRASH RECOVERY OVERVIEW



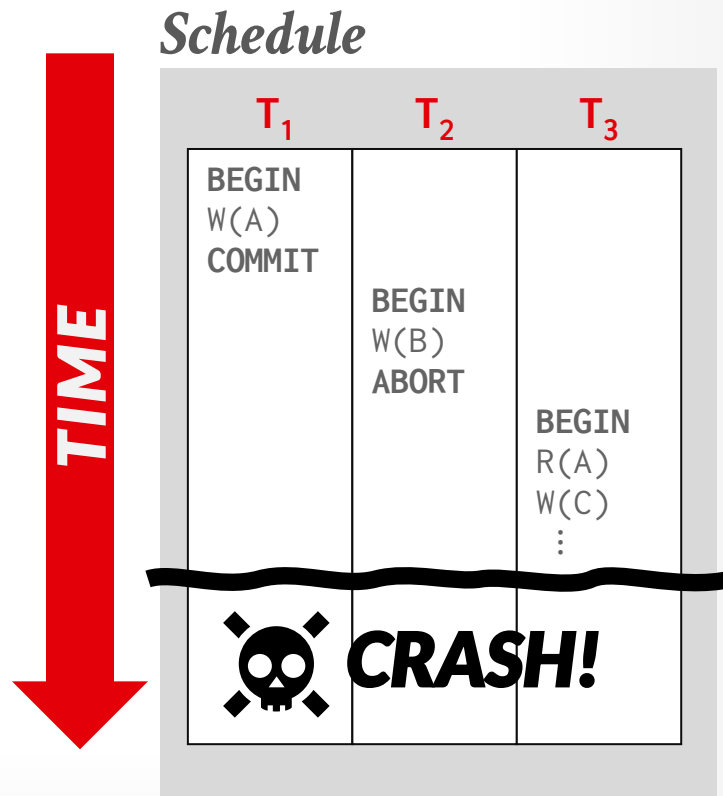
## STEAL + NO-FORCE

**Atomicity:** Txns may abort/fail.

**Durability:** Changes of committed txns should survive system failure.

Desired behavior after the DBMS restarts (i.e., the contents of volatile memory are lost):

- $T_1$  should be durable.
- $T_2 + T_3$  should be aborted.



## Algorithms for Recovery and Isolation Exploiting Semantics

Developed at IBM Research in early 1990s for the DB2 DBMS.

Not all systems implement ARIES exactly as defined in this paper but they're close enough.

### ARIES: A Transaction Recovery Method Supporting Fine-Granularity Locking and Partial Rollbacks Using Write-Ahead Logging

C. MOHAN  
IBM Almaden Research Center

and  
DON HADERLE  
IBM Santa Teresa Laboratory

and  
BRUCE LINDSAY, HAMID PIRAHESH and PETER SCHWARZ  
IBM Almaden Research Center

In this paper we present a simple and efficient method, called ARIES (*Algorithms for Recovery and Isolation Exploiting Semantics*), which supports partial rollbacks of transactions, fine-granularity (e.g., record) locking and recovery using write-ahead logging (WAL). We introduce the paradigm of *repeating history* to redo all missing updates *before* performing the rollbacks of the *lower* transactions during restart after a system failure. ARIES uses a log sequence number in each page to correlate the state of a page with respect to logged updates of that page. All updates of a transaction are logged, including those performed during rollbacks. By appropriate chaining of the log records written during rollbacks to those written during forward progress, a bounded amount of logging is ensured during rollbacks even in the face of repeated failures during restart or of nested rollbacks. We deal with a variety of features that are very important in building and operating an *industrial-strength* transaction processing system. ARIES supports fuzzy checkpoints, selective and deferred restart, fuzzy image copies, media recovery, and high concurrency lock modes (e.g., increment/decrement) which exploit the semantics of the operations and require the ability to perform operation logging. ARIES is flexible with respect to the kinds of buffer management policies that can be implemented. It supports objects of varying length efficiently. By enabling parallelism during restart, page-oriented redo, and logical undo, it enhances concurrency and performance. We show why some of the System R paradigms for logging and recovery, which were based on the shadow page technique, need to be changed in the context of WAL. We compare ARIES to the WAL-based recovery methods of

*Authors' addresses:* C. Mohan, Data Base Technology Institute, IBM Almaden Research Center, San Jose, CA 95120; D. Haderle, Data Base Technology Institute, IBM Santa Teresa Laboratory, San Jose, CA 95150; B. Lindsay, H. Pirahesh, and P. Schwarz, IBM Almaden Research Center, San Jose, CA 95120.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1992 0362-5915/92/0300-0094 \$1.50

ACM Transactions on Database Systems, Vol. 17, No. 1, March 1992, Pages 94-162

# ARIES: MAIN IDEAS

---



## Write-Ahead Logging:

- Flush WAL record(s) changes to disk before a database object is written to disk.
- Must use **STEAL** + **NO-FORCE** buffer pool policies.

## Repeating History During Redo:

- On DBMS restart, retrace actions and restore database to exact state before crash.

## Logging Changes During Undo:

- Record undo actions to log to ensure action is not repeated in the event of repeated failures.

# TODAY'S AGENDA

---



Log Sequence Numbers

Normal Commit & Abort Operations

Fuzzy Checkpointing

Recovery Algorithm

⚡DB Flash Talk: [ClickHouse](#)

# WAL RECORDS

---



We need to extend our log record format from last class to include additional info.

Every log record includes a globally unique **log sequence number** (LSN).

→ LSNs represent the physical order that txns make changes to the database.

Various components in the system keep track of **LSNs** that pertain to them...



# WAL BOOKKEEPING



Log Sequence Number (LSN).

→ Unique and monotonically increasing.

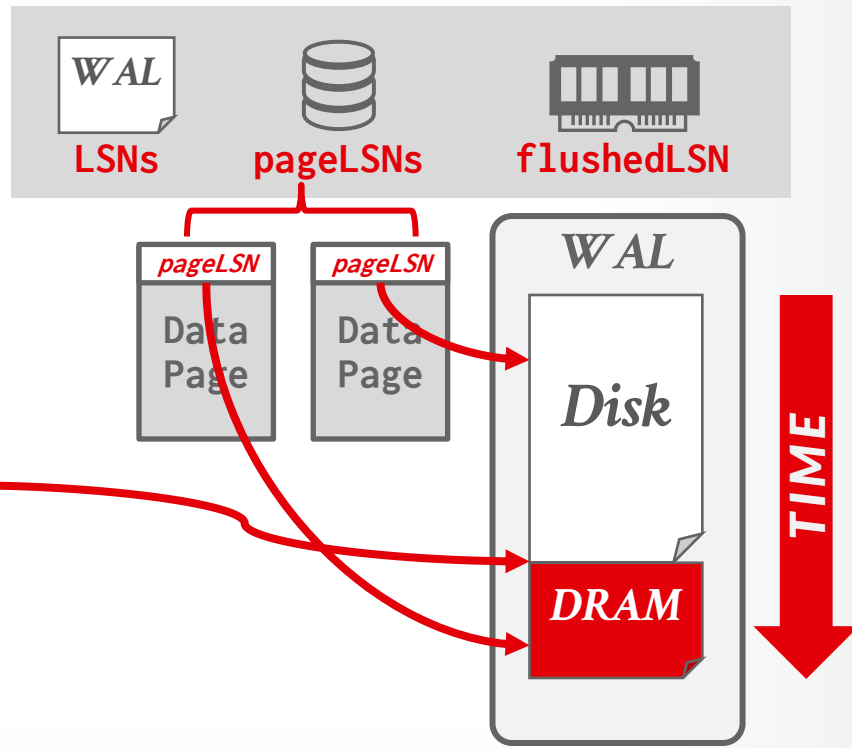
Each data page contains a **pageLSN**.

→ The LSN of the most recent log record that updated the page.

System keeps track of **flushedLSN**.

→ The max LSN flushed so far.

WAL: Before a  $\text{page}_x$  is written,  
 **$\text{pageLSN}_x \leq \text{flushedLSN}$**



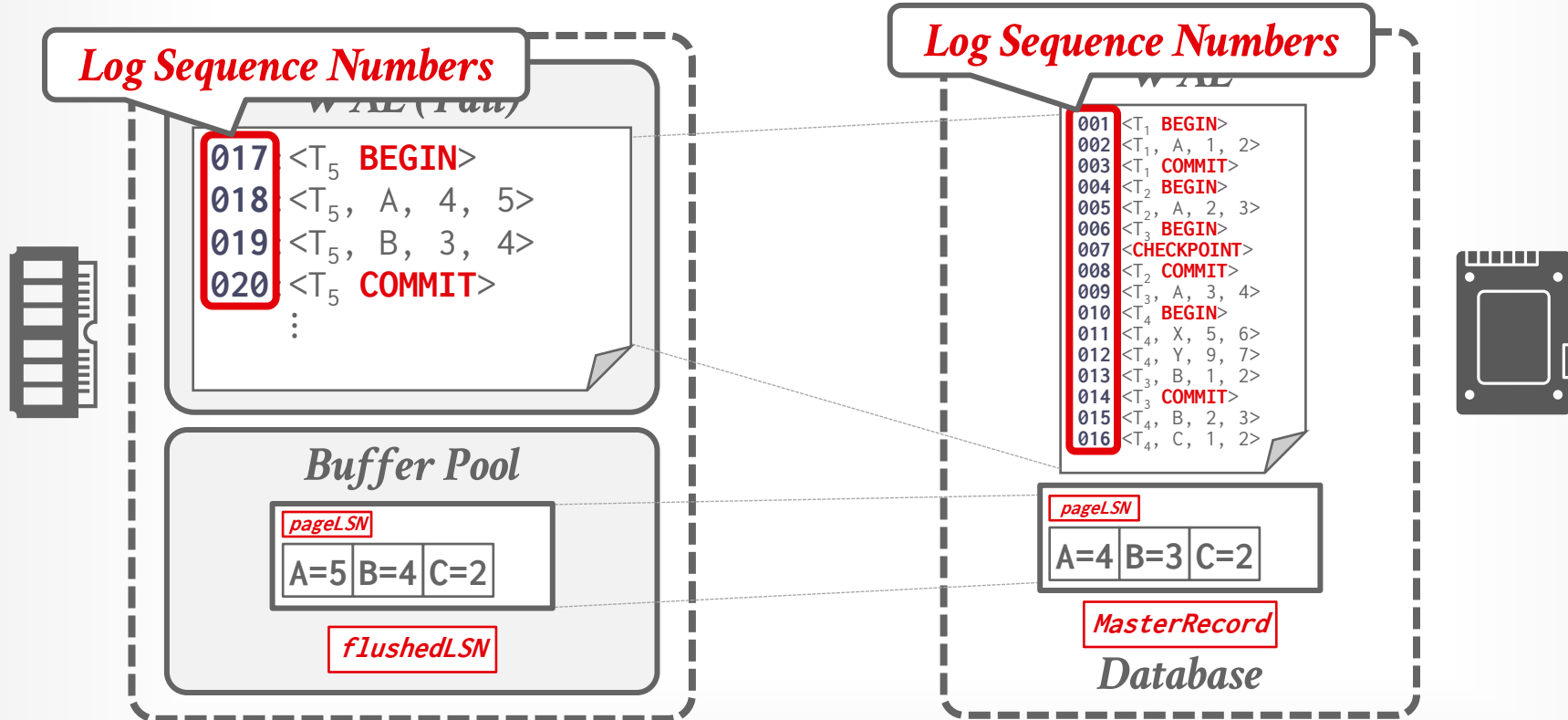
# LOG SEQUENCE NUMBERS

Name	Location	Definition
<b>flushedLSN</b>	Memory	Last LSN in log on disk
<b>pageLSN</b>	page <sub>x</sub>	Newest update to page <sub>x</sub>
<b>recLSN</b>	DPT <sup>♦</sup>	Oldest update to page <sub>x</sub> since it was last flushed
<b>lastLSN</b>	ATT <sup>♣</sup>	Latest record of txn T <sub>i</sub>
<b>MasterRecord</b>	Disk	LSN of latest checkpoint

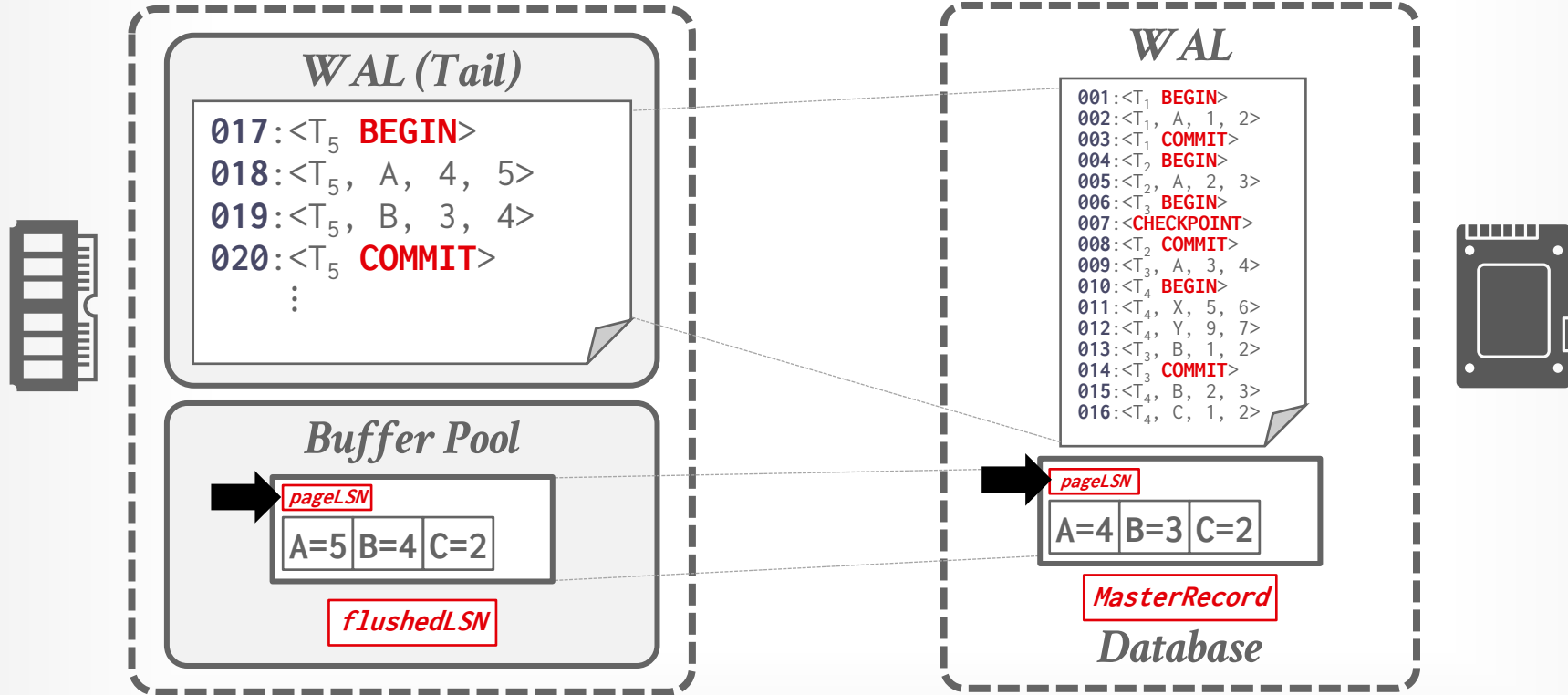
♦ DPT = Dirty Page Table.

♣ ATT = Active Transaction Table.

