



Intro to Database Systems (15-445/645)

20 Database Recovery

Carnegie
Mellon
University

FALL
2022

Andy
Pavlo

ADMINISTRIVIA

Project #3 is due **Wed Nov 16th @ 11:59pm**

Project #4 is due **Sun Dec 11th @ 11:59pm**

→ Zoom Q&A Session Thu Nov 17th @ 8:00pm

We are looking for spirited and impressionable TAs for 15-445/645 in Spring 2023.

- All BusTub projects will remain in C++.
- I will announce this on Piazza.

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, it is ensured that rollbacks of those written during rollbacks themselves, a bounded number of logins is incurred. In parallel, in the face of expected 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 levels of durability. By employing parallelism during both parallelized rollbacks and logical redo, it enhances consistency and performance. We show how 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

Changes are written
to log before they're
written to disk.

Force \Rightarrow every update has to be written back to disk
before commit

ARIES - MAIN IDEAS

write ahead log

use steal +

NO FORCE
policy

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.

These WAL
records

Repeating History During Redo:

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

We see them

Logging Changes During Undo:

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

After restart
check if any fails

retrace db.

= retrace actions?

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.

LSN → an identifier

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

of log record.

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

LSN → represent physical order that txns make changes to db.

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

various
cards of
LSN

maintained by
var part of
consumer

LOG SEQUENCE NUMBERS

Name	Location	Definition
<u>flushedLSN</u>	Memory	Last LSN in <u>log</u> 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

flushed LSN → last LSN
 stored in

*pageLSN
↓
most recent update
to page L := stored in
data page
itself*

WRITING LOG RECORDS

*pageLSN, newest
update
to page x.*

Each data page contains a **pageLSN**.

→ The *LSN* of the most recent update to that page.

flushed LSN + The

System keeps track of **flushedLSN**.

→ The max *LSN* flushed so far.

System tracks flushed LSN

*last LSN on
the disk*

Before the DBMS can write page **x** to disk, it must flush the log at least to the point where:

