# Carnegie Mellon University

# 21 Database Recovery

Intro to Database Systems
15-445/15-645
Fall 2020

AP Andy Pavlo
Computer Science
Carnegie Mellon University

# UPCOMING DATABASE TALKS

**Confluent ksqlDB (Kafka)**
→ Monday Nov 23rd @ 5pm ET

**Microsoft SQL Server Optimizer**
→ Monday Nov 30th @ 5pm ET

**Snowflake Lecture**
→ Monday Dec 7th @ 3:20pm ET

# 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*

# ARIES

**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 (Algorithm 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 loser 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

# ARIES – MAIN IDEAS

**Write-Ahead Logging:**
→ Any change is recorded in log on stable storage before the database change is written to disk.
→ Must use **STEAL** + **NO-FORCE** buffer pool policies.

**Repeating History During Redo:**
→ On 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

# WAL RECORDS

We need to extend our log record format from last class 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…

# LOG SEQUENCE NUMBERS

| Name | Where | Definition |
|---|---|---|
| **flushedLSN** | Memory | Last LSN in log on disk |
| **pageLSN** | $page_x$ | Newest update to $page_x$ |
| **recLSN** | $page_x$ | Oldest update to $page_x$ since it was last flushed |
| **lastLSN** | $T_i$ | Latest record of txn $T_i$ |
| **MasterRecord** | Disk | LSN of latest checkpoint |

# WRITING LOG RECORDS

Each data page contains a **pageLSN**.
→ The *LSN* of the most recent update to that page.

System keeps track of **flushedLSN**.
→ The max *LSN* flushed so far.

Before page **x** can be written to disk, we must flush log at least to the point where:
→ **pageLSN$_x$ ≤ flushedLSN**

# WRITING LOG RECORDS



**Log Sequence Numbers**

WAL (Tail)

017 <T₅ **BEGIN**>
018 <T₅, A, 9, 8>
019 <T₅, B, 5, 1>
020 <T₅ **COMMIT**>
⋮

**Log Sequence Numbers**

WAL

001 <T₁ **BEGIN**>
002 <T₁, A, 1, 2>
003 <T₁ **COMMIT**>
004 <T₂ **BEGIN**>
005 <T₂, A, 2, 3>
006 <T₃ **BEGIN**>
007 <**CHECKPOINT**>
008 <T₂ **COMMIT**>
009 <T₃, A, 3, 4>
010 <T₄ **BEGIN**>
011 <T₄, X, 5, 6>
012 <T₄, Y, 9, 7>
013 <T₃, B, 4, 2>
014 <T₃ **COMMIT**>
015 <T₄, B, 2, 3>
016 <T₄, C, 1, 2>

## Buffer Pool

| pageLSN | recLSN |
| --- | --- |
| A=9 | B=5 | C=2 |

*flushedLSN*

| pageLSN | recLSN |
| --- | --- |
| A=9 | B=5 | C=2 |

*MasterRecord*

Database

# WRITING LOG RECORDS



**WAL (Tail)**

```
017:<T₅ BEGIN>
018:<T₅, A, 9, 8>
019:<T₅, B, 5, 1>
020:<T₅ COMMIT>
    ⋮
```

**Buffer Pool**

| pageLSN | recLSN |
|---------|--------|

*__Not safe to unpin because__ pageLSN > flushedLSN*

**WAL**

```
001:<T₁ BEGIN>
002:<T₁, A, 1, 2>
003:<T₁ COMMIT>
004:<T₂ BEGIN>
005:<T₂, A, 2, 3>
006:<T₃ BEGIN>
007:<CHECKPOINT>
008:<T₂ COMMIT>
009:<T₃, A, 3, 4>
010:<T₄ BEGIN>
011:<T₄, X, 5, 6>
012:<T₄, Y, 9, 7>
013:<T₃, B, 4, 2>
014:<T₃ COMMIT>
015:<T₄, B, 2, 3>
016:<T₄, C, 1, 2>
```

| pageLSN | recLSN |
|---------|--------|
| A=9 | B=5 | C=2 |

*MasterRecord*

Database

# WRITING LOG RECORDS

All log records have an **LSN**.

Update the **pageLSN** every time a txn modifies a record in the page.
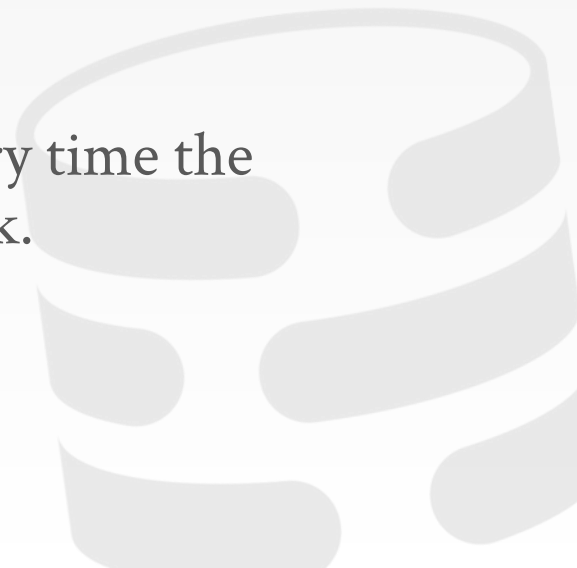
Update the **flushedLSN** in memory every time the DBMS writes out the WAL buffer to disk.

# NORMAL EXECUTION

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

Assumptions in this lecture:
→ All log records fit within a single page.
→ Disk writes are atomic.
→ Single-versioned tuples with Strict 2PL.
→ **STEAL** + **NO-FORCE** buffer management with WAL.

# TRANSACTION COMMIT

Write **COMMIT** record to log.