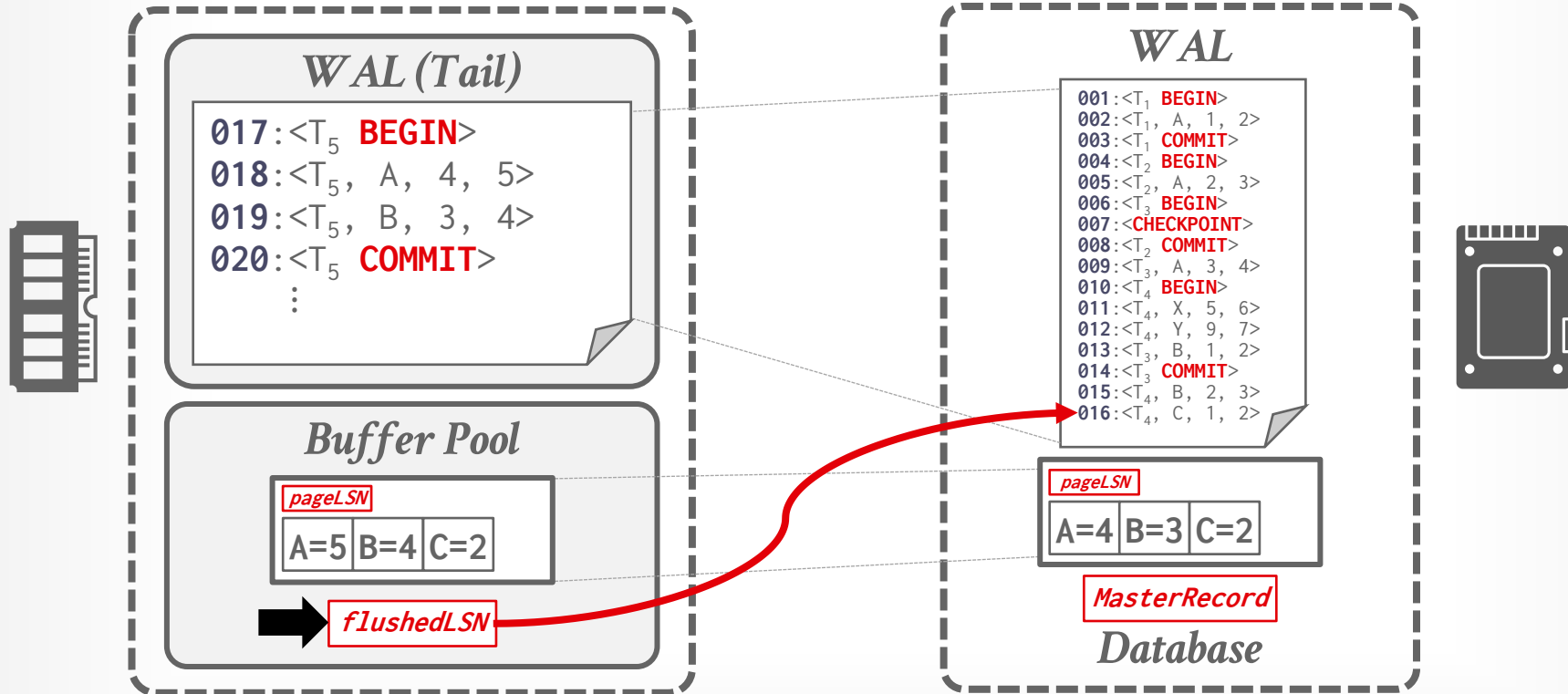
# WRITING LOG RECORDS



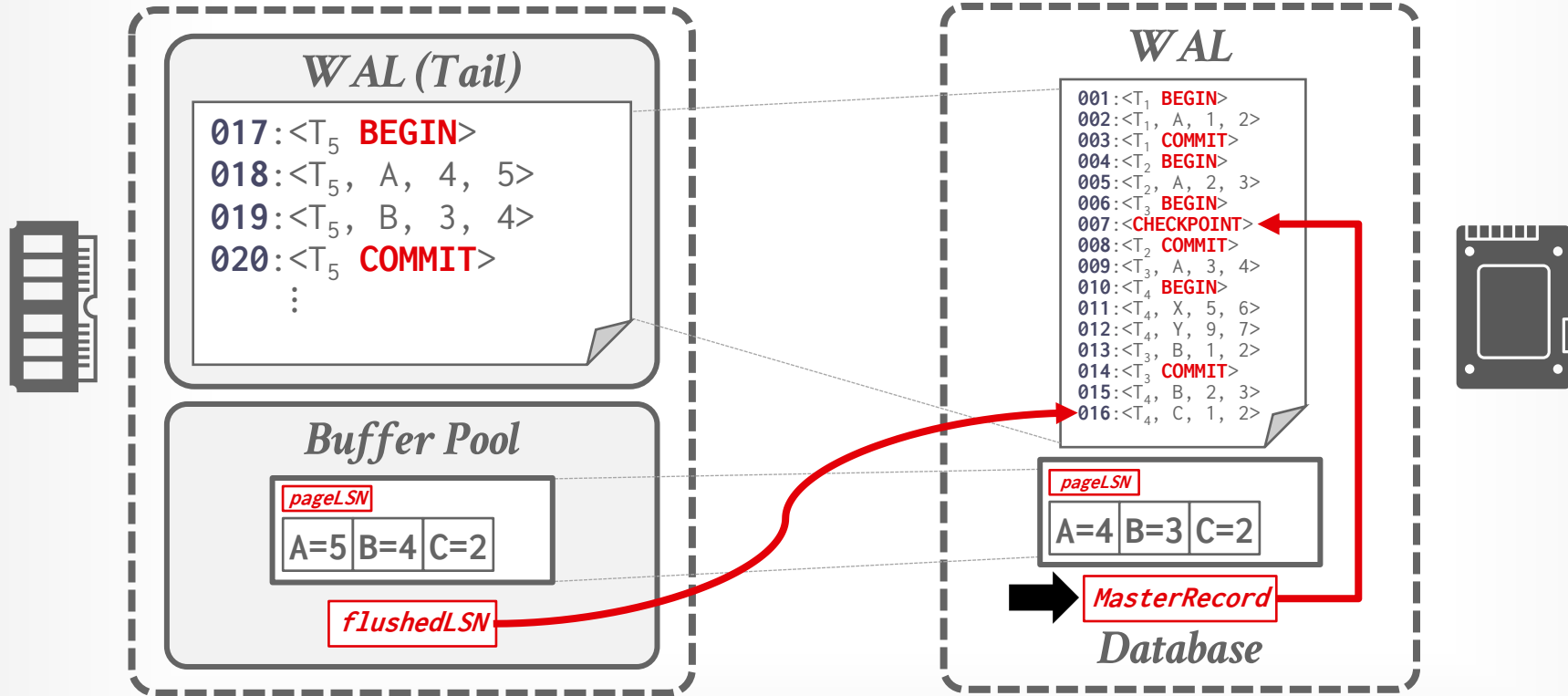
# WRITING LOG RECORDS



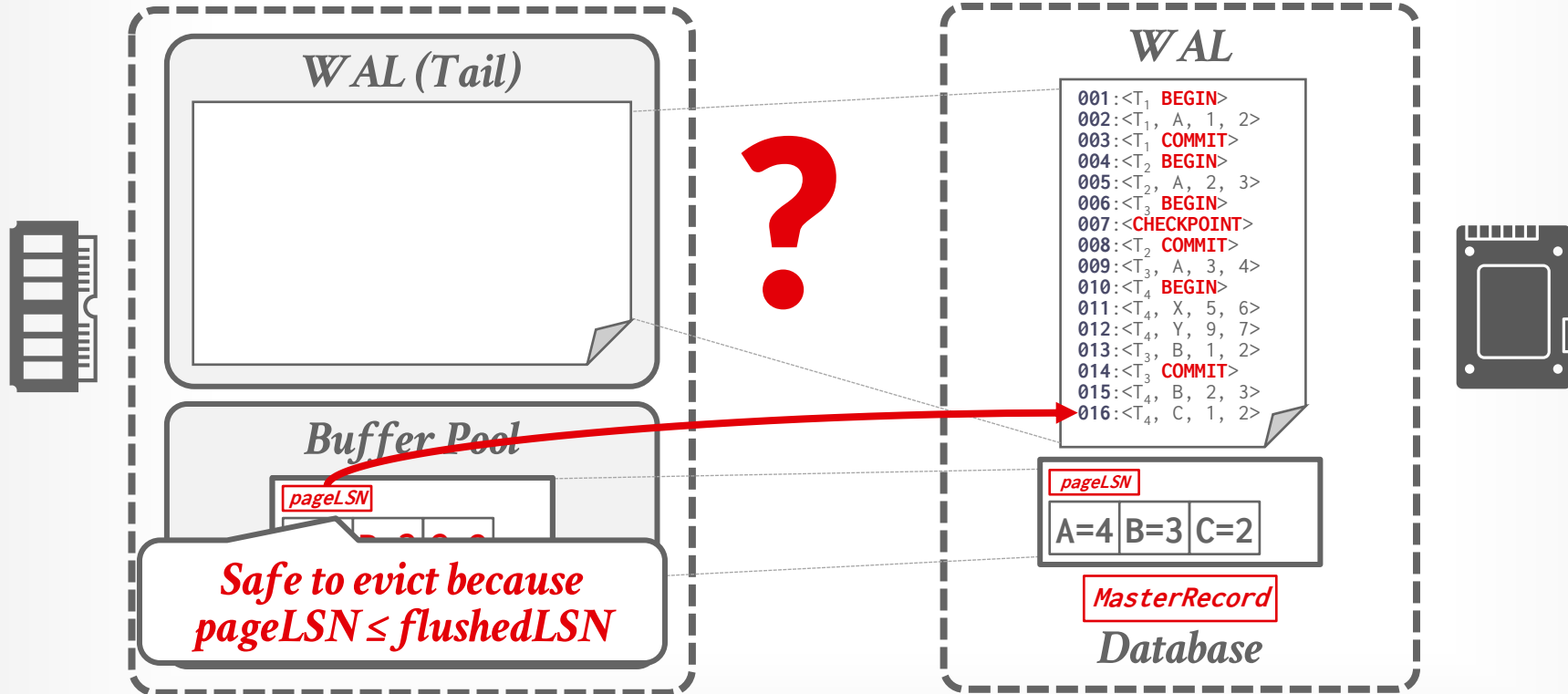
# WRITING LOG RECORDS



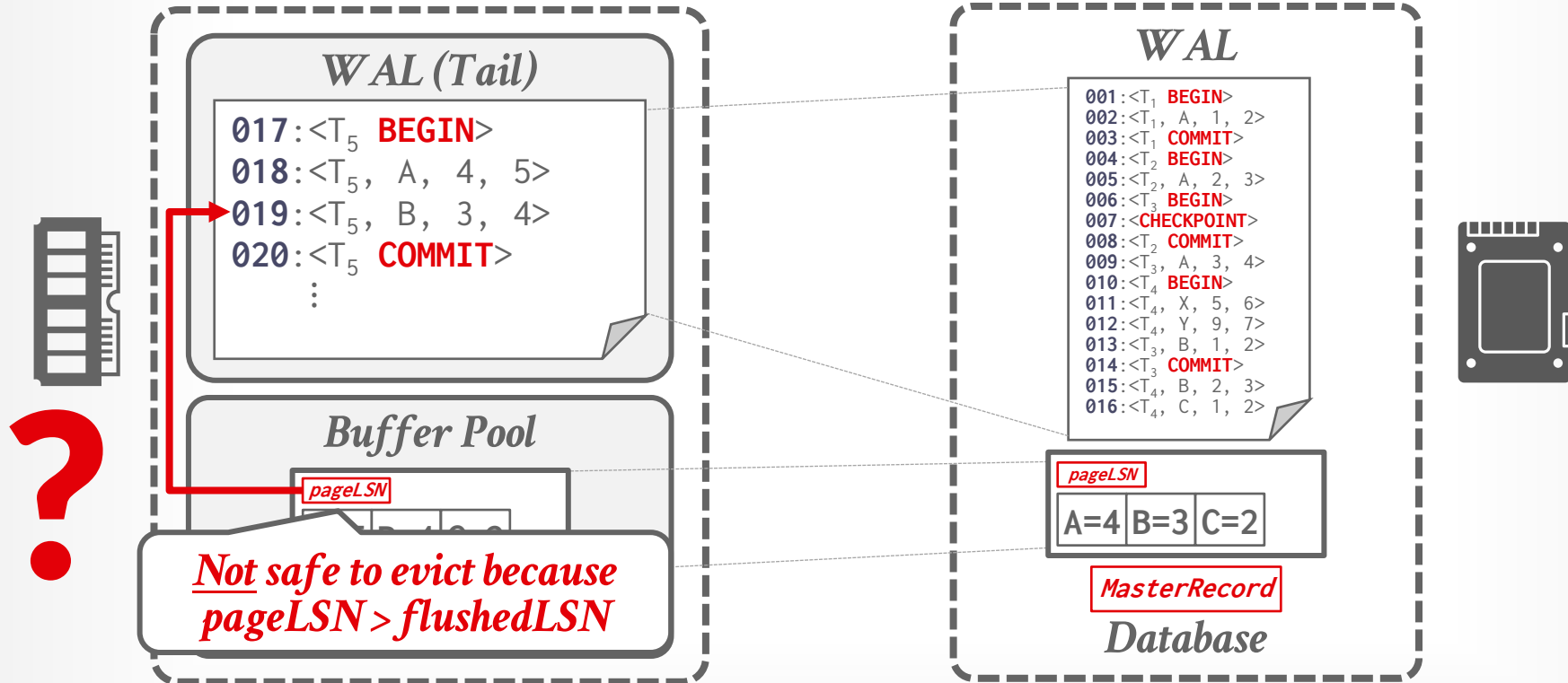
# WRITING LOG RECORDS



# WRITING LOG RECORDS



# WRITING LOG RECORDS





# WRITING LOG RECORDS

---

All log records have an LSN.

Update the pageLSN every time a txn modifies a record in the page.

Update the flushedLSN in memory every time the DBMS writes the WAL buffer to disk.

# NORMAL EXECUTION

---

Each txn invokes a sequence of reads and writes, followed by commit or rollback.

Assumptions in this lecture:

- All log records fit within a single page.
- Disk writes are atomic.
- Single-versioned tuples with Strong Strict 2PL.
- **STEAL** + **NO-FORCE** buffer management with WAL.
- Physical log record scheme.
- Omitting log records for indexes.

# TRANSACTION COMMIT

---

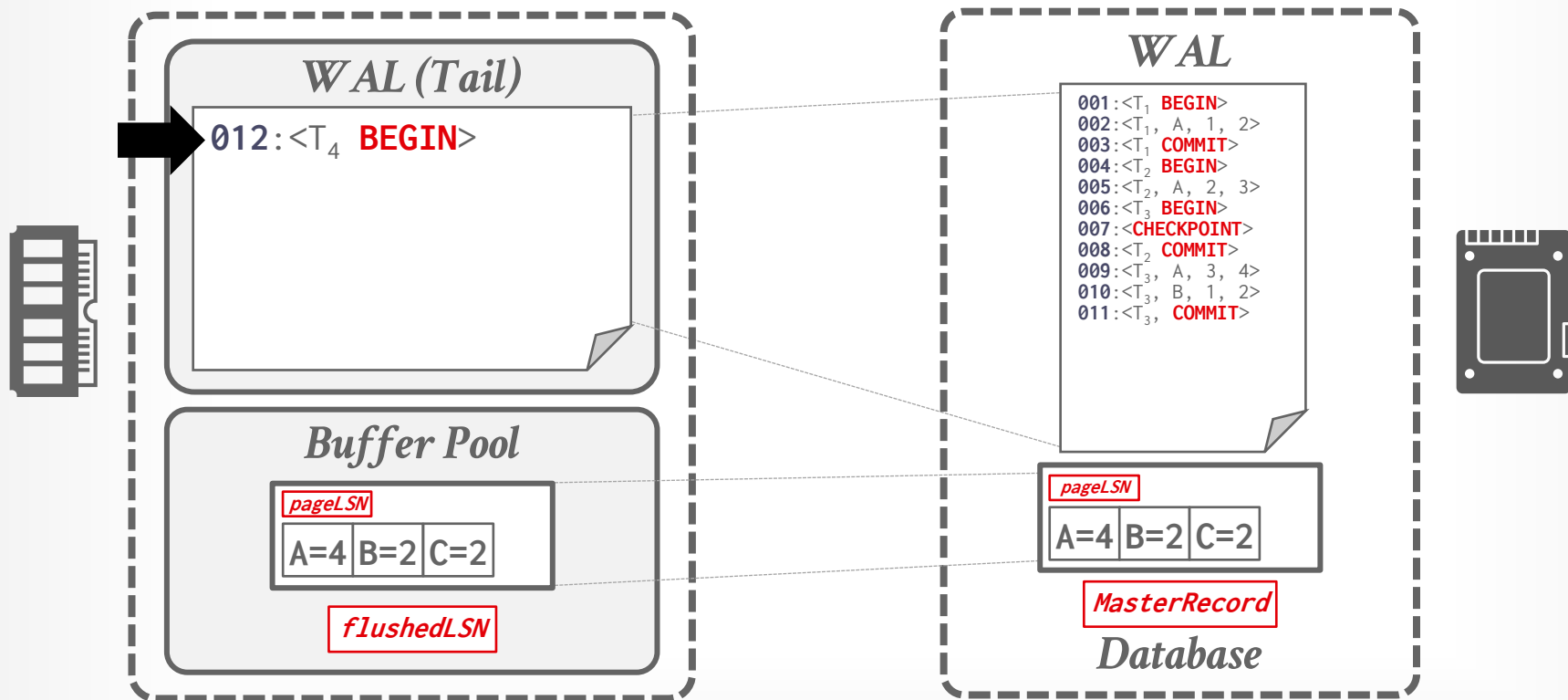
When a txn commits, the DBMS writes a **COMMIT** record to log and guarantees that all log records up to txn's **COMMIT** record are flushed to disk.