→ $\text{pageLSN}_x \leq \text{flushedLSN}$

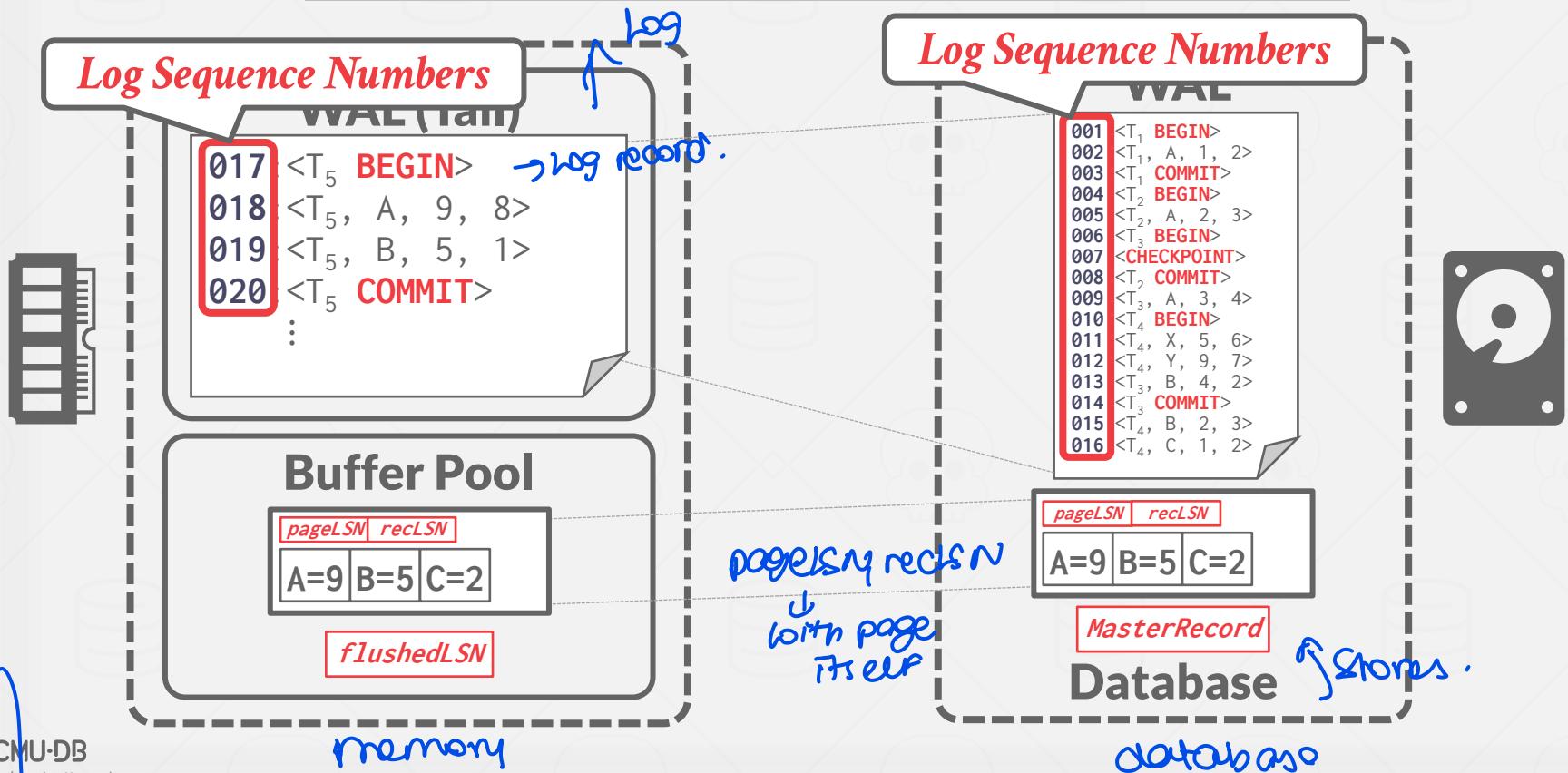


*Before writing we need
to check that*

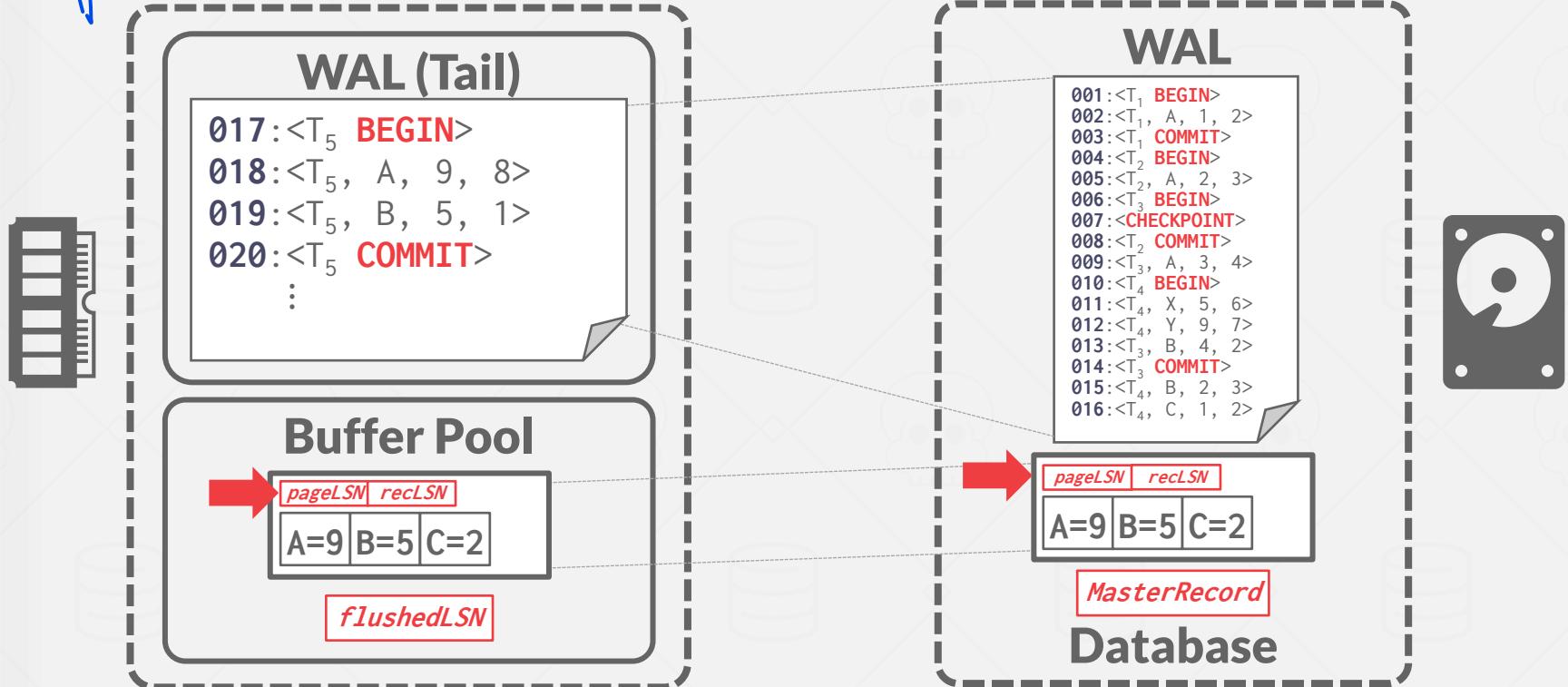
*that page from log
pool*

$\text{pageLSN}_x \leq \text{flushedLSN}$