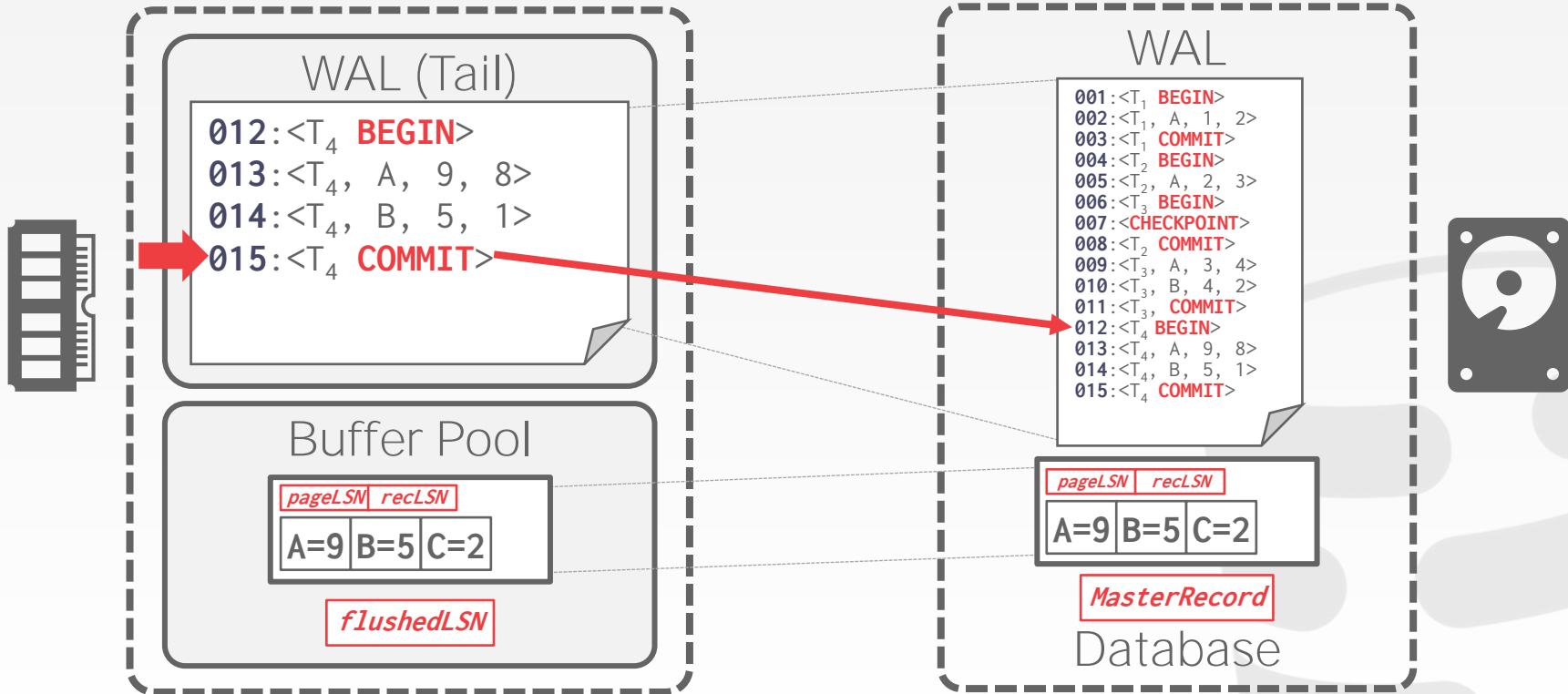All log records up to txn's **COMMIT** record are flushed to disk.
→ Log flushes are sequential, synchronous writes to disk.
→ Many log records per log page.

When the commit succeeds, write a special **TXN-END** record to log.
→ This does <u>not</u> need to be flushed immediately.

# TRANSACTION COMMIT

## WAL (Tail)

```
012:<T₄ BEGIN>
013:<T₄, A, 9, 8>
014:<T₄, B, 5, 1>
015:<T₄ COMMIT>
```

## Buffer Pool

| pageLSN | recLSN |
| --- | --- |

| A=9 | B=5 | C=2 |
| --- | --- | --- |

*flushedLSN*

## WAL

```
001:<T₁ BEGIN>
002:<T₁, A, 1, 2>
003:<T₁ COMMIT>
004:<T₂ BEGIN>
005:<T₂, A, 2, 3>
006:<T₃ BEGIN>
007:<CHECKPOINT>
008:<T₂ COMMIT>
009:<T₃, A, 3, 4>
010:<T₃, B, 4, 2>
011:<T₃, COMMIT>
012:<T₄ BEGIN>
013:<T₄, A, 9, 8>
014:<T₄, B, 5, 1>
015:<T₄ COMMIT>
```

| pageLSN | recLSN |
| --- | --- |

| A=9 | B=5 | C=2 |
| --- | --- | --- |

*MasterRecord*

Database

# TRANSACTION COMMIT



WAL (Tail)

```
012:<T₄ BEGIN>
013:<T₄, A, 9, 8>
014:<T₄, B, 5, 1>
015:<T₄ COMMIT>
```

Buffer Pool

| pageLSN | recLSN |
| --- | --- |

| A=9 | B=5 | C=2 |
| --- | --- | --- |

flushedLSN

**flushedLSN = 015**

WAL

```
001:<T₁ BEGIN>
002:<T₁, A, 1, 2>
003:<T₁ COMMIT>
004:<T₂ BEGIN>
005:<T₂, A, 2, 3>
006:<T₃ BEGIN>
007:<CHECKPOINT>
008:<T₂ COMMIT>
009:<T₃, A, 3, 4>
010:<T₃, B, 4, 2>
011:<T₃, COMMIT>
012:<T₄ BEGIN>
013:<T₄, A, 9, 8>
014:<T₄, B, 5, 1>
015:<T₄ COMMIT>
```