- Log flushes are sequential, synchronous writes to disk.
- Many log records per log page.

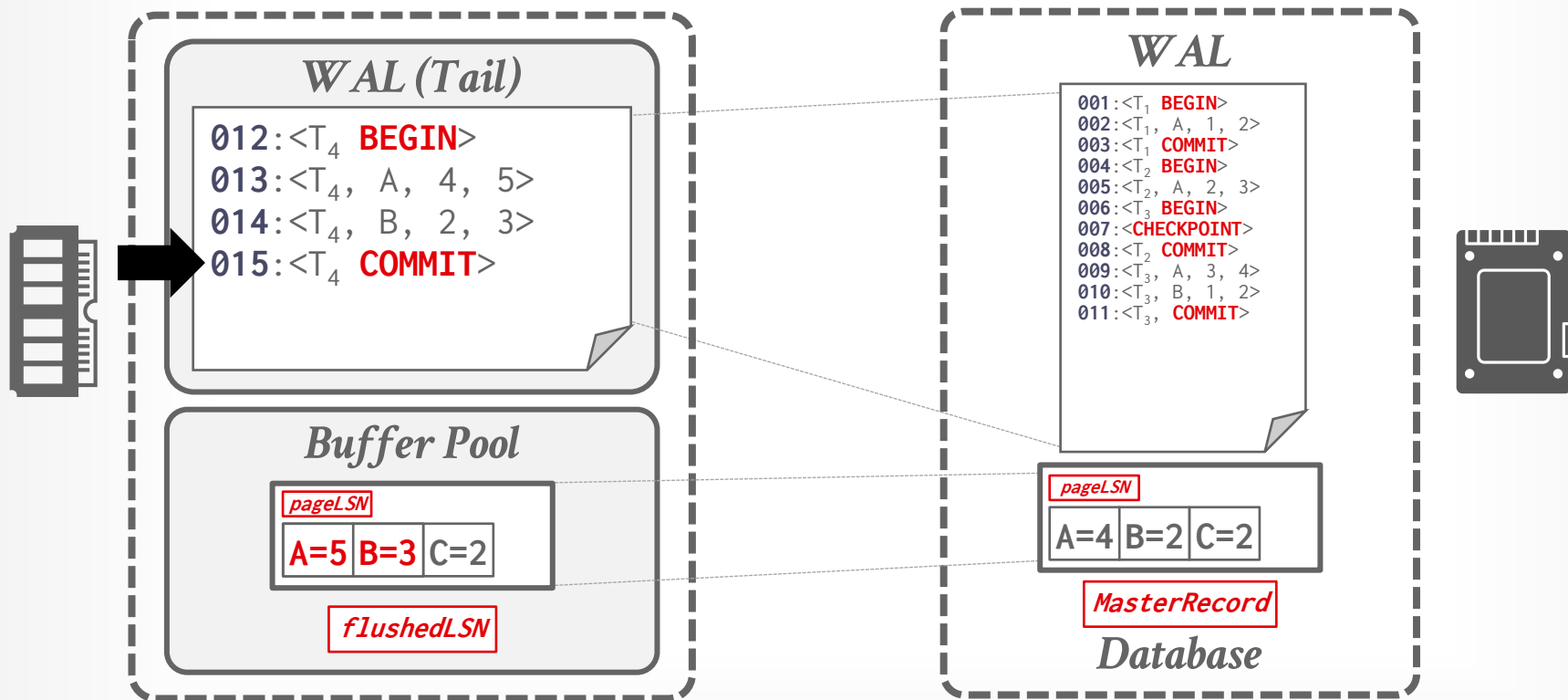
When commit succeeds, write **TXN-END** record to log.

- Indicates that no new log record for that txn will appear in the log ever again.
- DBMS does not need to flush these records immediately.