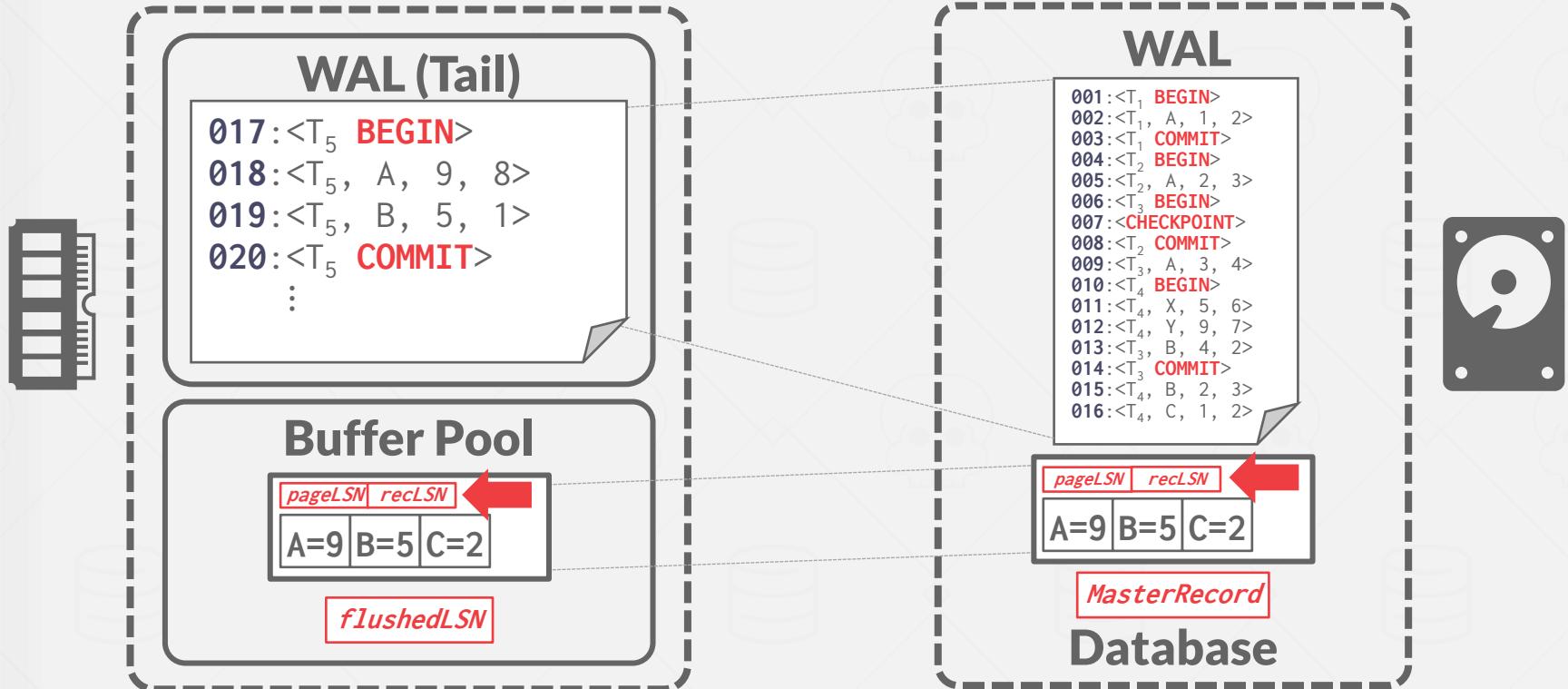
WRITING LOG RECORDS



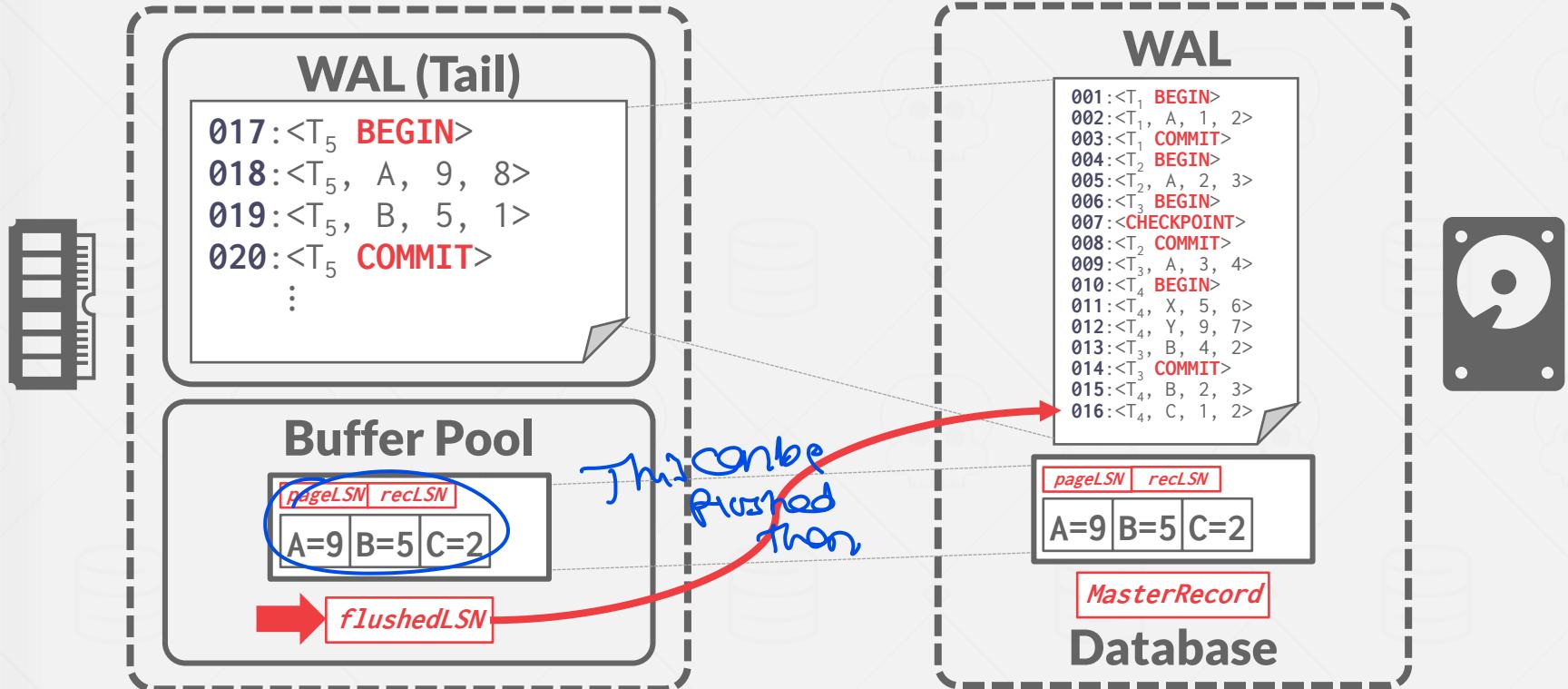
WRITING LOG RECORDS



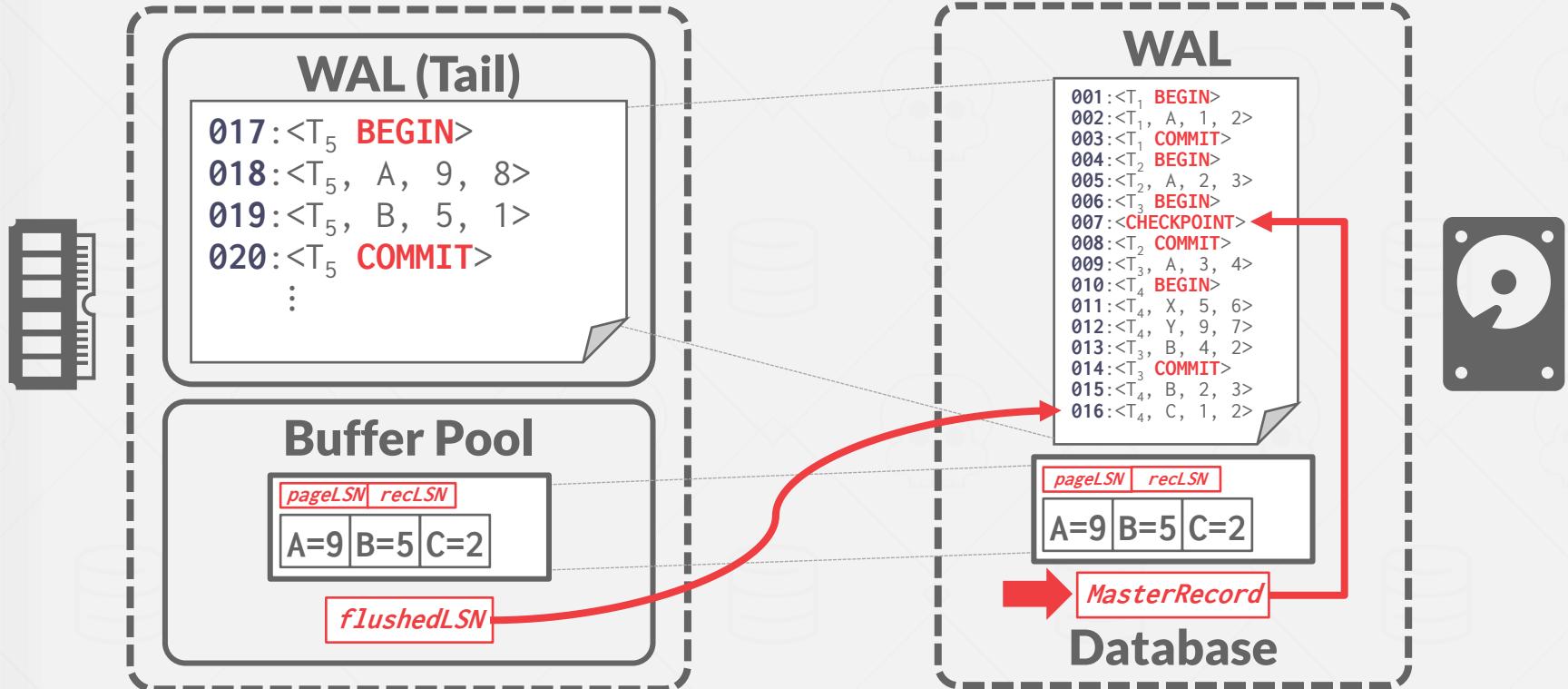
WRITING LOG RECORDS



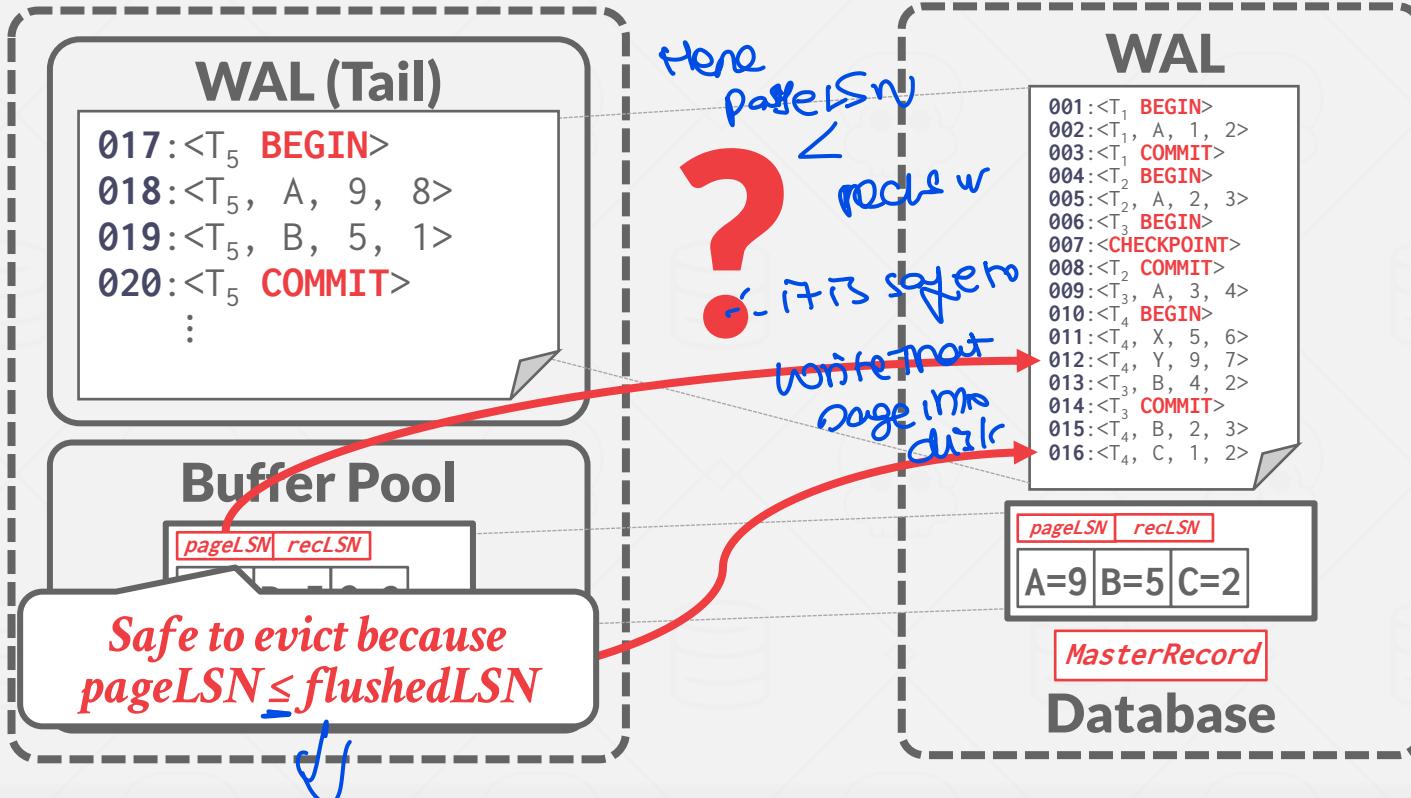