| pageLSN | recLSN |
| --- | --- |

| A=9 | B=5 | C=2 |
| --- | --- | --- |

*MasterRecord*

Database

# TRANSACTION COMMIT

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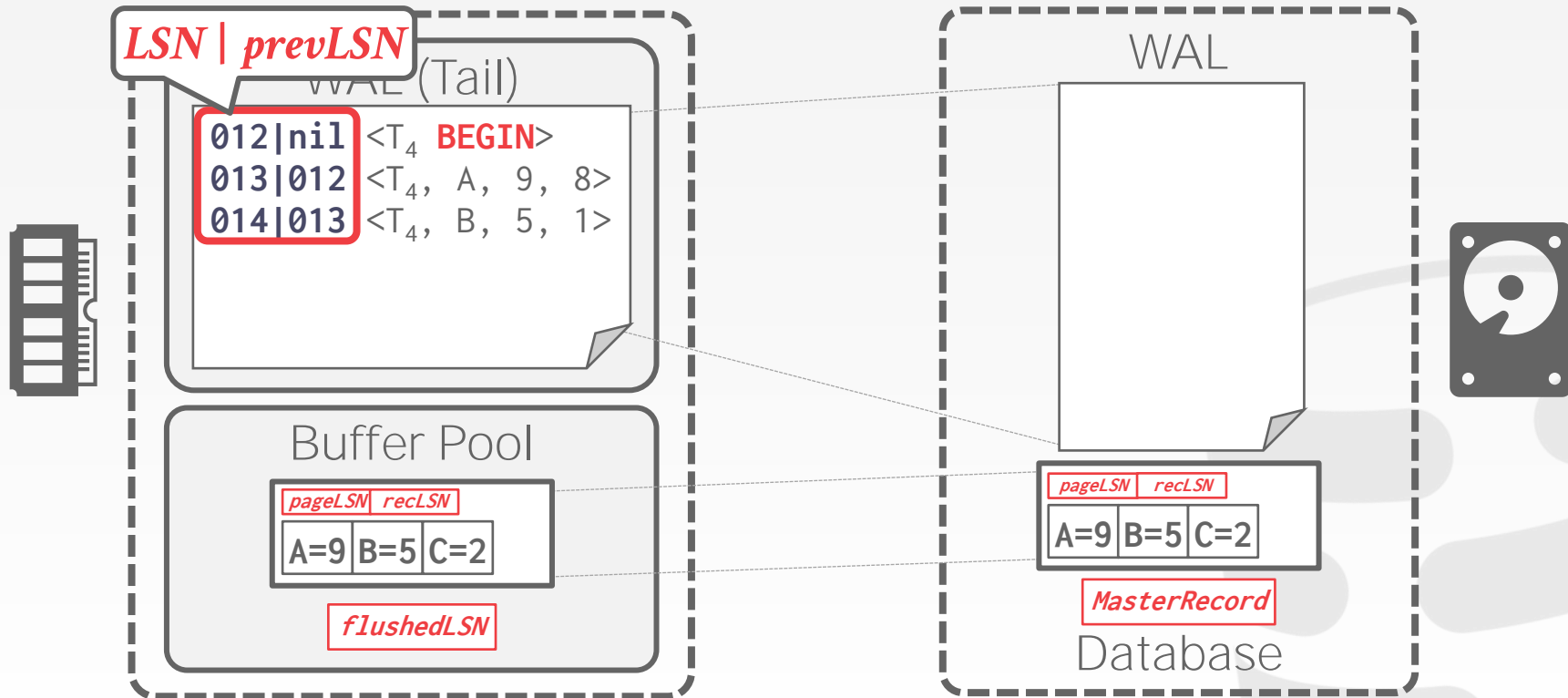
Aborting a txn is a special case of the ARIES undo operation applied to only one txn.

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 that makes it easy to walk through its records.

# TRANSACTION ABORT



**LSN | prevLSN**

WAL (Tail)

```
012|nil <T4  BEGIN>
013|012 <T4, A, 9, 8>
014|013 <T4, B, 5, 1>
```

WAL

Buffer Pool

| pageLSN | recLSN |
|---------|--------|

| A=9 | B=5 | C=2 |

*flushedLSN*

| pageLSN | recLSN |
|---------|--------|

| A=9 | B=5 | C=2 |

*MasterRecord*

Database

# TRANSACTION ABORT

# TRANSACTION ABORT

# COMPENSATION LOG RECORDS

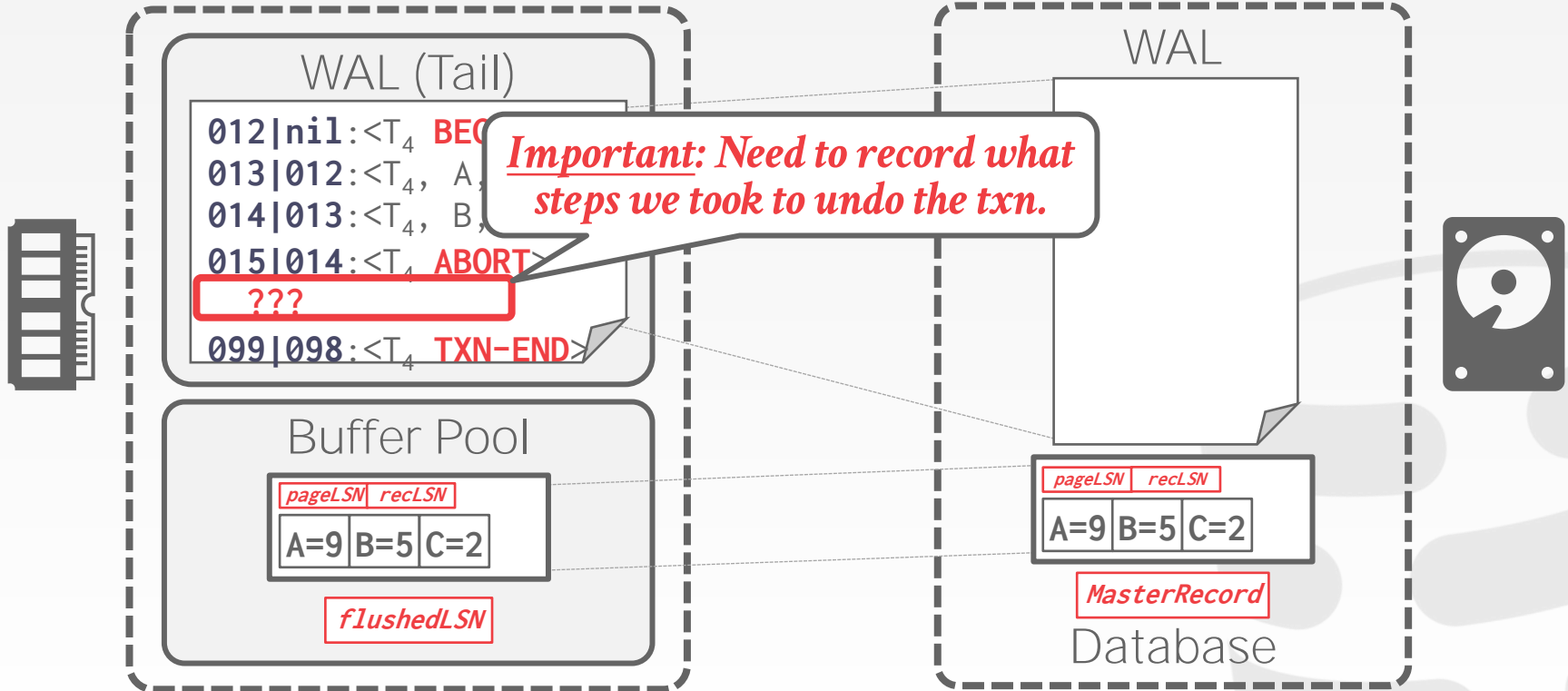A **CLR** describes the actions taken to undo the actions of a previous update record.