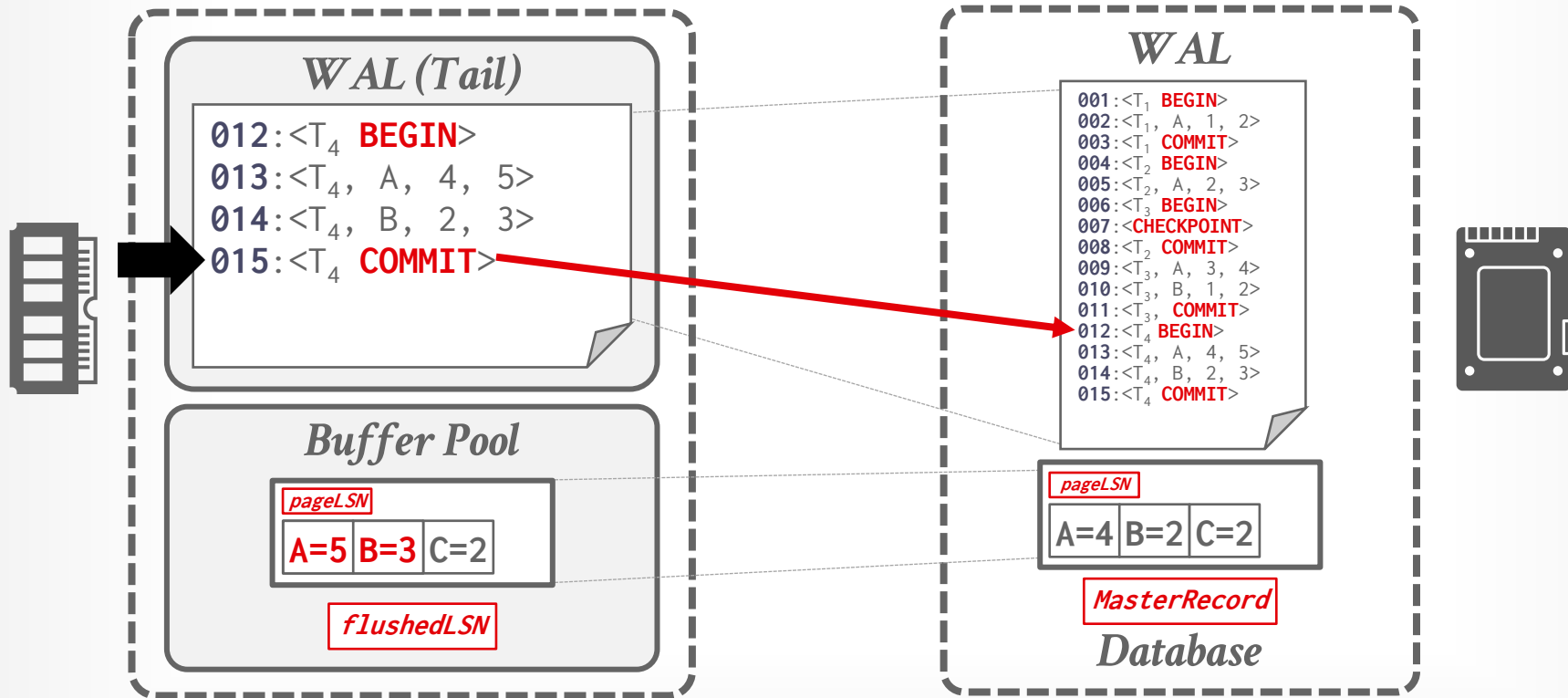
# TRANSACTION COMMIT



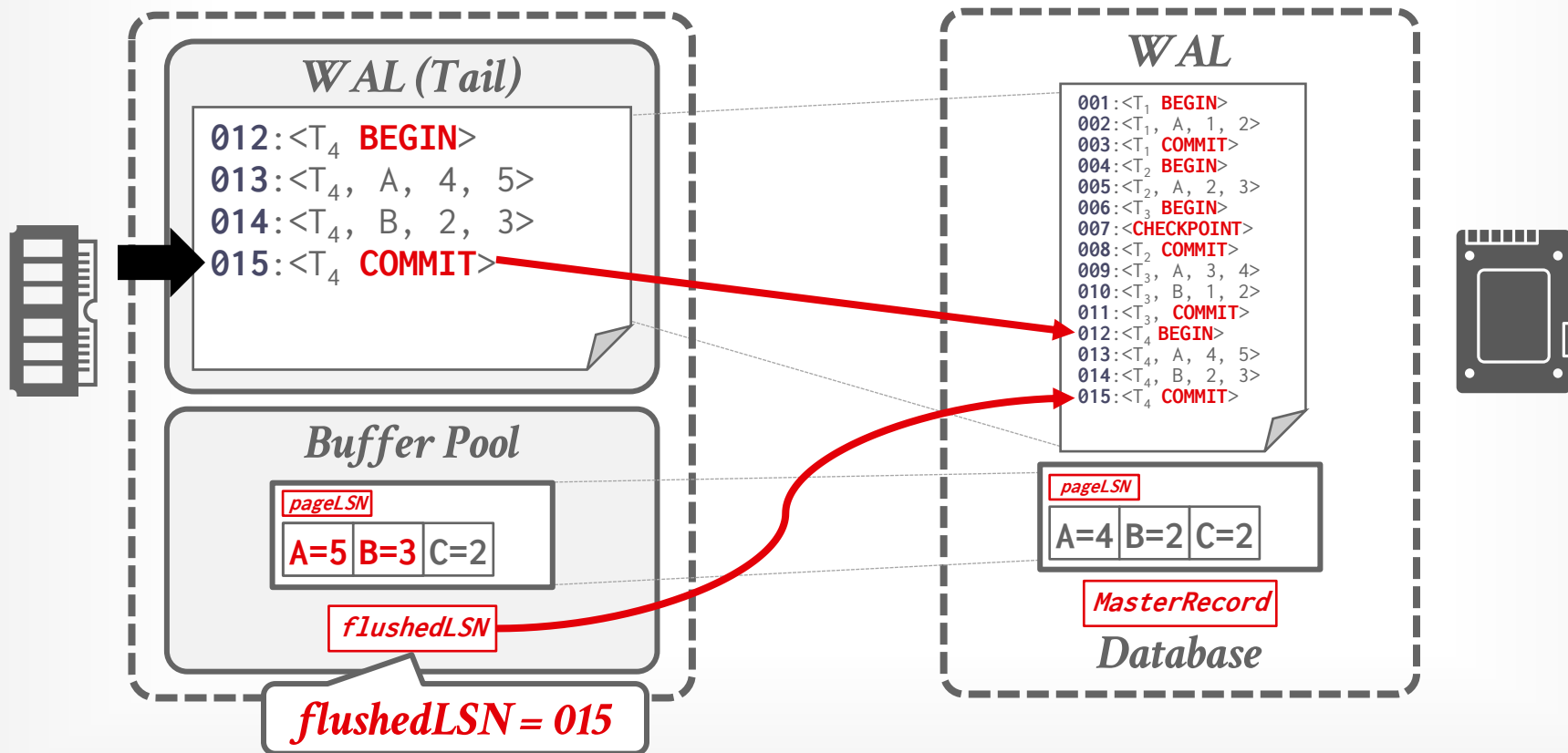
# TRANSACTION COMMIT



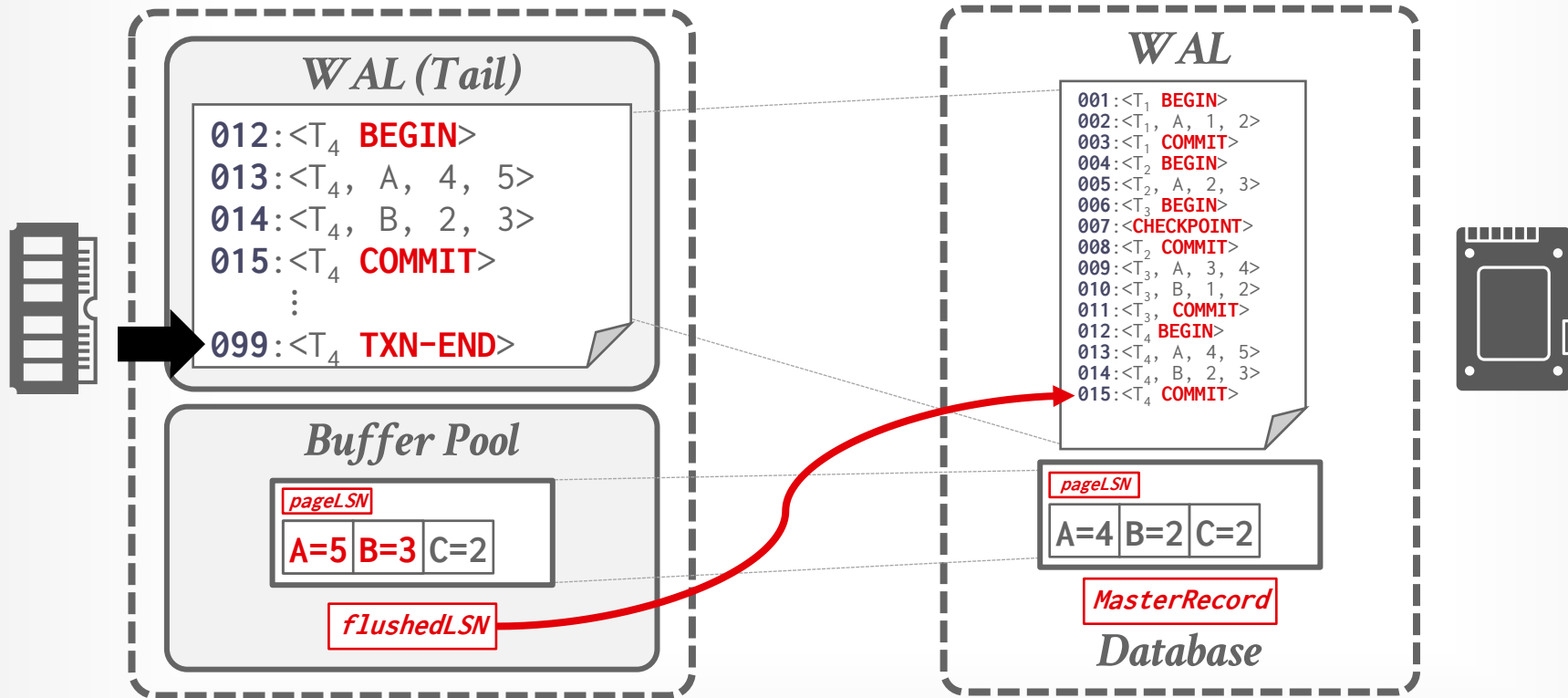
# TRANSACTION COMMIT



# TRANSACTION COMMIT

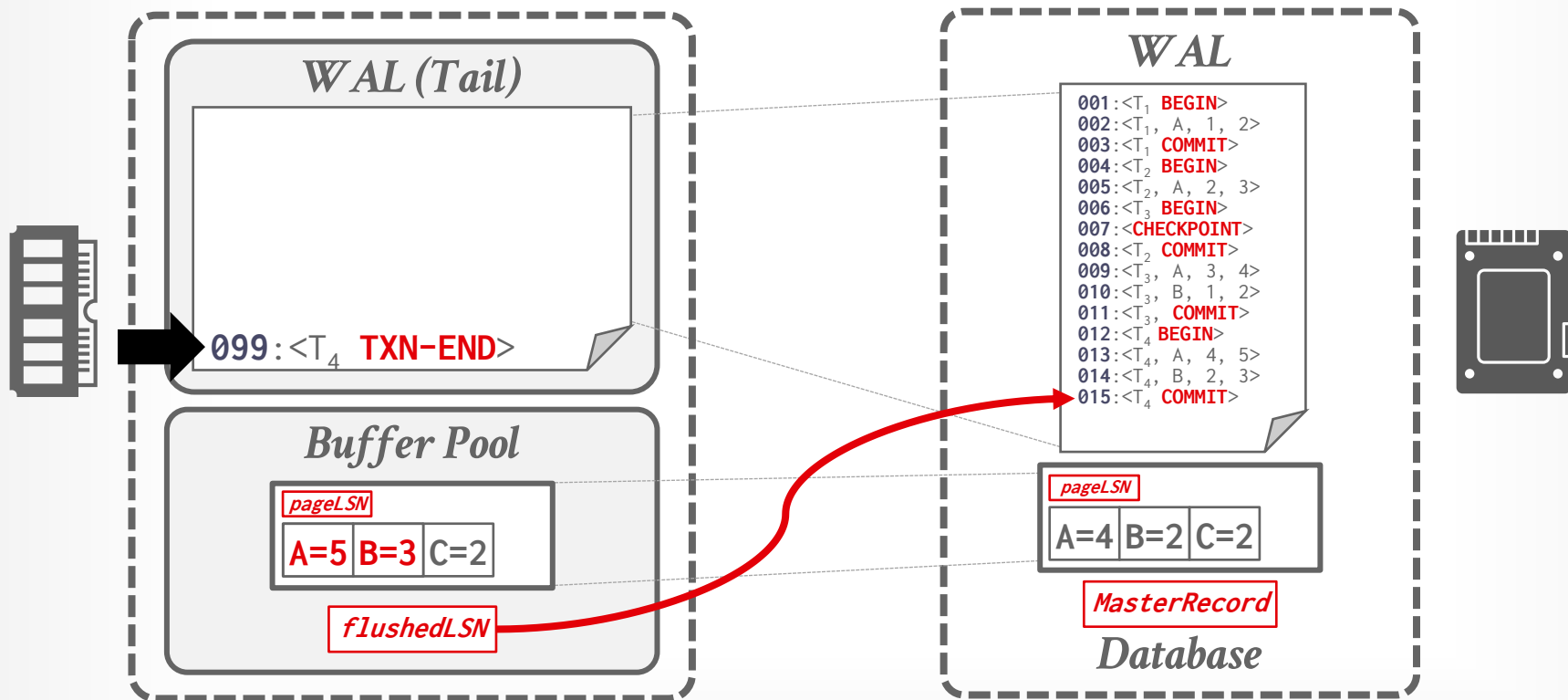


# TRANSACTION COMMIT





# TRANSACTION COMMIT



# TRANSACTION ABORT

---

Aborting a txn is a special case of the ARIES undo operation applied to only one txn.

We need to add another field to our log records:

- **prevLSN**: The previous **LSN** for the txn.
- This maintains a linked-list for each txn to make it easier to walk through its log records.

A txn does not release its locks until it has successfully reverted all its changes.

# TRANSACTION ABORT

---

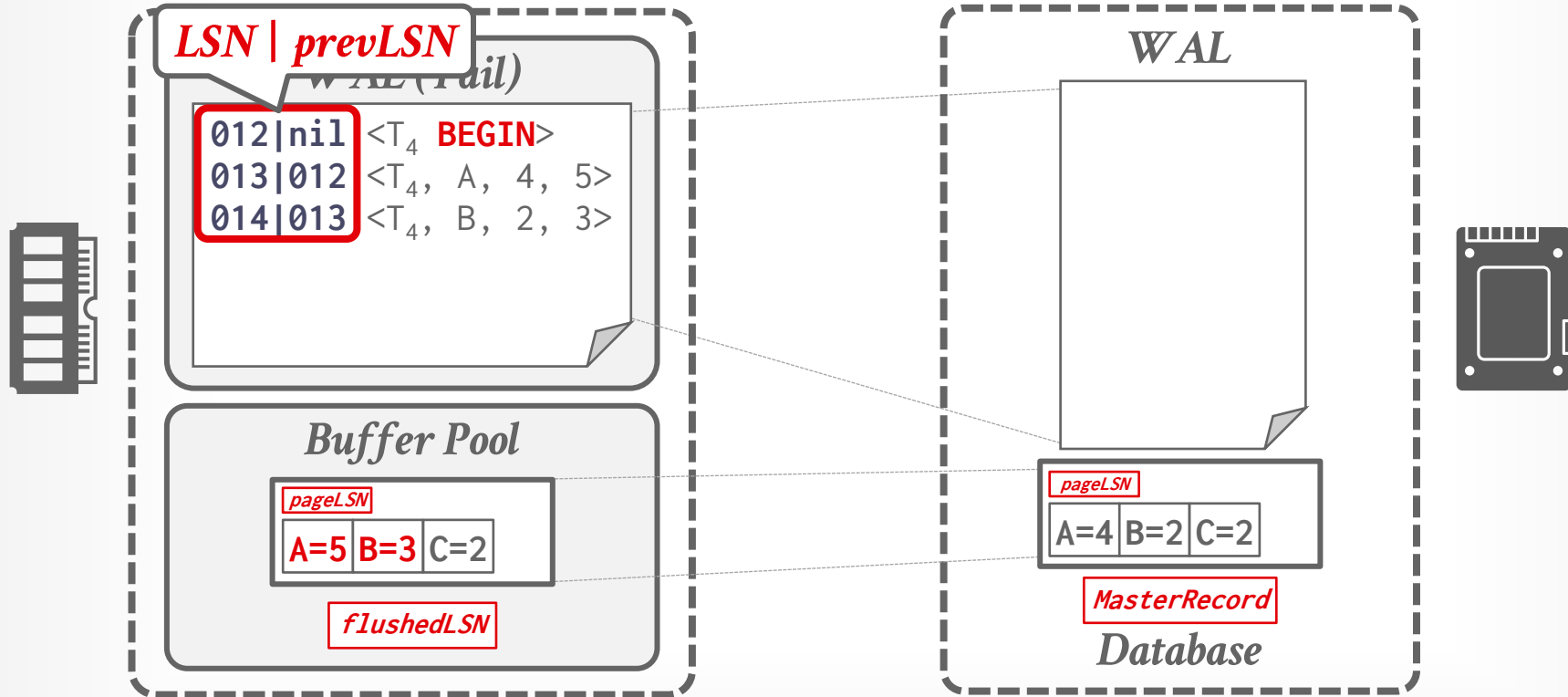
Aborting a txn is a special case of the ARIES undo operation applied to only one txn.

We need to add another field to our log records:

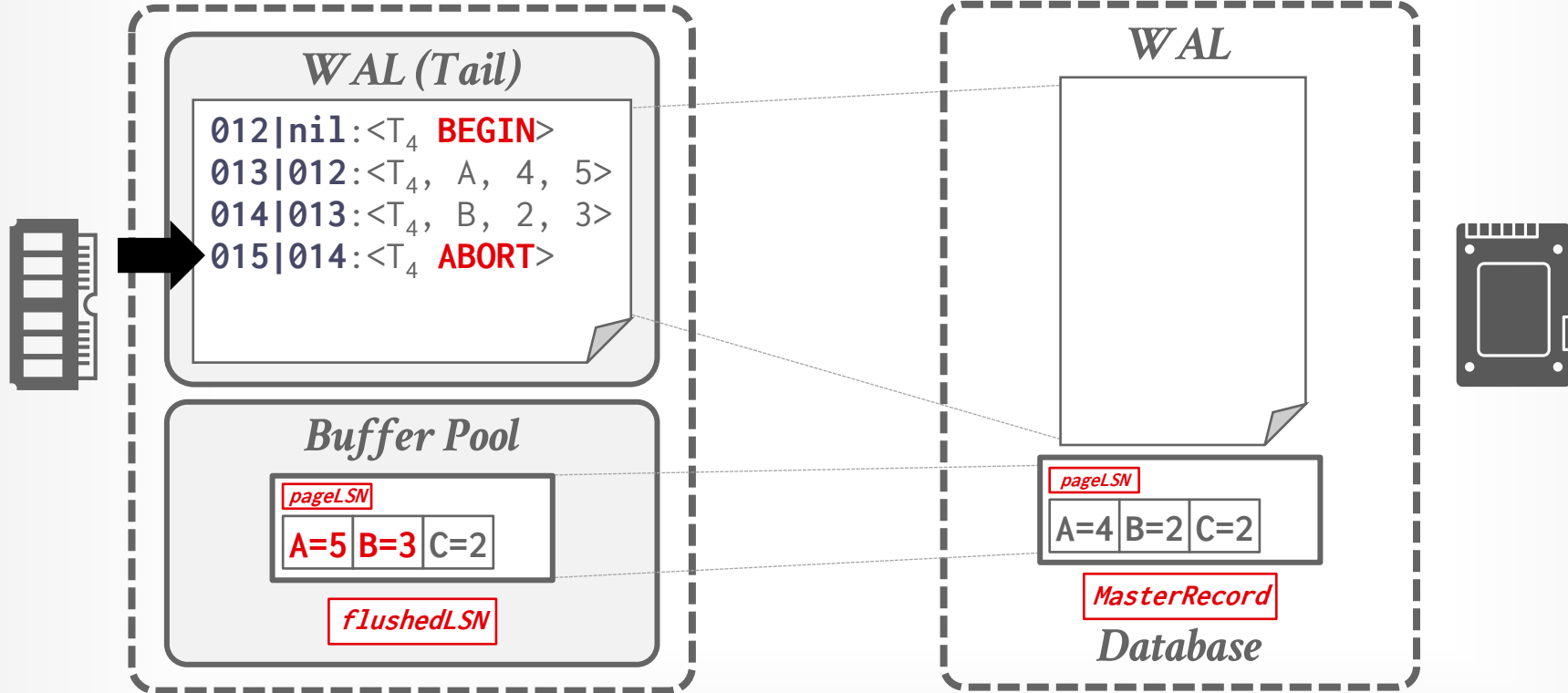
- **prevLSN**: The previous **LSN** for the txn.
- This maintains a linked-list for each txn to make it easier to walk through its log records.

A txn does not release its locks until it has successfully reverted all its changes.

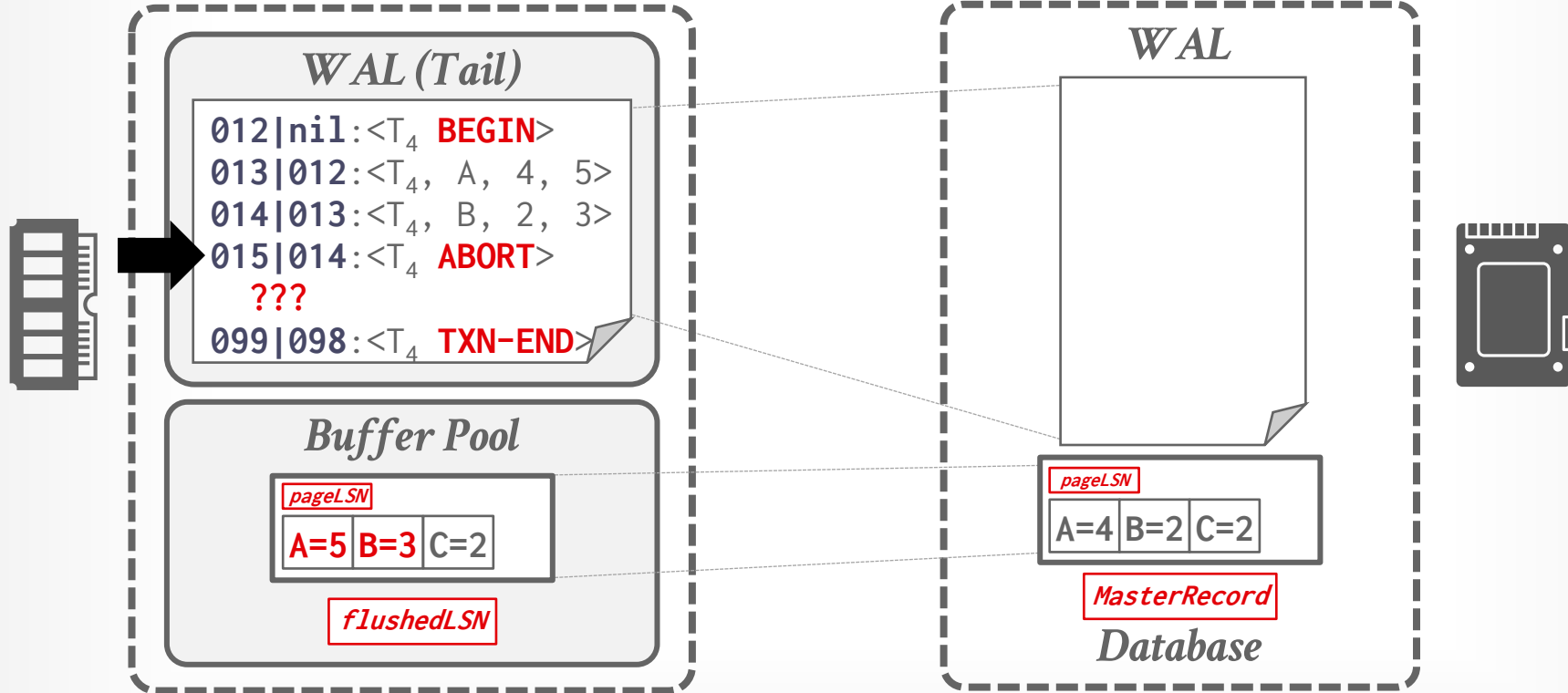
# TRANSACTION ABORT



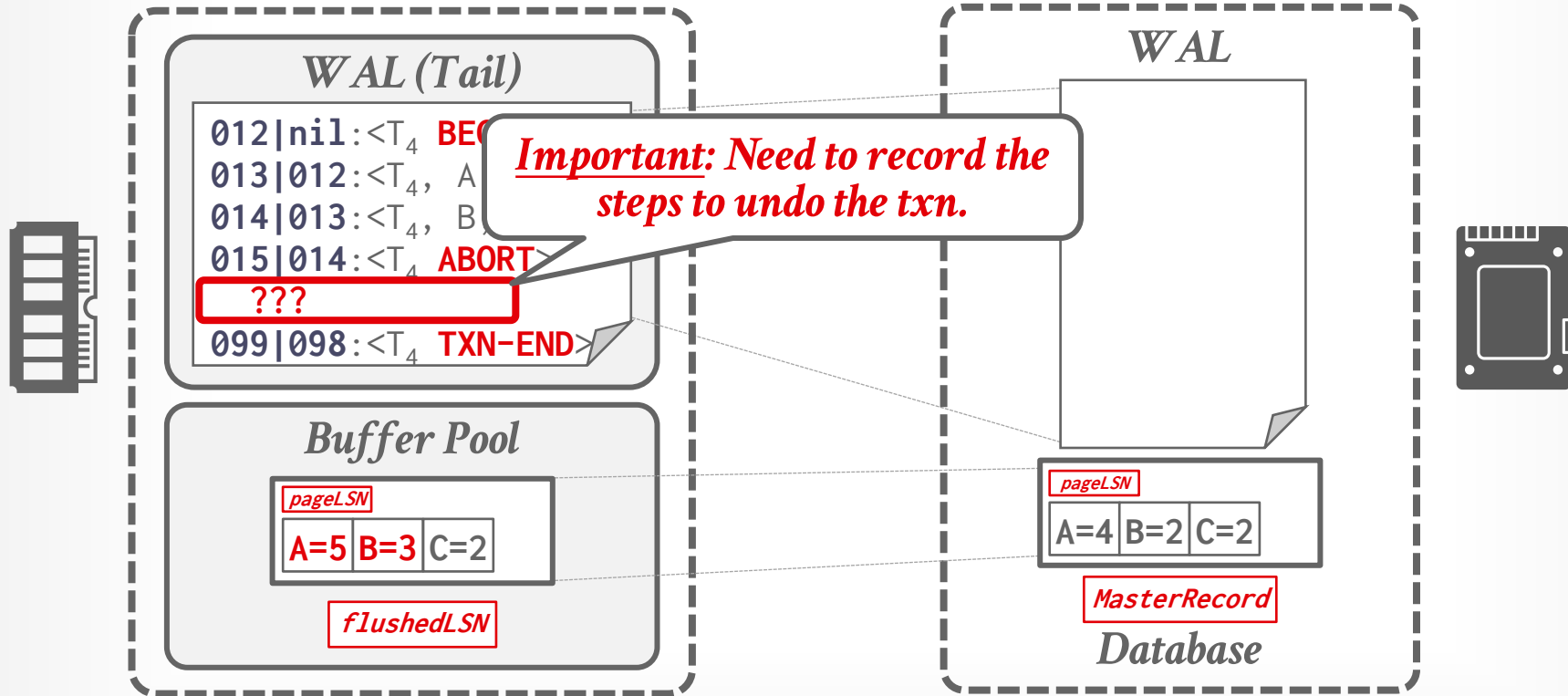
# TRANSACTION ABORT



# TRANSACTION ABORT



# TRANSACTION ABORT



# COMPENSATION LOG RECORDS

---

A compensation log record (CLR) describes the actions taken to reverse (i.e., undo) the changes of a previous update log record.