It has all the fields of an update log record plus the **undoNext** pointer (the next-to-be-undone LSN).

*CLRs* are added to log other records but the DBMS does not wait for them to be flushed before notifying the application that the txn aborted.

# TRANSACTION ABORT – CLR EXAMPLE

TIME

| LSN | prevLSN | TxnId | Type | Object | Before | After | UndoNext |
|-----|---------|-------|------|--------|--------|-------|----------|
| 001 | nil | $T_1$ | BEGIN | – | – | – | – |
| 002 | 001 | $T_1$ | UPDATE | A | 30 | 40 | – |
| ⋮ | | | | | | | |
| 011 | 002 | $T_1$ | ABORT | – | – | – | – |

# TRANSACTION ABORT – CLR EXAMPLE

TIME

| LSN | prevLSN | TxnId | Type | Object | Before | After | UndoNext |
|-----|---------|-------|------|--------|--------|-------|----------|
| 001 | nil | $T_1$ | BEGIN | – | – | – | – |
| 002 | 001 | $T_1$ | UPDATE | A | 30 | 40 | – |
| ⋮ | | | | | | | |
| 011 | 002 | $T_1$ | ABORT | – | – | – | – |
| ⋮ | | | | | | | |
| 026 | 011 | $T_1$ | CLR-002 | A | 40 | 30 | 001 |

# TRANSACTION ABORT — CLR EXAMPLE

TIME

| LSN | prevLSN | TxnId | Type | Object | Before | After | UndoNext |
|-----|---------|-------|------|--------|--------|-------|----------|
| 001 | nil | $T_1$ | BEGIN | – | – | – | – |
| 002 | 001 | $T_1$ | UPDATE | A | 30 | 40 | – |
| ⋮ | | | | | | | |
| 011 | 002 | $T_1$ | ABORT | – | – | – | – |
| ⋮ | | | | | | | |
| 026 | 011 | $T_1$ | CLR-002 | A | 40 | 30 | 001 |

# TRANSACTION ABORT – CLR EXAMPLE



| LSN | prevLSN | TxnId | Type | Object | Before | After | UndoNext |
|-----|---------|-------|------|--------|--------|-------|----------|
| 001 | nil | $T_1$ | BEGIN | – | – | – | – |
| 002 | 001 | $T_1$ | UPDATE | A | 30 | 40 | – |
| ⋮ | | | | | | | |
| 011 | 002 | $T_1$ | ABORT | – | – | – | – |
| ⋮ | | | | | | | |
| 026 | 011 | $T_1$ | CLR-002 | A | 40 | 30 | 001 |

TIME

# TRANSACTION ABORT – CLR EXAMPLE

TIME

| LSN | prevLSN | TxnId | Type | Object | Before | After | UndoNext |
|-----|---------|-------|------|--------|--------|-------|----------|
| 001 | nil | $T_1$ | BEGIN | – | – | – | – |
| 002 | 001 | $T_1$ | UPDATE | A | 30 | 40 | – |
| ⋮ | | | | | | | |
| 011 | 002 | $T_1$ | ABORT | – | – | – | – |
| ⋮ | | | | | | | |
| 026 | 011 | $T_1$ | CLR-002 | A | 40 | 30 | 001 |

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

# TRANSACTION ABORT – CLR EXAMPLE

TIME

| LSN | prevLSN | TxnId | Type | Object | Before | After | UndoNext |
|-----|---------|-------|------|--------|--------|-------|----------|
| 001 | nil | $T_1$ | BEGIN | – | – | – | – |
| 002 | 001 | $T_1$ | UPDATE | A | 30 | 40 | – |
| ⋮ | | | | | | | |
| 011 | 002 | $T_1$ | ABORT | – | – | – | – |
| ⋮ | | | | | | | |
| 026 | 011 | $T_1$ | CLR-002 | A | 40 | 30 | 001 |
| 027 | 026 | $T_1$ | TXN-END | – | – | – | nil |

# ABORT ALGORITHM

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

Then play back the txn's updates in reverse order.
For each update record:
→ Write a **CLR** entry to the log.
→ Restore old value.

At end, write a **TXN-END** log record.

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.
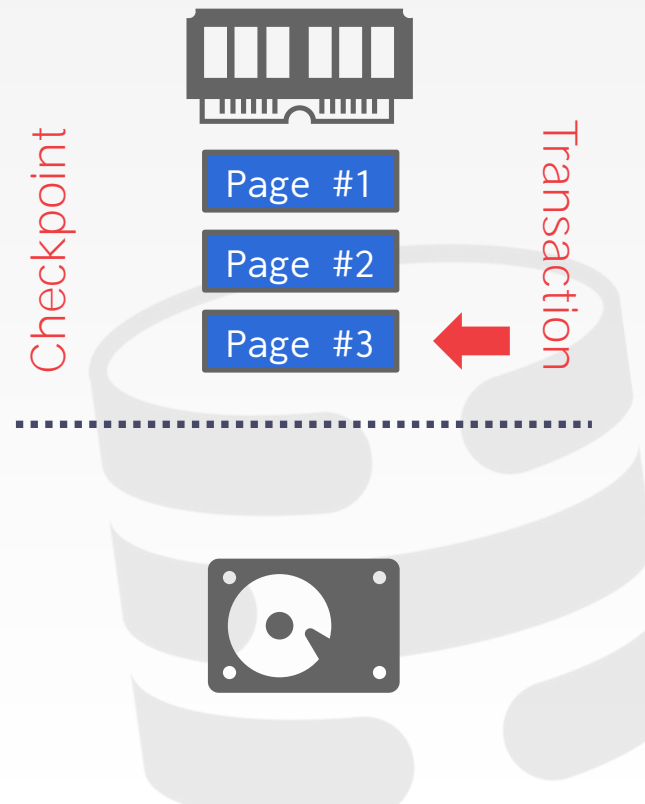→ Flushes dirty pages on disk.