→ Each CLR contains all the fields of an update log record plus the **undoNextLSN** pointer (i.e., next-to-be-undone LSN).

DBMS adds CLR's to in-memory WAL buffer but it does not wait for them to be flushed before notifying the application the txn aborted.



# TRANSACTION ABORT: CLR EXAMPLE

TIME



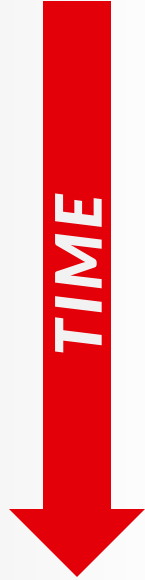
LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	$T_1$	BEGIN	–	–	–	–
002	001	$T_1$	UPDATE	A	30	40	–
003	002	$T_1$	UPDATE	B	10	24	–
⋮							
011	003	$T_1$	ABORT	–	–	–	–

# TRANSACTION ABORT: CLR EXAMPLE



LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	T <sub>1</sub>	BEGIN	–	–	–	–
002	001	T <sub>1</sub>	UPDATE	A	30	40	–
003	002	T <sub>1</sub>	UPDATE	B	10	24	–
⋮							
011	003	T <sub>1</sub>	ABORT	–	–	–	–
⋮							
026	011	T <sub>1</sub>	CLR-003	B	24	10	002

# TRANSACTION ABORT: CLR EXAMPLE



LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	T <sub>1</sub>	BEGIN	–	–	–	–
002	001	T <sub>1</sub>	UPDATE	A	30	40	–
003	002	T <sub>1</sub>	UPDATE	B	10	24	–
⋮							
011	003	T <sub>1</sub>	ABORT	–	–	–	–
⋮							
026	011	T <sub>1</sub>	CLR-003	B	24	10	002

# TRANSACTION ABORT: CLR EXAMPLE



LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	T <sub>1</sub>	BEGIN	–	–	–	–
002	001	T <sub>1</sub>	UPDATE	A	30	40	–
003	002	T <sub>1</sub>	UPDATE	B	10	24	–
⋮							
011	003	T <sub>1</sub>	ABORT	–	–	–	–
⋮							
026	011	T <sub>1</sub>	CLR-003	B	24	10	002

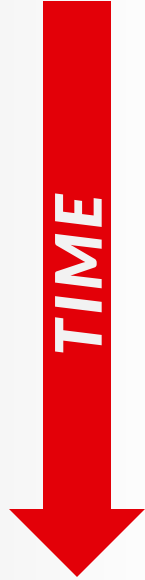
# TRANSACTION ABORT: CLR EXAMPLE



LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	T <sub>1</sub>	BEGIN	–	–	–	–
002	001	T <sub>1</sub>	UPDATE	A	30	40	–
003	002	T <sub>1</sub>	UPDATE	B	10	24	–
⋮							
011	003	T <sub>1</sub>	ABORT	–	–	–	–
⋮							
026	011	T <sub>1</sub>	CLR-003	B	24	10	002

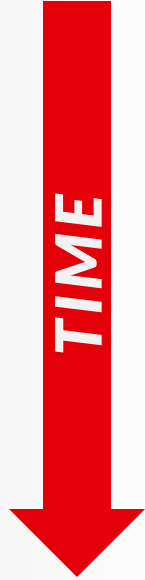
*The LSN of the next log record to be undone.*

# TRANSACTION ABORT: CLR EXAMPLE




LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	T <sub>1</sub>	BEGIN	–	–	–	–
002	001	T <sub>1</sub>	UPDATE	A	30	40	–
003	002	T <sub>1</sub>	UPDATE	B	10	24	–
⋮							
011	003	T <sub>1</sub>	ABORT	–	–	–	–
⋮							
026	011	T <sub>1</sub>	CLR-003	B	24	10	002
027	026	T <sub>1</sub>	CLR-002	A	40	30	001

# TRANSACTION ABORT: CLR EXAMPLE



LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	T <sub>1</sub>	BEGIN	–	–	–	–
002	001	T <sub>1</sub>	UPDATE	A	30	40	–
003	002	T <sub>1</sub>	UPDATE	B	10	24	–
⋮							
011	003	T <sub>1</sub>	ABORT	–	–	–	–
⋮							
026	011	T <sub>1</sub>	CLR-003	B	24	10	002
027	026	T <sub>1</sub>	CLR-002	A	40	30	001

# TRANSACTION ABORT: CLR EXAMPLE

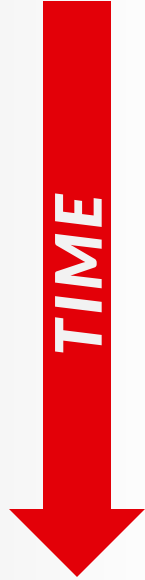


LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	T <sub>1</sub>	BEGIN	–	–	–	–
002	001	T <sub>1</sub>	UPDATE	A	30	40	–
003	002	T <sub>1</sub>	UPDATE	B	10	24	–
⋮							
011	003	T <sub>1</sub>	ABORT	–	–	–	–
⋮							
026	011	T <sub>1</sub>	CLR-003	B	24	10	002
027	026	T <sub>1</sub>	CLR-002	A	40	30	001

Diagram illustrating a transaction abort (CLR) example. The table shows log records for transaction T<sub>1</sub>. The 'Before' and 'After' values for the CLR records (026 and 027) are swapped, indicating a rollback. Red arrows and boxes highlight the swap: the 'Before' value 24 from record 026 is moved to the 'After' position of record 027, and the 'After' value 10 from record 026 is moved to the 'Before' position of record 027. The 'UndoNextLSN' for record 026 is 002, and for record 027 is 001.



# TRANSACTION ABORT: CLR EXAMPLE



LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	T <sub>1</sub>	BEGIN	–	–	–	–
002	001	T <sub>1</sub>	UPDATE	A	30	40	–
003	002	T <sub>1</sub>	UPDATE	B	10	24	–
⋮							
011	003	T <sub>1</sub>	ABORT	–	–	–	–
⋮							
026	011	T <sub>1</sub>	CLR-003	B	24	10	002
027	026	T <sub>1</sub>	CLR-002	A	40	30	001

# TRANSACTION ABORT: CLR EXAMPLE



LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	$T_1$	BEGIN	–	–	–	–
002	001	$T_1$	UPDATE	A	30	40	–
003	002	$T_1$	UPDATE	B	10	24	–
⋮							
011	003	$T_1$	ABORT	–	–	–	–
⋮							
026	011	$T_1$	CLR-003	B	24	10	002
027	026	$T_1$	CLR-002	A	40	30	001
028	027	$T_1$	TXN-END	–	–	–	nil

# TRANSACTION ABORT: CLR EXAMPLE



LSN	prevLSN	TxnId	Type	Object	Before	After	UndoNextLSN
001	nil	$T_1$	BEGIN	–	–	–	–
002	001	$T_1$	UPDATE	A	30	40	–
003	002	$T_1$	UPDATE	B	10	24	–
⋮							
011	003	$T_1$	ABORT	–	–	–	–
⋮							
026	011	$T_1$	CLR-003	B	24	10	002
027	026	$T_1$	CLR-002	A	40	30	001
028	027	$T_1$	TXN-END	–	–	–	nil

# ABORT ALGORITHM

---

First write an **ABORT** record to log for the txn.

Then analyze the txn's updates in reverse order. For each update record:

- Write a **CLR** entry to the log.
- Restore old value.

Lastly, write a **TXN-END** record and release locks.

Notice: CLRs never need to be undone.

# TODAY'S AGENDA

---



~~Log Sequence Numbers~~

~~Normal Commit & Abort Operations~~

Fuzzy Checkpointing

Recovery Algorithm

# NON-FUZZY CHECKPOINTS

---

The DBMS halts everything when it takes a checkpoint to ensure a consistent snapshot:

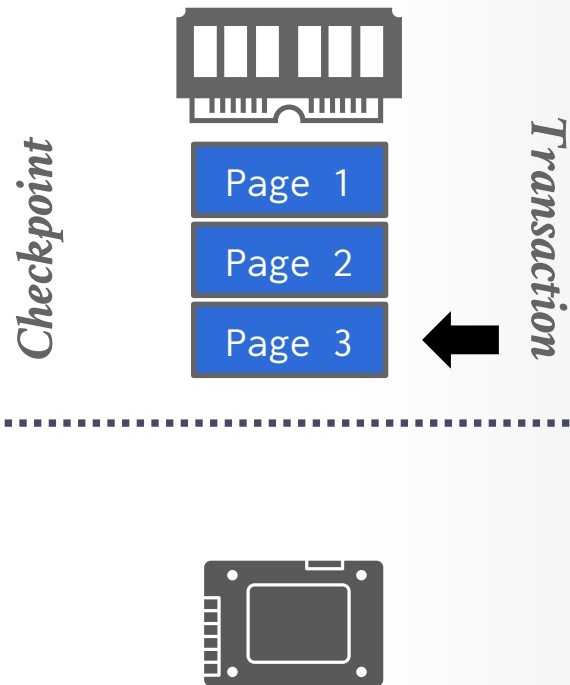
- Halt the start of any new txns.
- Wait until all active txns finish executing.
- Flushes dirty pages on disk.

This checkpoint implementation is bad for runtime performance, but it makes recovery easy.

# SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the DBMS takes the checkpoint.

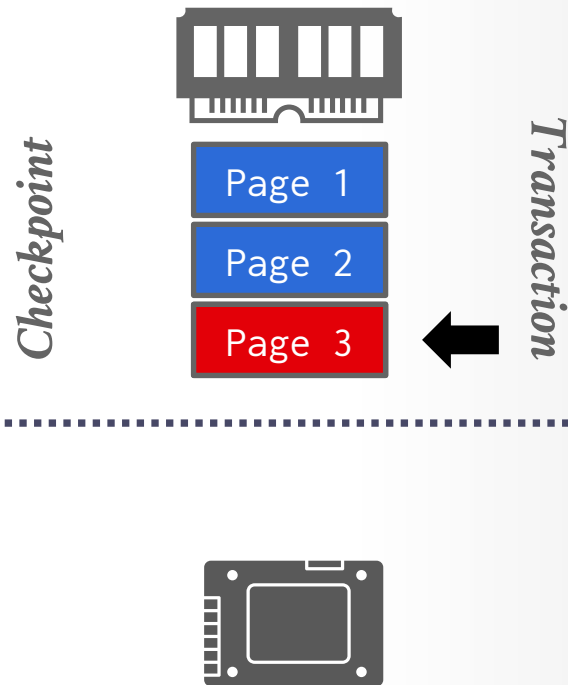
- Flushes dirty pages on disk.
- Block queries from acquiring write latch on pages.
- Do not wait until all txns finish before taking the checkpoint.



# SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the DBMS takes the checkpoint.

- Flushes dirty pages on disk.
- Block queries from acquiring write latch on pages.
- Do not wait until all txns finish before taking the checkpoint.

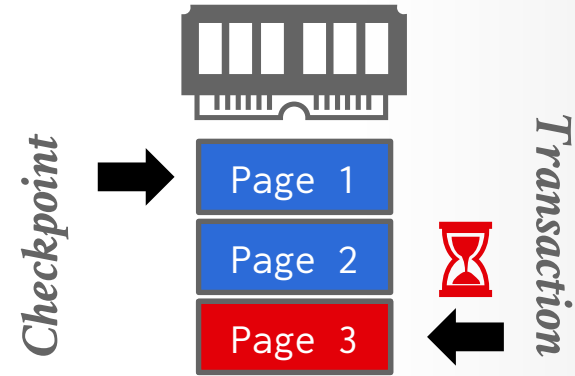




# SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the DBMS takes the checkpoint.

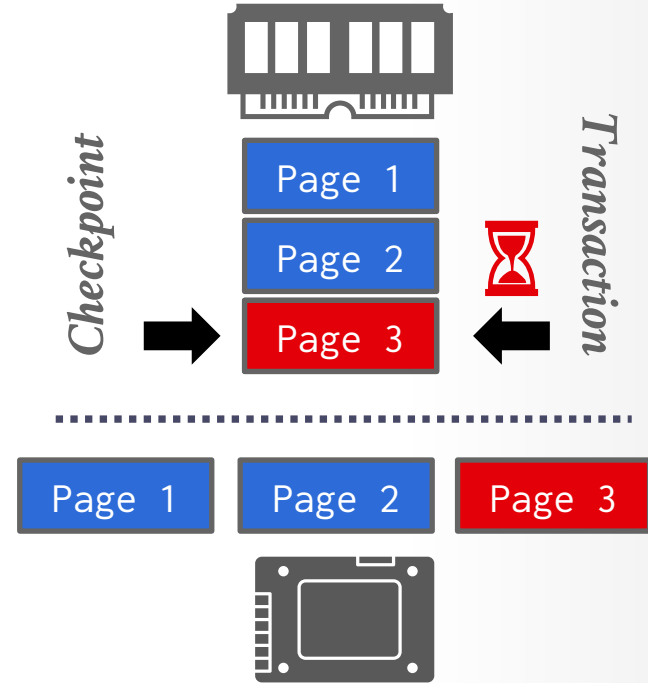
- Flushes dirty pages on disk.
- Block queries from acquiring write latch on pages.
- Do not wait until all txns finish before taking the checkpoint.



# SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the DBMS takes the checkpoint.

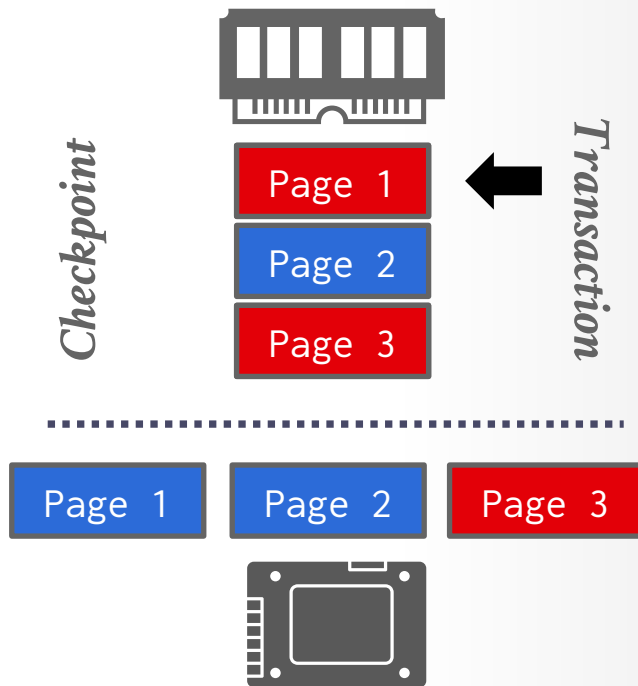
- Flushes dirty pages on disk.
- Block queries from acquiring write latch on pages.
- Do not wait until all txns finish before taking the checkpoint.



# SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the DBMS takes the checkpoint.

- Flushes dirty pages on disk.
- Block queries from acquiring write latch on pages.
- Do not wait until all txns finish before taking the checkpoint.



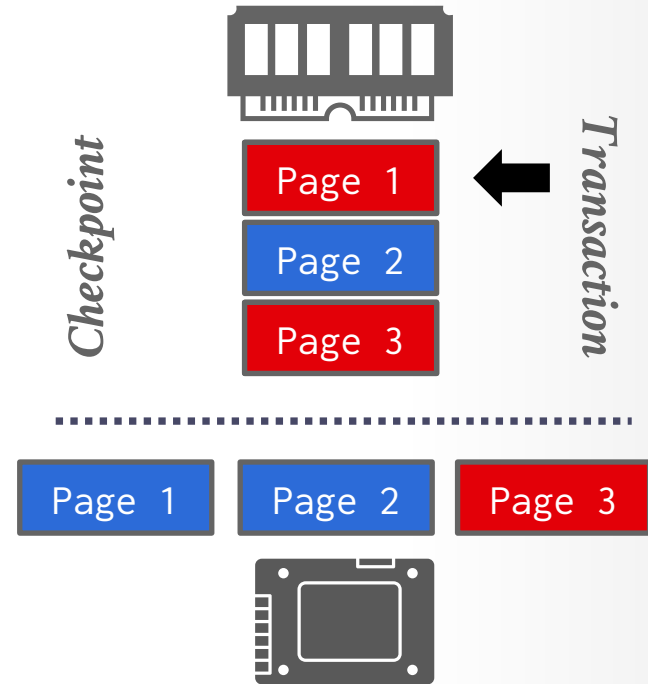
# SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the DBMS takes the checkpoint.

- Flushes dirty pages on disk.
- Block queries from acquiring write latch on pages.
- Do not wait until all txns finish before taking the checkpoint.

We must record internal state as of the beginning of the checkpoint.

- **Active Transaction Table (ATT)**
- **Dirty Page Table (DPT)**



# ACTIVE TRANSACTION TABLE (ATT)

---

One entry per currently active txn.

- **txnId**: Unique txn identifier.
- **status**: The current status mode of the txn.
- **lastLSN**: Most recent LSN created by the txn.

Remove a txn's ATT entry after appending its the **TXN-END** record to the in-memory WAL buffer.

## Txn Status Codes:

- Running (**R**)
- Committing (**C**)
- Candidate for Undo (**U**)

# DIRTY PAGE TABLE (DPT)

---

Separate data structure to track which pages in the buffer pool contain changes that the DBMS has not flushed to disk yet.

One entry per dirty page in the buffer pool:

- **recLSN**: The LSN of the oldest log record that modified the page since the last time the DBMS wrote the page to disk. This allows the DBMS to reason about whether the page was modified before or after the checkpoint started.

# SLIGHTLY BETTER CHECKPOINTS

At the first checkpoint, assuming  $P_{11}$  was flushed,  $T_2$  is still running and there is only one dirty page ( $P_{22}$ ).

```

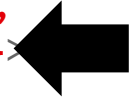
001:<T1 BEGIN>
002:<T2 BEGIN>
003:<T1, A→P11, 100, 120>
004:<T1 COMMIT>
005:<T2, C→P22, 100, 120>
006:<T1 TXN-END >
007:<CHECKPOINT
    ↗ ATT={T2},
    ↘ DPT={P22}>
008:<T3 BEGIN>
009:<T2, A→P11, 120, 130>
010:<T2 COMMIT>
011:<T3, B→P33, 200, 400>
012:<CHECKPOINT
    ↗ ATT={T2, T3},
    ↘ DPT={P11, P33}>
013:<T3, B→P33, 400, 600>
  
```

# SLIGHTLY BETTER CHECKPOINTS

At the first checkpoint, assuming  $P_{11}$  was flushed,  $T_2$  is still running and there is only one dirty page ( $P_{22}$ ).

```

001:<T1 BEGIN>
002:<T2 BEGIN>
003:<T1, A→P11, 100, 120>
004:<T1 COMMIT>
005:<T2, C→P22, 100, 120>
006:<T1 TXN-END >
007:<CHECKPOINT
      ↗ ATT={T2},
      ↗ DPT={P22}>
008:<T3 BEGIN>
009:<T2, A→P11, 120, 130>
010:<T2 COMMIT>
011:<T3, B→P33, 200, 400>
012:<CHECKPOINT
      ↗ ATT={T2, T3},
      ↗ DPT={P11, P33}>
013:<T3, B→P33, 400, 600>
  
```





# SLIGHTLY BETTER CHECKPOINTS

At the first checkpoint, assuming  $P_{11}$  was flushed,  $T_2$  is still running and there is only one dirty page ( $P_{22}$ ).

At the second checkpoint, assuming  $P_{22}$  was flushed,  $T_2$  and  $T_3$  are active and the dirty pages are ( $P_{11}$ ,  $P_{33}$ ).

```

001:<T1 BEGIN>
002:<T2 BEGIN>
003:<T1, A→P11, 100, 120>
004:<T1 COMMIT>
005:<T2, C→P22, 100, 120>
006:<T1 TXN-END >
007:<CHECKPOINT
      ↗ ATT={T2},
      ↗ DPT={P22}>
008:<T3 BEGIN>
009:<T2, A→P11, 120, 130>
010:<T2 COMMIT>
011:<T3, B→P33, 200, 400>
012:<CHECKPOINT
      ↗ ATT={T2, T3},
      ↗ DPT={P11, P33}>
013:<T3, B→P33, 400, 600>
  
```

# SLIGHTLY BETTER CHECKPOINTS


At the first checkpoint, assuming  $P_{11}$  was flushed,  $T_2$  is still running and there is only one dirty page ( $P_{22}$ ).

At the second checkpoint, assuming  $P_{22}$  was flushed,  $T_2$  and  $T_3$  are active and the dirty pages are ( $P_{11}$ ,  $P_{33}$ ).

This protocol still is not ideal because the DBMS stalls txns during the checkpoint...

```

001:<T1 BEGIN>
002:<T2 BEGIN>
003:<T1, A→P11, 100, 120>
004:<T1 COMMIT>
005:<T2, C→P22, 100, 120>
006:<T1 TXN-END >
007:<CHECKPOINT
      ↗ ATT={T2},
      ↗ DPT={P22}>
008:<T3 BEGIN>
009:<T2, A→P11, 120, 130>
010:<T2 COMMIT>
011:<T3, B→P33, 200, 400>
012:<CHECKPOINT
      ↗ ATT={T2, T3},
      ↗ DPT={P11, P33}>
013:<T3, B→P33, 400, 600>
  
```



# FUZZY CHECKPOINTS

---

A fuzzy checkpoint is where the DBMS allows active txns to continue to run while the system writes the log records for checkpoint.

→ No attempt to force dirty pages to disk.

New log records to track checkpoint boundaries:


- **CHECKPOINT-BEGIN**: Indicates start of checkpoint
- **CHECKPOINT-END**: Contains the state of the ATT + DPT from when the checkpoint started.

# FUZZY CHECKPOINTS

Assume the DBMS flushes **P<sub>11</sub>** before the first checkpoint starts.

Any txn that begins after the checkpoint starts is excluded from the ATT in the **CHECKPOINT-END** record.

The LSN of the **CHECKPOINT-BEGIN** record is written to the **MasterRecord** when it completes.



```


001:<T1 BEGIN>
002:<T2 BEGIN>
003:<T1, A→P11, 100, 120>
004:<T1 COMMIT>
005:<T2, C→P22, 100, 120>
006:<T1 TXN-END>
007:<CHECKPOINT-BEGIN>
008:<T3 BEGIN>
009:<T2, A→P11, 120, 130>
010:<CHECKPOINT-END
    ↪ ATT={T2},
    ↪ DPT={P22} >
011:<T2 COMMIT>
012:<T3, B→P33, 200, 400>
013:<CHECKPOINT-BEGIN>
014:<T3, B→P33, 10, 12>
015:<CHECKPOINT-END
    ↪ ATT={T2, T3},
    ↪ DPT={P11, P33}>
  
```

# FUZZY CHECKPOINTS

Assume the DBMS flushes **P<sub>11</sub>** before the first checkpoint starts.

Any txn that begins after the checkpoint starts is excluded from the ATT in the **CHECKPOINT-END** record.

The LSN of the **CHECKPOINT-BEGIN** record is written to the **MasterRecord** when it completes.



```

001:<T1 BEGIN>
002:<T2 BEGIN>
003:<T1, A→P11, 100, 120>
004:<T1 COMMIT>
005:<T2, C→P22, 100, 120>
006:<T1 TXN-END>
007:<CHECKPOINT-BEGIN>
008:<T3 BEGIN>
009:<T2, A→P11, 120, 130>
010:<CHECKPOINT-END
    ↪ ATT={T2},
    ↪ DPT={P22} >
011:<T2 COMMIT>
012:<T3, B→P33, 200, 400>
013:<CHECKPOINT-BEGIN>
014:<T3, B→P33, 10, 12>
015:<CHECKPOINT-END
    ↪ ATT={T2, T3},
    ↪ DPT={P11, P33}>
  
```

# FUZZY CHECKPOINTS

Assume the DBMS flushes  $P_{11}$  before the first checkpoint starts.

Any txn that begins after the checkpoint starts is excluded from the ATT in the **CHECKPOINT-END** record.

The LSN of the **CHECKPOINT-BEGIN** record is written to the **MasterRecord** when it completes.



```


001:<T1 BEGIN>
002:<T2 BEGIN>
003:<T1, A→P11, 100, 120>
004:<T1 COMMIT>
005:<T2, C→P22, 100, 120>
006:<T1 TXN-END>
007:<CHECKPOINT-BEGIN>
008:<T3 BEGIN>
009:<T2, A→P11, 120, 130>
010:<CHECKPOINT-END
    ↳ ATT={T2},
    ↳ DPT={P22} >
011:<T2 COMMIT>
012:<T3, B→P33, 200, 400>
013:<CHECKPOINT-BEGIN>
014:<T3, B→P33, 10, 12>
015:<CHECKPOINT-END
    ↳ ATT={T2, T3},
    ↳ DPT={P11, P33}>
  
```

# FUZZY CHECKPOINTS

Assume the DBMS flushes **P<sub>11</sub>** before the first checkpoint starts.

Any txn that begins after the checkpoint starts is excluded from the ATT in the **CHECKPOINT-END** record.

The LSN of the **CHECKPOINT-BEGIN** record is written to the **MasterRecord** when it completes.



```

001:<T1 BEGIN>
002:<T2 BEGIN>
003:<T1, A→P11, 100, 120>
004:<T1 COMMIT>
005:<T2, C→P22, 100, 120>
006:<T1 TXN-END>
007:<CHECKPOINT-BEGIN>
008:<T3 BEGIN>
009:<T2, A→P11, 120, 130>
010:<CHECKPOINT-END
    ↳ ATT={T2},
    ↳ DPT={P22}>
011:<T2 COMMIT>
012:<T3, B→P33, 200, 400>
013:<CHECKPOINT-BEGIN>
014:<T3, B→P33, 10, 12>
015:<CHECKPOINT-END
    ↳ ATT={T2, T3},
    ↳ DPT={P11, P33}>
  
```

# ARIES: RECOVERY PHASES

---

## Phase #1: Analysis

- Examine the WAL in forward direction starting at MasterRecord 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 (even txns that will abort).

## Phase #3: Undo

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



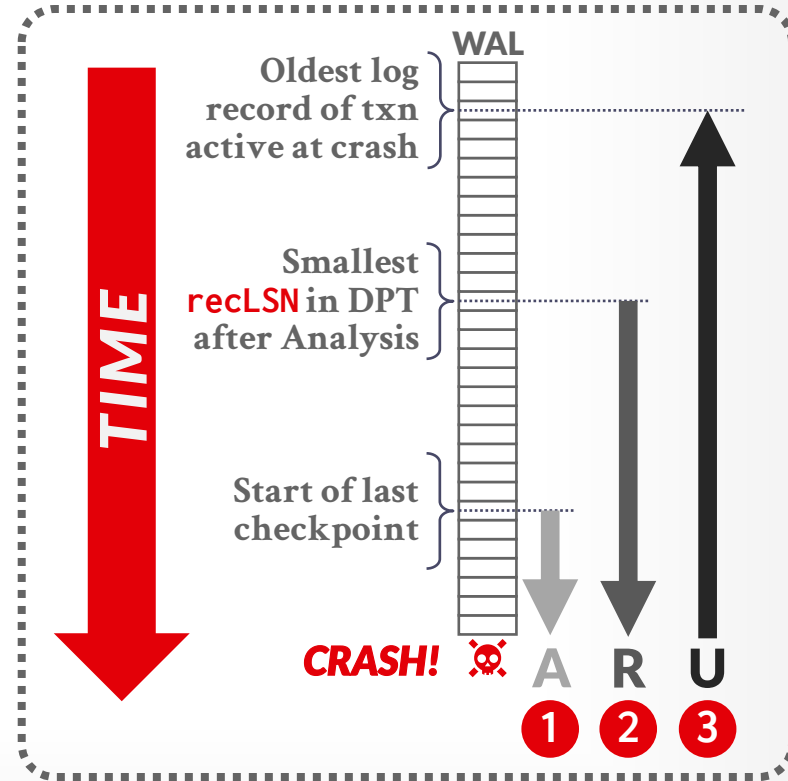
# ARIES: OVERVIEW

Start from last **BEGIN-CHECKPOINT** found via **MasterRecord**.

**Analysis:** Figure out which txns committed or failed since checkpoint.

**Redo:** Repeat all actions.

**Undo:** Reverse effects of failed txns.



# ANALYSIS PHASE

---

Scan log forward from the beginning of the last successful checkpoint to the end of the log and construct the **ATT**.

If the DBMS finds a **TXN-END** record, remove its corresponding txn from **ATT**.

All other records:

- If txn not in **ATT**, add it with undo candidate status (**U**).
- On commit, change the txn's status in **ATT** to commit (**C**).

For update log records:

- If page **P<sub>x</sub>** not in **DPT**, add **P<sub>x</sub>** to **DPT**, set its **recLSN=LSN**.