This is bad for runtime performance but makes recovery trivially easy.

# SLIGHTLY BETTER CHECKPOINTS

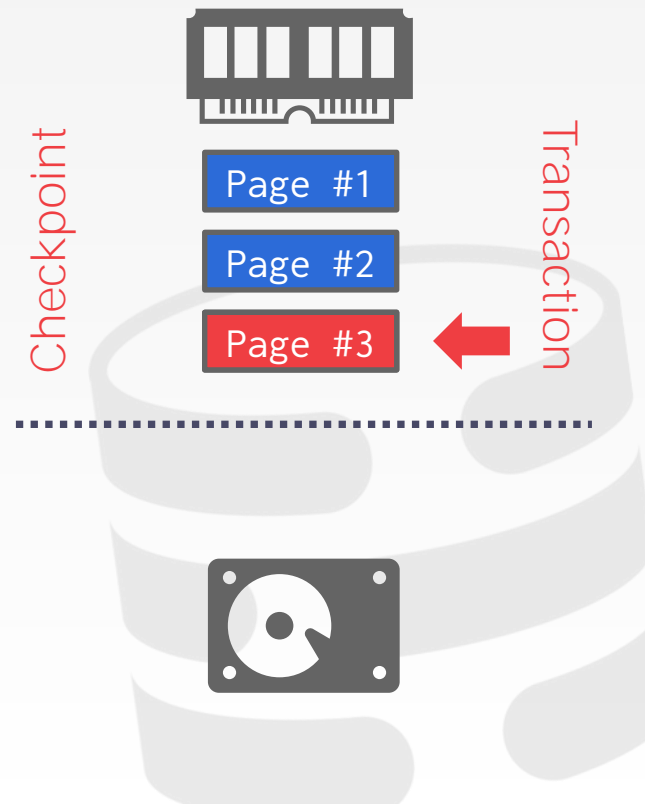Pause modifying txns while the
DBMS takes the checkpoint.
→ Prevent queries from acquiring write latch
  on table/index pages.
→ Don't have to wait until all txns finish
  before taking the checkpoint.

# SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the
DBMS takes the checkpoint.

→ Prevent queries from acquiring write latch
  on table/index pages.

→ Don't have to wait until all txns finish
  before taking the checkpoint.

# SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the
DBMS takes the checkpoint.
→ Prevent queries from acquiring write latch
on table/index pages.
→ Don't have to wait until all txns finish
before taking the checkpoint.

# SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the
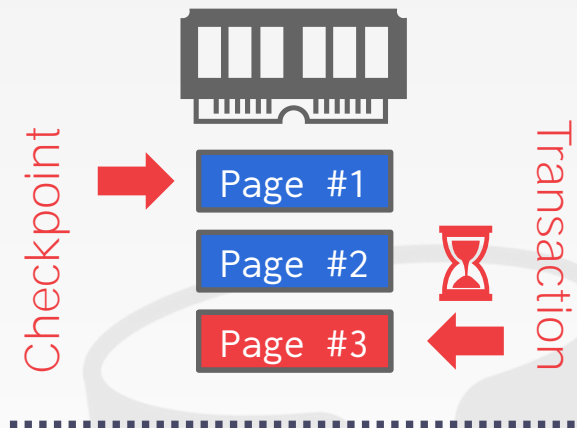DBMS takes the checkpoint.
→ Prevent queries from acquiring write latch
  on table/index pages.
→ Don't have to wait until all txns finish
  before taking the checkpoint.

# SLIGHTLY BETTER CHECKPOINTS

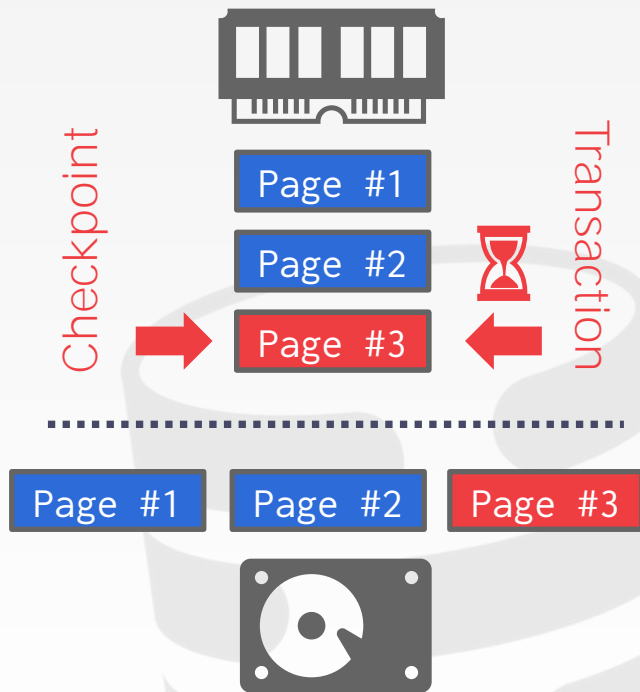Pause modifying txns while the DBMS takes the checkpoint.
→ Prevent queries from acquiring write latch on table/index pages.
→ Don't have to wait until all txns finish before taking the checkpoint.

We must record internal state as of the beginning of the checkpoint.
→ **Active Transaction Table (ATT)**
→ **Dirty Page Table (DPT)**

# ACTIVE TRANSACTION TABLE

One entry per currently active txn.
→ **txnId**: Unique txn identifier.
→ **status**: The current "mode" of the txn.
→ **lastLSN**: Most recent *LSN* created by txn.

Entry removed after the txn commits or aborts.

Txn Status Codes:
→ **R** → Running
→ **C** → Committing
→ **U** → Candidate for Undo

# DIRTY PAGE TABLE

Keep track of which pages in the buffer pool contain changes from uncommitted transactions.

One entry per dirty page in the buffer pool:
→ **recLSN**: The *LSN* of the log record that first caused the page to be dirty.

# SLIGHTLY BETTER CHECKPOINTS

At the first checkpoint, $T_2$ is still running and there are two dirty pages ($P_{11}$, $P_{22}$).