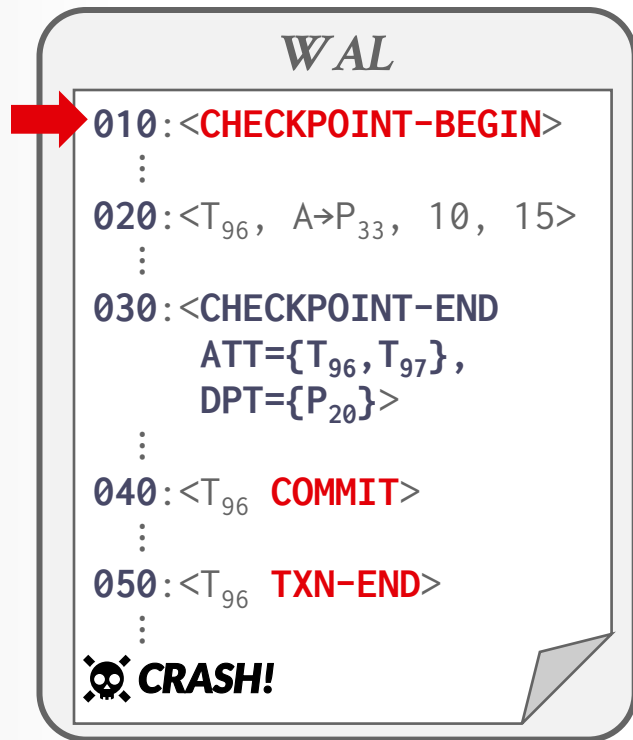
# ANALYSIS PHASE

---

At end of the Analysis Phase:

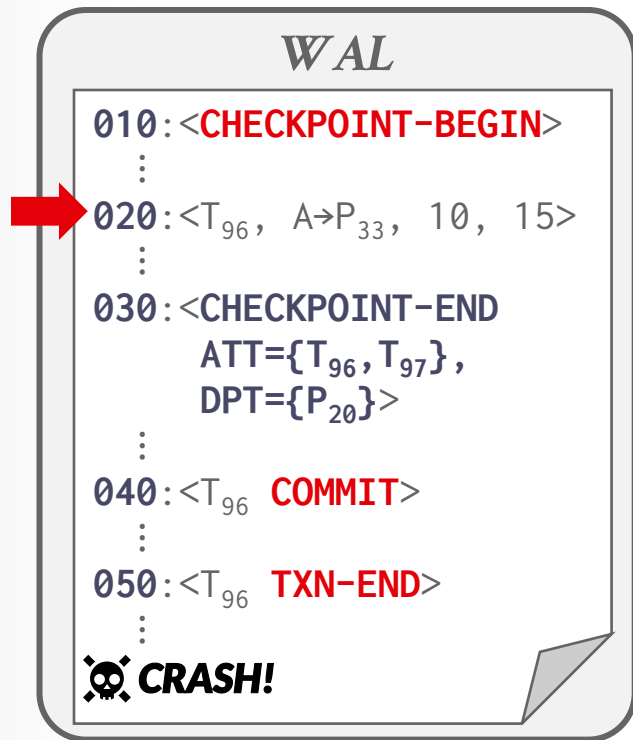
- **ATT** identifies which txns were active at time of crash.
- **DPT** identifies which dirty pages might not have made it to disk.

# ANALYSIS PHASE EXAMPLE



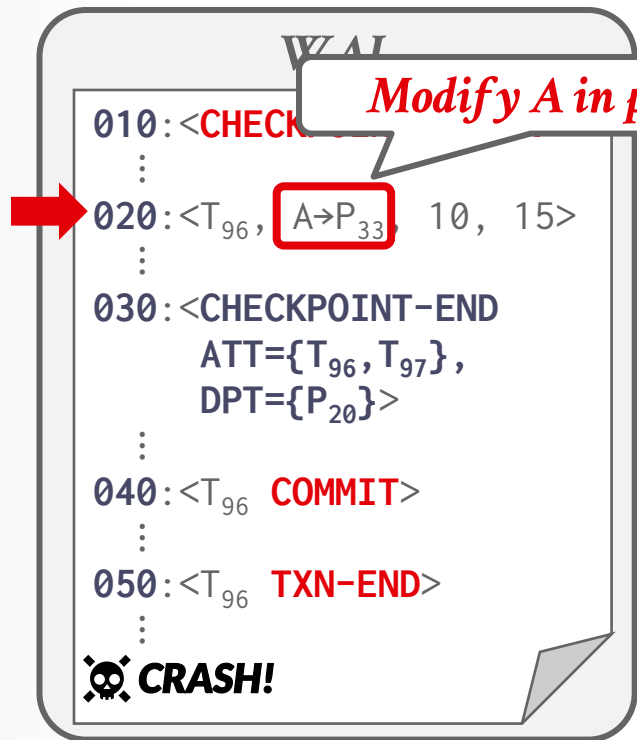
LSN	ATT	DPT
010		
020		
030		
040		
050		

# ANALYSIS PHASE EXAMPLE



LSN	ATT	DPT
010	<i>(TxnId, Status)</i>	
020		
030	(T <sub>96</sub> , U)	
040		
050		

# ANALYSIS PHASE EXAMPLE



LSN	ATT	DPT
010		
020	(T <sub>96</sub> , U)	(P <sub>33</sub> , 020)
030		
040		
050		

*(PageId, RecLSN)*

# ANALYSIS PHASE EXAMPLE

## WAL

010: <CHECKPOINT-BEGIN>

⋮

020: <T<sub>96</sub>, A→P<sub>33</sub>, 10, 15>

⋮

030: <CHECKPOINT-END  
ATT={T<sub>96</sub>, T<sub>97</sub>},  
DPT={P<sub>20</sub>}>

⋮

040: <T<sub>96</sub> COMMIT>

⋮

050: <T<sub>96</sub> TXN-END>

⋮

💣 CRASH!

LSN	ATT	DPT
010		
020	(T <sub>96</sub> , U)	(P <sub>33</sub> , 020)
030	(T <sub>96</sub> , U), (T <sub>97</sub> , U)	(P <sub>33</sub> , 020), (P <sub>20</sub> , 008)
040		
050		

# ANALYSIS PHASE EXAMPLE

## WAL

010: <CHECKPOINT-BEGIN>

⋮

020: <T<sub>96</sub>, A→P<sub>33</sub>, 10, 15>

⋮

030: <CHECKPOINT-END

ATT={T<sub>96</sub>, T<sub>97</sub>} ,

DPT={P<sub>20</sub>}>

⋮

040: <T<sub>96</sub> COMMIT>

⋮

050: <T<sub>96</sub> TXN-END>

⋮

💥 CRASH!

LSN	ATT	DPT
010		
020	(T <sub>96</sub> , U)	(P <sub>33</sub> , 020)
030	(T <sub>96</sub> , U), (T <sub>97</sub> , U)	(P <sub>33</sub> , 020), (P <sub>20</sub> , 008)
040	(T <sub>96</sub> , C), (T <sub>97</sub> , U)	(P <sub>33</sub> , 020), (P <sub>20</sub> , 008)
050		



# ANALYSIS PHASE EXAMPLE

## WAL

010: <CHECKPOINT-BEGIN>

⋮

020: <T<sub>96</sub>, A→P<sub>33</sub>, 10, 15>

⋮

030: <CHECKPOINT-END  
ATT={T<sub>96</sub>, T<sub>97</sub>},  
DPT={P<sub>20</sub>}>

⋮

040: <T<sub>96</sub> COMMIT>

⋮

050: <T<sub>96</sub> TXN-END>

⋮

**💣 CRASH!**

LSN	ATT	DPT
010		
020	(T <sub>96</sub> , <b>U</b> )	(P <sub>33</sub> , 020)
030	(T <sub>96</sub> , <b>U</b> ), (T <sub>97</sub> , <b>U</b> )	(P <sub>33</sub> , 020), (P <sub>20</sub> , 008)
040	(T <sub>96</sub> , <b>C</b> ), (T <sub>97</sub> , <b>U</b> )	(P <sub>33</sub> , 020), (P <sub>20</sub> , 008)
050	(T <sub>97</sub> , <b>U</b> )	(P <sub>33</sub> , 020), (P <sub>20</sub> , 008)

# REDO PHASE

---

The goal is to repeat history to reconstruct the database state at the moment of the crash:

→ Reapply all updates (even aborted txns!) and redo CLRs.

There are techniques that allow the DBMS to avoid unnecessary reads/writes, but we will ignore them in this lecture...

# REDO PHASE

---

Scan forward from the log record containing smallest **recLSN** in **DPT**.

For each update log record or CLR with a **LSN**, redo action unless one of the following conditions are true:

- Target page is not in **DPT**.
- Target page is in **DPT**, but that log record's **LSN** is less than the page's **recLSN**. DBMS does not need to retrieve page from disk for this check.
- Target page is in **DPT**, but that log record's **LSN**  $\leq$  **pageLSN**. DBMS must retrieve the page from disk for this check.

# REDO PHASE

---

To redo an action:

- Reapply logged update.
- Set **pageLSN** to log record's **LSN**.
- No additional logging, no forced flushes!

At the end of Redo Phase, write **TXN-END** log records for all txns with commit status (**C**) and remove them from the **ATT**.

# UNDO PHASE

---

Undo all txns that were active at the time of crash and therefore will never commit.

→ These are all the txns with undo candidate status (**U**) in the **ATT** after the Analysis Phase.

Process them in reverse LSN order using the **lastLSN** to speed up traversal.

→ At each step, pick the largest **lastLSN** across all transactions in the **ATT**.

→ Traverse **lastLSNs** in the same order, but in reverse, for how the updates happened originally.

Write a **CLR** for every modification.

# FULL EXAMPLE

011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> ABORT>

# FULL EXAMPLE

011:<CHECKPOINT-BEGIN>

012:<CHECKPOINT-END>

013:<T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014:<T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015:<T<sub>1</sub> ABORT>

016:<T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

# FULL EXAMPLE

prevLSNs

011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> ABORT>

016: <T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> TXN-END>



# FULL EXAMPLE

011:<CHECKPOINT-BEGIN>

012:<CHECKPOINT-END>

013:<T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014:<T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015:<T<sub>1</sub> ABORT>

016:<T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

017:<T<sub>1</sub> TXN-END>

# FULL EXAMPLE

011:<CHECKPOINT-BEGIN>

012:<CHECKPOINT-END>

013:<T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014:<T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015:<T<sub>1</sub> ABORT>

016:<T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

017:<T<sub>1</sub> TXN-END>

018:<T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

# FULL EXAMPLE

011:<CHECKPOINT-BEGIN>

012:<CHECKPOINT-END>

013:<T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014:<T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015:<T<sub>1</sub> ABORT>

016:<T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

017:<T<sub>1</sub> TXN-END>

018:<T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019:<T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

# FULL EXAMPLE

011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> ABORT>

016: <T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

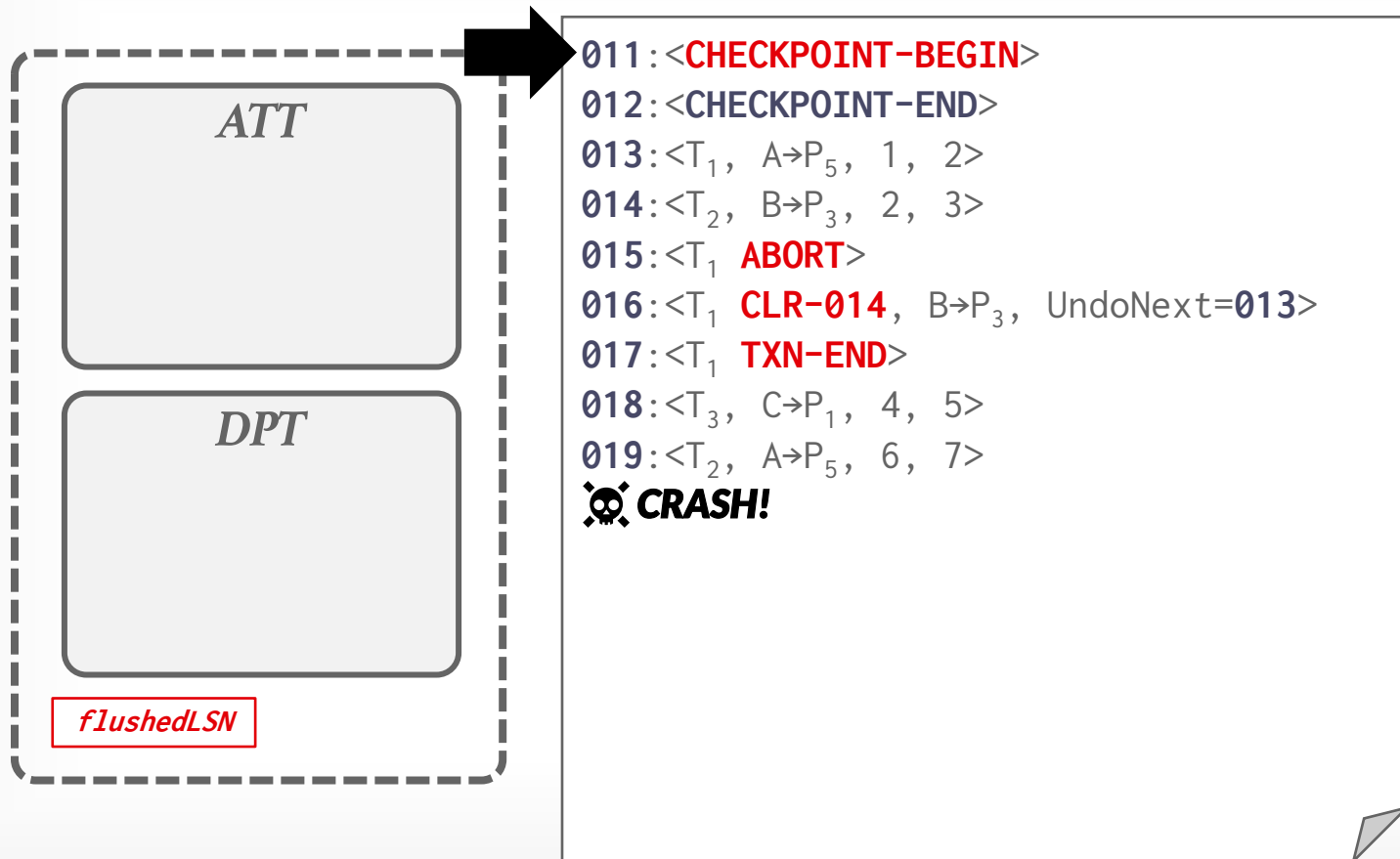
017: <T<sub>1</sub> TXN-END>

018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

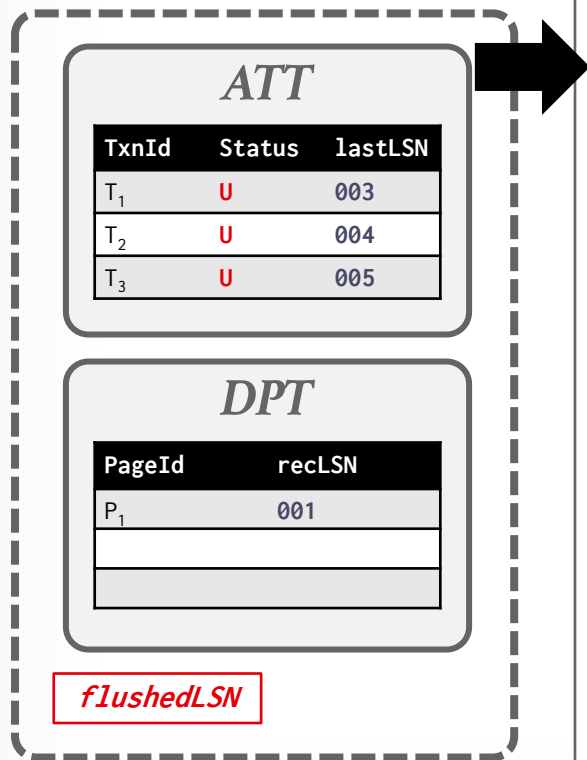
019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

💀 CRASH!