WAL

```
<T1 BEGIN>
<T2 BEGIN>
<T1, A→P11, 100, 120>
<T1 COMMIT>
<T2, C→P22, 100, 120>
<CHECKPOINT
  ATT={T2},
  DPT={P11,P22}>
<T3 START>
<T2, A→P11, 120, 130>
<T2 COMMIT>
<T3, B→P33, 200, 400>
<CHECKPOINT
  ATT={T3},
  DPT={P11,P33}>
<T3, B→P33, 400, 600>
```

# SLIGHTLY BETTER CHECKPOINTS

At the first checkpoint, $T_2$ is still running and there are two dirty pages ($P_{11}$, $P_{22}$).

At the second checkpoint, $T_3$ is active and there are two dirty pages ($P_{11}$, $P_{33}$).

This still is not ideal because the DBMS must stall txns during checkpoint…

**WAL**

```
<T1 BEGIN>
<T2 BEGIN>
<T1, A→P11, 100, 120>
<T1 COMMIT>
<T2, C→P22, 100, 120>
<CHECKPOINT
  ATT={T2},
  DPT={P11,P22}>
<T3 START>
<T2, A→P11, 120, 130>
<T2 COMMIT>
<T3, B→P33, 200, 400>
<CHECKPOINT
  ATT={T3},
  DPT={P11,P33}>
<T3, B→P33, 400, 600>
```

# FUZZY CHECKPOINTS

A *fuzzy checkpoint* is where the DBMS allows active txns to continue the run while the system flushes dirty pages to disk.

New log records to track checkpoint boundaries:
→ CHECKPOINT-BEGIN: Indicates start of checkpoint
→ CHECKPOINT-END: Contains **ATT** + **DPT**.

# FUZZY CHECKPOINT

The **LSN** of the <span style="color:red">**CHECKPOINT-BEGIN**</span> record is written to the database's <span style="color:red">**MasterRecord**</span> entry on disk when the checkpoint successfully completes.

Any txn that starts <u>after</u> the checkpoint is excluded from the ATT in the <span style="color:red">**CHECKPOINT-END**</span> record.

**WAL**

```
<T₁ BEGIN>
<T₂ BEGIN>
<T₁, A→P₁₁, 100, 120>
<T₁ COMMIT>
<T₂, C→P₂₂, 100, 120>
<CHECKPOINT-BEGIN>
<T₃ START>
<T₂, A→P₁₁, 120, 130>
<CHECKPOINT-END
    ATT={T₂},
    DPT={P₁₁}>
<T₂ COMMIT>
<T₃, B→P₃₃, 200, 400>
<CHECKPOINT-BEGIN>
<T₃, B→P₃₃, 10, 12>
<CHECKPOINT-END
    ATT={T₃},
    DPT={P₃₃}>
```

# ARIES – RECOVERY PHASES

**Phase #1 – Analysis**
→ Read WAL from last `CHECKPOINT-END` to identify dirty pages in the buffer pool and active txns at the time of the crash.

**Phase #2 – Redo**
→ Repeat <u>all</u> 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 <u>all</u> actions.

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

# ANALYSIS PHASE

Scan log forward from last successful checkpoint.

If you find a **TXN-END** record, remove its corresponding txn from **ATT**.

All other records:
→ Add txn to **ATT** with status **UNDO**.
→ On commit, change txn status to **COMMIT**.

For **UPDATE** records:
→ If page **P** not in **DPT**, add **P** 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

### WAL

```
010:<CHECKPOINT-BEGIN>
    ⋮
020:<T_96, A→P_33, 10, 15>
    ⋮
030:<CHECKPOINT-END
     ATT={T_96,T_97},
     DPT={P_20,P_33}>
    ⋮
040:<T_96 COMMIT>
    ⋮
050:<T_96 TXN-END>
    ⋮
CRASH!
```

| LSN | ATT | DPT |
|-----|-----|-----|
| 010 | | |
| 020 | | |
| 030 | | |
| 040 | | |
| 050 | | |

# ANALYSIS PHASE EXAMPLE

**WAL**

```
010:<CHECKPOINT-BEGIN>
    ⋮
020:<T96, A→P33, 10, 15>
    ⋮
030:<CHECKPOINT-END
      ATT={T96,T97},
      DPT={P20,P33}>
    ⋮
040:<T96 COMMIT>
    ⋮
050:<T96 TXN-END>
    ⋮
CRASH!
```

| LSN | ATT | DPT |
|-----|-----|-----|
| 010 | | |
| 020 | $(T_{96}, U)$ | |
| 030 | | |
| 040 | | |
| 050 | | |

*(TxnId, Status)*

# ANALYSIS PHASE EXAMPLE

WAL

**010**:<**CHECK**...

*Modify A in page $P_{33}$*

**020**:<$T_{96}$, A→$P_{33}$, 10, 15>

**030**:<**CHECKPOINT-END**
        ATT={$T_{96}$,$T_{97}$},
        DPT={$P_{20}$,$P_{33}$}>

**040**:<$T_{96}$ **COMMIT**>

**050**:<$T_{96}$ **TXN-END**>

**CRASH!**

| | ATT | DPT |
|---|---|---|
| **010** | | |
| **020** | ($T_{96}$,**U**) | ($P_{33}$,**020**) |
| **030** | | |
| **040** | | |
| **050** | | |

*(PageId, RecLSN)*

# ANALYSIS PHASE EXAMPLE



WAL

```
010:<CHECKPOINT-BEGIN>
   ⋮
020:<T96, A→P33, 10, 15>
   ⋮
030:<CHECKPOINT-END
       ATT={T96,T97},
       DPT={P20,P33}>
   ⋮
040:<T96 COMMIT>
   ⋮
050:<T96 TXN-END>
   ⋮
CRASH!
```

| LSN | ATT | DPT |
|-----|-----|-----|
| 010 | | |
| 020 | $(T_{96}, U)$ | $(P_{33}, 020)$ |
| 030 | $(T_{96}, U)$, $(T_{97}, U)$ | $(P_{33}, 020)$, $(P_{20}, 022)$ |
| 040 | | |
| 050 | | |

# ANALYSIS PHASE EXAMPLE

### WAL

```
010:<CHECKPOINT-BEGIN>
    ⋮
020:<T96, A→P33, 10, 15>
    ⋮
030:<CHECKPOINT-END
    ATT={T96,T97},
    DPT={P20,P33}>
    ⋮
040:<T96 COMMIT>
    ⋮
050:<T96 TXN-END>
    ⋮
CRASH!
```

| LSN | ATT | DPT |
|-----|-----|-----|
| 010 | | |
| 020 | $(T_{96}, U)$ | $(P_{33}, 020)$ |
| 030 | $(T_{96}, U)$, $(T_{97}, U)$ | $(P_{33}, 020)$, $(P_{20}, 022)$ |
| 040 | $(T_{96}, C)$, $(T_{97}, U)$ | $(P_{33}, 020)$, $(P_{20}, 022)$ |
| 050 | | |

# ANALYSIS PHASE EXAMPLE

**WAL**

```
010:<CHECKPOINT-BEGIN>
   ⋮
020:<T96, A→P33, 10, 15>
   ⋮
030:<CHECKPOINT-END
      ATT={T96,T97},
      DPT={P20,P33}>
   ⋮
040:<T96 COMMIT>
   ⋮
050:<T96 TXN-END>
   ⋮
CRASH!
```

| LSN | ATT | DPT |
|-----|-----|-----|
| 010 | | |
| 020 | $(T_{96},U)$ | $(P_{33},020)$ |
| 030 | $(T_{96},U)$, $(T_{97},U)$ | $(P_{33},020)$, $(P_{20},022)$ |
| 040 | $(T_{96},C)$, $(T_{97},U)$ | $(P_{33},020)$, $(P_{20},022)$ |
| 050 | $(T_{97},U)$ | $(P_{33},020)$, $(P_{20},022)$ |

# REDO PHASE

The goal is to repeat history to reconstruct 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 that 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 given **LSN**, redo the action unless:
→ Affected page is not in **DPT**, <u>or</u>
→ Affected page is in **DPT** but that record's **LSN** is less than the page's **recLSN**.

# REDO PHASE

To redo an action:
→ Reapply logged action.
→ 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 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 **U** status in the **ATT** after the Analysis Phase.

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

Write a **CLR** for every modification.

# FULL EXAMPLE



**LSN**     **LOG**

00 — <**CHECKPOINT-BEGIN**>

05 — <**CHECKPOINT-END**>

10 — <$T_1$, A→$P_5$, 1, 2>

20 — <$T_2$, B→$P_3$, 2, 3>

30 — <$T_1$ **ABORT**>

40 — <**CLR**: Undo $T_1$ LSN **10**>

45 — <$T_1$ **TXN-END**>

50 —

60 —

TIME

prevLSNs

# FULL EXAMPLE



| LSN | LOG |
|---|---|
| 00 | <CHECKPOINT-BEGIN> |
| 05 | <CHECKPOINT-END> |
| 10 | <$T_1$, A→$P_5$, 1, 2> |
| 20 | <$T_2$, B→$P_3$, 2, 3> |
| 30 | <$T_1$ ABORT> |
| 40 | <CLR: Undo $T_1$ LSN 10> |
| 45 | <$T_1$ TXN-END> |
| 50 | <$T_3$, C→$P_1$, 4, 5> |
| 60 | <$T_2$, D→$P_5$, 6, 7> |
| ✗ | CRASH! |

TIME

# FULL EXAMPLE



**LSN**    **LOG**

00,05   `<CHECKPOINT-BEGIN>, <CHECKPOINT-END>`

10   `<T₁, A➜P₅, 1, 2>`

20   `<T₂, B➜P₃, 2, 3>`

30   `<T₁ ABORT>`

40,45   `<CLR: Undo T₁ LSN 10>, <T₁ TXN-END>`

50   `<T₃, C➜P₁, 4, 5>`

60   `<T₂, D➜P₅, 6, 7>`

✗ *CRASH! RESTART!*

### ATT

| TxnId | Status | lastLSN |
|-------|--------|---------|
| T₂ | U | 60 |
| T₃ | U | 50 |
| – | – | – |

### DPT

| PageId | recLSN |
|--------|--------|
| P₁ | 50 |
| P₃ | 08 |
| P₅ | 10 |

*flushedLSN*

# FULL EXAMPLE

**LSN**  **LOG**

**ATT**

| TxnId | Status | lastLSN |
|-------|--------|---------|
| $T_2$ | U | 60 |
| $T_3$ | U | 50 |
| – | – | – |

**DPT**

| PageId | recLSN |
|--------|--------|
| $P_1$ | 50 |
| $P_3$ | 08 |
| $P_5$ | 10 |

*flushedLSN*

**00,05** ── <**CHECKPOINT-BEGIN**>, <**CHECKPOINT-END**>

**10** ── <$T_1$, A→$P_5$, 1, 2>

**20** ── <$T_2$, B→$P_3$, 2, 3>

**30** ── <$T_1$ **ABORT**>

**40,45** ── <**CLR**: Undo $T_1$ LSN **10**>, <$T_1$ **TXN-END**>

**50** ── <$T_3$, C→$P_1$, 4, 5>

**60** ── <$T_2$, D→$P_5$, 6, 7>

✗ *CRASH! RESTART!*

**70** ── <**CLR**: Undo $T_2$ LSN **60**, UndoNext **20**>

CMU·DB

# FULL EXAMPLE

**LSN**    **LOG**

**ATT**

| TxnId | Status | lastLSN |
|-------|--------|---------|
| T$_2$ | U | 60 |
| T$_3$ | U | 50 |
| – | – | – |

**DPT**

| PageId | recLSN |
|--------|--------|
| P$_1$ | 50 |
| P$_3$ | 08 |
| P$_5$ | 10 |