# FULL EXAMPLE



# FULL EXAMPLE



011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> **ABORT**>

016: <T<sub>1</sub> **CLR-014**, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> **TXN-END**>

018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

💀 **CRASH!**

# FULL EXAMPLE



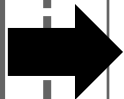
*ATT*

TxnId	Status	lastLSN
T <sub>1</sub>	U	017
T <sub>2</sub>	U	019
T <sub>3</sub>	U	018

*DPT*

PageId	recLSN
P <sub>1</sub>	001
P <sub>3</sub>	014
P <sub>5</sub>	013

*flushedLSN*



011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> ABORT>

016: <T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> TXN-END>

018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

💀 **CRASH!**

# FULL EXAMPLE



*ATT*

TxnId	Status	lastLSN
T <sub>1</sub>	U	017
T <sub>2</sub>	U	019
T <sub>3</sub>	U	018

*DPT*

PageId	recLSN
P <sub>1</sub>	001
P <sub>3</sub>	014
P <sub>5</sub>	013

*flushedLSN*

011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> ABORT>

016: <T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> TXN-END>

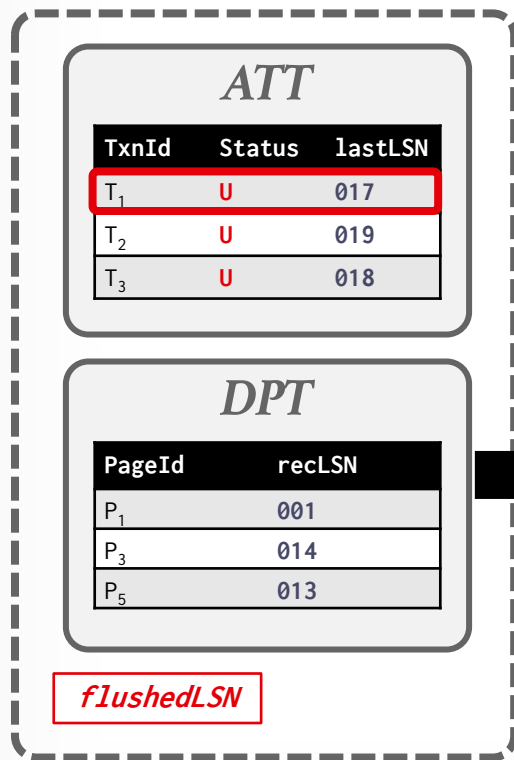
018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

💀 **CRASH!**



# FULL EXAMPLE



011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> ABORT>

016: <T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

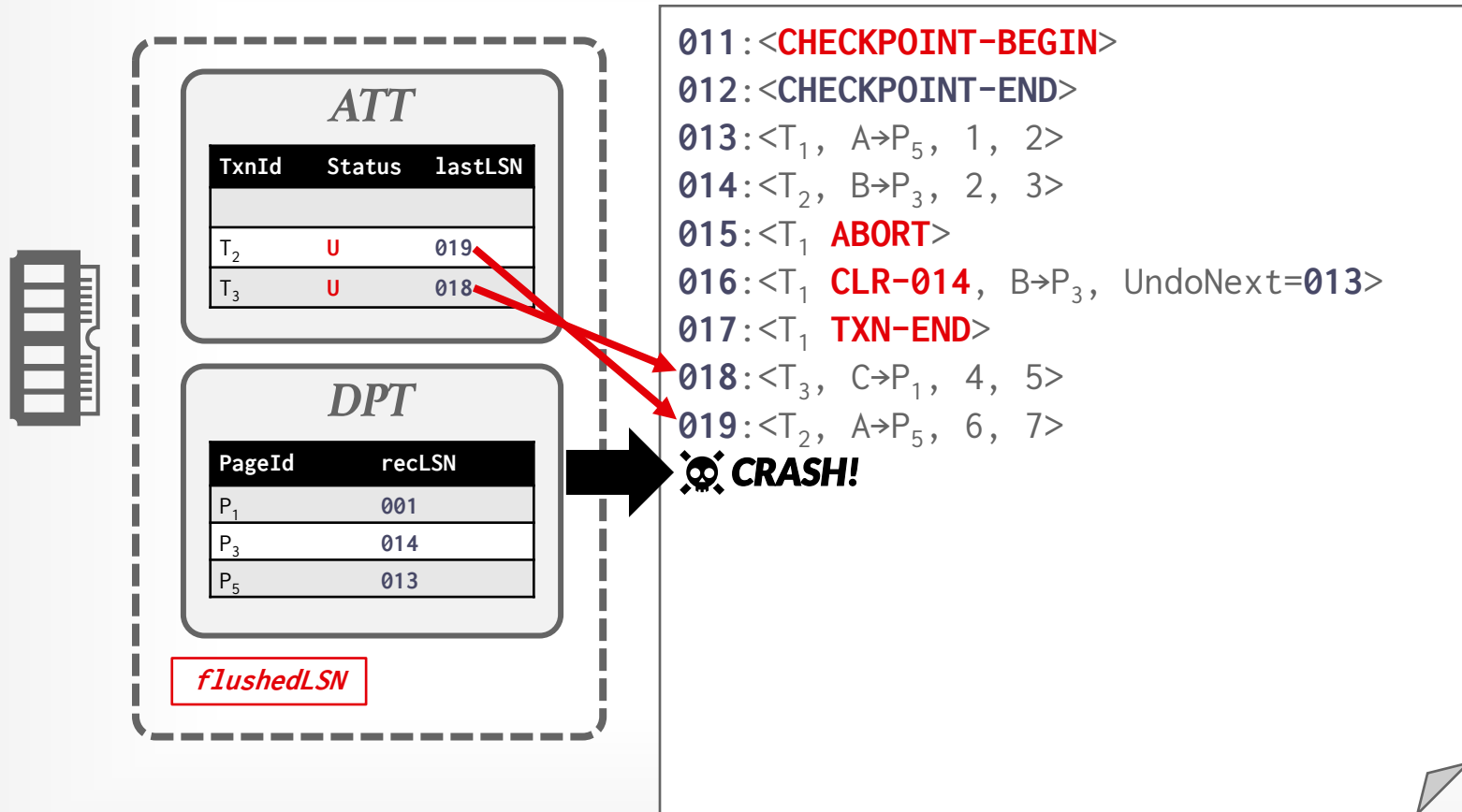
017: <T<sub>1</sub> TXN-END>

018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

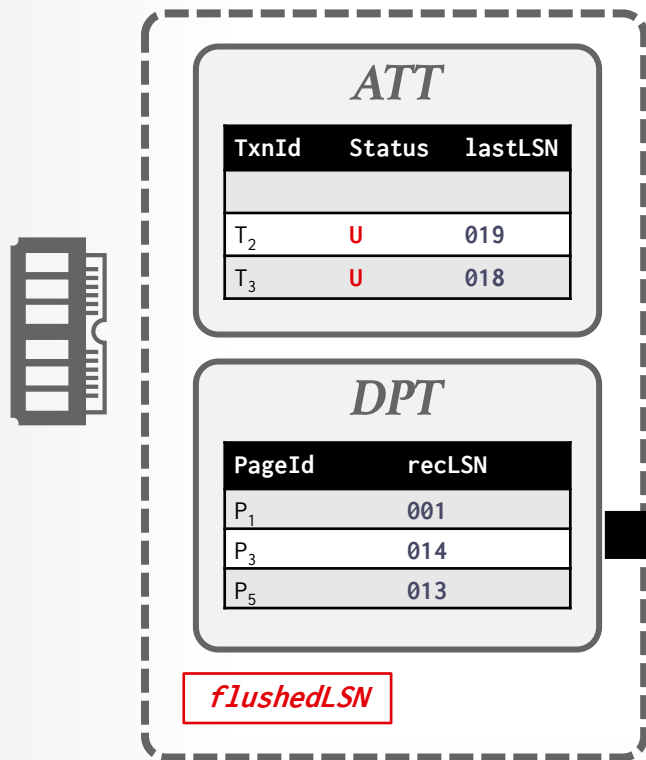
019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

💀 **CRASH!**

# FULL EXAMPLE



# FULL EXAMPLE



011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> ABORT>

016: <T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> TXN-END>

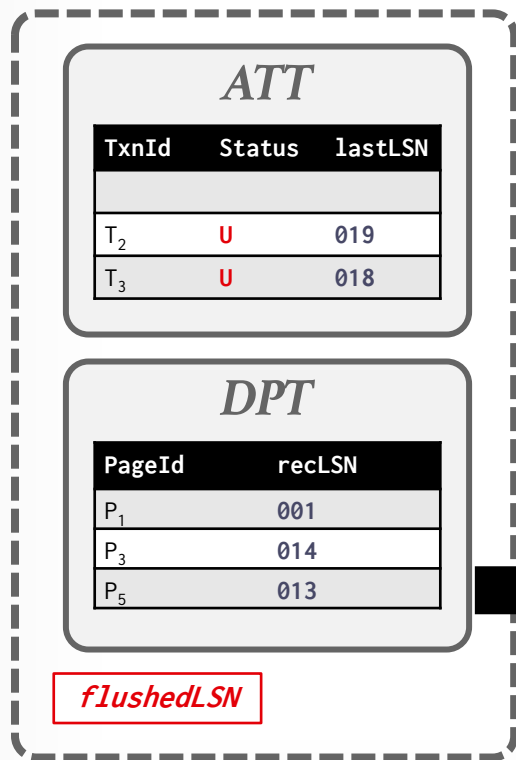
018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

💀 **CRASH!**

101: <T<sub>2</sub> CLR-019, D→P<sub>5</sub>, UndoNext=014>

# FULL EXAMPLE



011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> ABORT>

016: <T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> TXN-END>

018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

💀 **CRASH!**

101: <T<sub>2</sub> CLR-019, D→P<sub>5</sub>, UndoNext=014>

102: <T<sub>3</sub> CLR-018, C→P<sub>1</sub>, UndoNext=005>

005: <T<sub>3</sub> BEGIN>

# FULL EXAMPLE



*ATT*

TxnId	Status	lastLSN
T <sub>2</sub>	U	019
T <sub>3</sub>	U	018

*DPT*

PageId	recLSN
P <sub>1</sub>	001
P <sub>3</sub>	014
P <sub>5</sub>	013

*flushedLSN*

011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> ABORT>

016: <T<sub>1</sub> CLR-014, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> TXN-END>

018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019: <T<sub>2</sub>, A→

☠ CRASH!

*Flush dirty pages  
+ WAL to disk!*

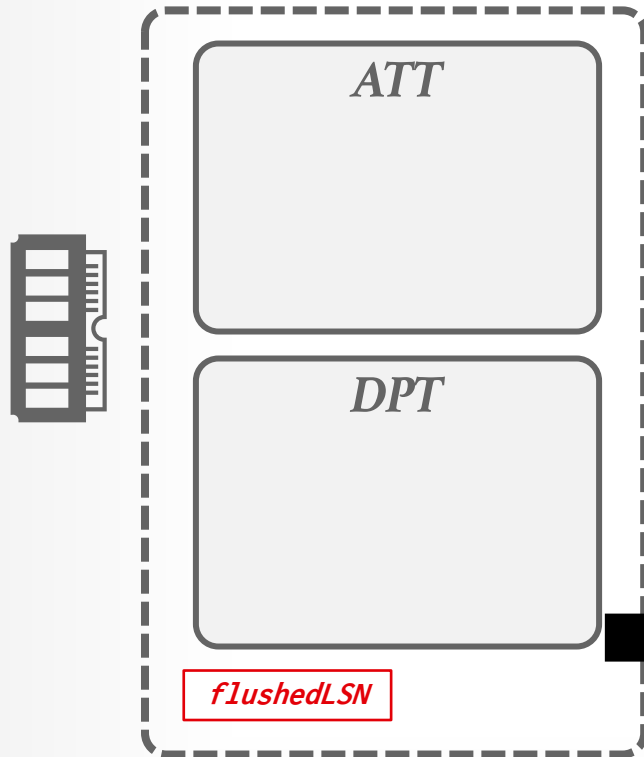
101: <T<sub>2</sub> CLR-014, B→P<sub>5</sub>, UndoNext=014>

102: <T<sub>3</sub> CLR-018, C→P<sub>1</sub>, UndoNext=005>

103: <T<sub>3</sub> TXN-END>

005: <T<sub>3</sub> BEGIN>

# FULL EXAMPLE



011: <CHECKPOINT-BEGIN>  
 012: <CHECKPOINT-END>  
 013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>  
 014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>  
 015: <T<sub>1</sub> **ABORT**>  
 016: <T<sub>1</sub> **CLR-014**, B→P<sub>3</sub>, UndoNext=013>  
 017: <T<sub>1</sub> **TXN-END**>  
 018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>  
 019: <T<sub>2</sub>, A→  
 ☠ **CRASH!**  
 101: <T<sub>2</sub> **CLR-014**, B→P<sub>5</sub>, UndoNext=014>  
 102: <T<sub>3</sub> **CLR-014**, C→P<sub>1</sub>, UndoNext=005>  
 103: <T<sub>3</sub> **TXN-END**>  
 ☠ **CRASH!**

*Flush dirty pages  
+ WAL to disk!*

005: <T<sub>3</sub> **BEGIN**>

# FULL EXAMPLE



*ATT*

TxnId	Status	lastLSN
T <sub>2</sub>	<b>U</b>	101

*DPT*

PageId	recLSN
P <sub>1</sub>	001
P <sub>3</sub>	014
P <sub>5</sub>	013

*flushedLSN*

011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> **ABORT**>

016: <T<sub>1</sub> **CLR-014**, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> **TXN-END**>

018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

💀 **CRASH!**

101: <T<sub>2</sub> **CLR-019**, D→P<sub>5</sub>, UndoNext=014>

102: <T<sub>3</sub> **CLR-018**, C→P<sub>1</sub>, UndoNext=005>

103: <T<sub>3</sub> **TXN-END**>

💀 **CRASH!**

# FULL EXAMPLE



*ATT*

TxnId	Status	lastLSN
T <sub>2</sub>	<b>U</b>	101

*DPT*

PageId	recLSN
P <sub>1</sub>	001
P <sub>3</sub>	014
P <sub>5</sub>	013

*flushedLSN*

011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> **ABORT**>

016: <T<sub>1</sub> **CLR-014**, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> **TXN-END**>

018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

**💀 CRASH!**

101: <T<sub>2</sub> **CLR-019**, D→P<sub>5</sub>, UndoNext=014>

102: <T<sub>3</sub> **CLR-018**, C→P<sub>1</sub>, UndoNext=005>

103: <T<sub>3</sub> **TXN-END**>

**💀 CRASH!**

201: <T<sub>2</sub> **CLR-014**, B→P<sub>3</sub>, UndoNext=004>



# FULL EXAMPLE



*ATT*

TxnId	Status	lastLSN
T <sub>2</sub>	<b>U</b>	101

*DPT*

PageId	recLSN
P <sub>1</sub>	001
P <sub>3</sub>	014
P <sub>5</sub>	013

*flushedLSN*

011: <CHECKPOINT-BEGIN>

012: <CHECKPOINT-END>

013: <T<sub>1</sub>, A→P<sub>5</sub>, 1, 2>

014: <T<sub>2</sub>, B→P<sub>3</sub>, 2, 3>

015: <T<sub>1</sub> **ABORT**>

016: <T<sub>1</sub> **CLR-014**, B→P<sub>3</sub>, UndoNext=013>

017: <T<sub>1</sub> **TXN-END**>

018: <T<sub>3</sub>, C→P<sub>1</sub>, 4, 5>

019: <T<sub>2</sub>, A→P<sub>5</sub>, 6, 7>

☠ **CRASH!**

101: <T<sub>2</sub> **CLR-019**, D→P<sub>5</sub>, UndoNext=014>

102: <T<sub>3</sub> **CLR-018**, C→P<sub>1</sub>, UndoNext=005>

103: <T<sub>3</sub> **TXN-END**>

☠ **CRASH!**

201: <T<sub>2</sub> **CLR-014**, B→P<sub>3</sub>, UndoNext=004>

202: <T<sub>2</sub> **TXN-END**>

# ADDITIONAL CRASH ISSUES (1)

---

*What does the DBMS do if it crashes during recovery in the Analysis Phase?*

→ Nothing. Just run recovery again.

*What does the DBMS do if it crashes during recovery in the Redo Phase?*

→ Again nothing. Redo everything again.

# ADDITIONAL CRASH ISSUES (2)

---

*How can the DBMS improve performance during recovery in the Redo Phase?*

- Assume that it is not going to crash again and flush all changes to disk asynchronously in the background.

*How can the DBMS improve performance during recovery in the Undo Phase?*

- Lazily rollback changes before new txns access pages.
- Rewrite the application to avoid long-running txns.

# CONCLUSION

---

## Mains ideas of ARIES:

- WAL with **STEAL** + **NO-FORCE**
- Fuzzy Checkpoints (snapshot of dirty page ids)
- Redo everything since the earliest dirty page
- Undo txns that never commit
- Write CLRs when undoing, to survive failures during restarts

## Log Sequence Numbers:

- LSNs identify log records; linked into backwards chains per transaction via **prevLSN**.
- **pageLSN** allows comparison of data page and log records.

# NEXT CLASS

---

You now know how to build a single-node DBMS.

Let's make it even more challenging and start talking about distributed DBMSs!