*flushedLSN*

**00,05** — <**CHECKPOINT-BEGIN**>, <**CHECKPOINT-END**>

**10** — <T$_1$, A➡P$_5$, 1, 2>

**20** — <T$_2$, B➡P$_3$, 2, 3>

**30** — <T$_1$ **ABORT**>

**40,45** — <**CLR**: Undo T$_1$ LSN **10**>, <T$_1$ **TXN-END**>

**50** — <T$_3$, C➡P$_1$, 4, 5>

**60** — <T$_2$, D➡P$_5$, 6, 7>

✗ *CRASH! RESTART!*

**70** — <**CLR**: Undo T$_2$ LSN **60**, UndoNext ... >

**80,85** — <**CLR**: Undo T$_3$ LSN **50**>, <T$_3$ **TXN-END**>
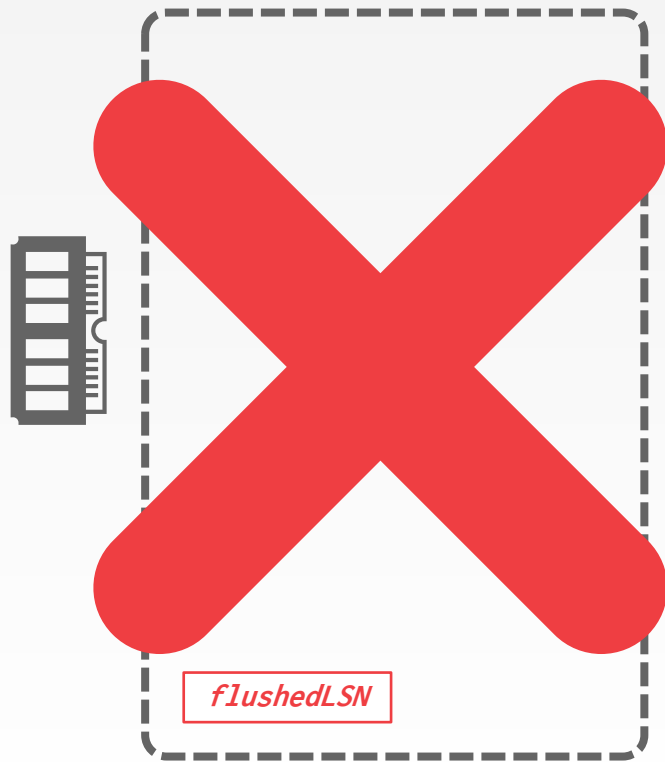
*Flush dirty pages + WAL to disk!*

# FULL EXAMPLE



**LSN**     **LOG**

**00,05**   <CHECKPOINT-BEGIN>, <CHECKPOINT-END>

**10**   <$T_1$, A→$P_5$, 1, 2>

**20**   <$T_2$, B→$P_3$, 2, 3>

**30**   <$T_1$ ABORT>

**40,45**   <CLR: Undo $T_1$ LSN **10**>, <$T_1$ TXN-END>

**50**   <$T_3$, C→$P_1$, 4, 5>

**60**   <$T_2$, D→$P_5$, 6, 7>

    CRASH! RESTART!

**70**   <CLR: Undo $T_2$ LSN **60**, UndoNext ...>

**80,85**   <CLR: Undo $T_3$ LSN **50**>, <$T_3$ TXN-END>

    CRASH! RESTART!

**ATT**

| TxnId | Status | lastLSN |
|-------|--------|---------|
| $T_2$ | U | 60 |
| $T_3$ | U | 50 |
| – | – | – |

**DPT**

| PageId | recLSN |
|--------|--------|
| $P_1$ | 50 |
| $P_3$ | 08 |
| $P_5$ | 10 |

*flushedLSN*

*Flush dirty pages + WAL to disk!*

# FULL EXAMPLE

# FULL EXAMPLE

**LSN**    **LOG**

**00,05**   <**CHECKPOINT-BEGIN**>, <**CHECKPOINT-END**>

**10**   <$T_1$, A→$P_5$, 1, 2>

**20**   <$T_2$, B→$P_3$, 2, 3>

**30**   <$T_1$ **ABORT**>

**40,45**   <**CLR**: Undo $T_1$ LSN **10**>, <$T_1$ **TXN-END**>

**50**   <$T_3$, C→$P_1$, 4, 5>

**60**   <$T_2$, D→$P_5$, 6, 7>

   *CRASH! RESTART!*

**70**   <**CLR**: Undo $T_2$ LSN **60**, UndoNext **20**>

**80,85**   <**CLR**: Undo $T_3$ LSN **50**>, <$T_3$ **TXN-END**>

   *CRASH! RESTART!*

## ATT

| TxnId | Status | lastLSN |
|-------|--------|---------|
| $T_2$ | U | 70 |
| – | – | – |
| – | – | – |

## DPT

| PageId | recLSN |
|--------|--------|
| $P_1$ | 50 |
| $P_3$ | 08 |
| $P_5$ | 10 |

*flushedLSN*

# FULL EXAMPLE



| ATT | | |
|---|---|---|
| **TxnId** | **Status** | **lastLSN** |
| $T_2$ | U | 70 |
| – | – | – |
| – | – | – |

| DPT | |
|---|---|
| **PageId** | **recLSN** |
| $P_1$ | 50 |
| $P_3$ | 08 |
| $P_5$ | 10 |

*flushedLSN*

| LSN | LOG |
|---|---|
| 00,05 | `<CHECKPOINT-BEGIN>, <CHECKPOINT-END>` |
| 10 | `<T1, A→P5, 1, 2>` |
| 20 | `<T2, B→P3, 2, 3>` |
| 30 | `<T1 ABORT>` |
| 40,45 | `<CLR: Undo T1 LSN 10>, <T1 TXN-END>` |
| 50 | `<T3, C→P1, 4, 5>` |
| 60 | `<T2, D→P5, 6, 7>` |
| | *CRASH! RESTART!* |
| 70 | `<CLR: Undo T2 LSN 60, UndoNext 20>` |
| 80,85 | `<CLR: Undo T3 LSN 50>, <T3 TXN-END>` |
| | *CRASH! RESTART!* |

# FULL EXAMPLE

**LSN**  **LOG**

**00,05** — <**CHECKPOINT-BEGIN**>, <**CHECKPOINT-END**>

**10** — <$T_1$, A→$P_5$, 1, 2>

**20** — <$T_2$, B→$P_3$, 2, 3>

**30** — <$T_1$ **ABORT**>

**40,45** — <**CLR**: Undo $T_1$ LSN **10**>, <$T_1$ **TXN-END**>

**50** — <$T_3$, C→$P_1$, 4, 5>

**60** — <$T_2$, D→$P_5$, 6, 7>

✗ *CRASH! RESTART!*

**70** — <**CLR**: Undo $T_2$ LSN **60**, UndoNext **20**>

**80,85** — <**CLR**: Undo $T_3$ LSN **50**>, <$T_3$ **TXN-END**>

✗ *CRASH! RESTART!*

## ATT

| TxnId | Status | lastLSN |
|-------|--------|---------|
| $T_2$ | U | 70 |
| – | – | – |
| – | – | – |

## DPT

| PageId | recLSN |
|--------|--------|
| $P_1$ | 50 |
| $P_3$ | 08 |
| $P_5$ | 10 |

*flushedLSN*

# 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.

So now we can talk about distributed databases!