WRITING LOG RECORDS



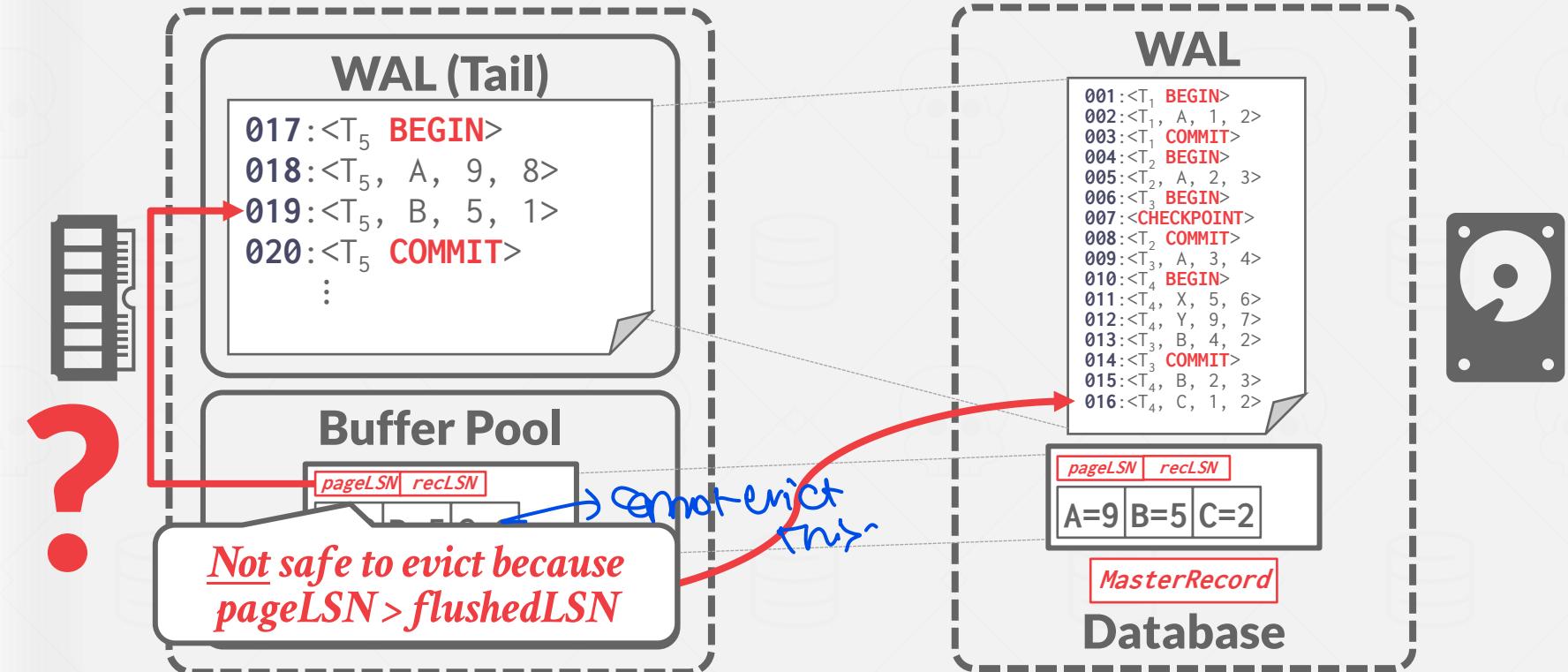
WRITING LOG RECORDS



WRITING LOG RECORDS



WRITING LOG RECORDS



WRITING LOG RECORDS

All log records have an *LSN*.

-> replay of the log page.

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

*Ample'n page
each time a page is modified
update the page*

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

When doing writes out several WAL records to disk choose the largest LSN record among them and that will be our flushed LSN.

NORMAL EXECUTION

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

we assume no writes can be stopped.

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.

↓ Best for recovery → Buffer policy used in WAL as well.

Before committing

txn we add
commit record to
log and then

flush all
log records to
disk.

↓
synchronous
flushing.

TRANSACTION COMMIT

Async flushing

↓
after txn

committed after

some time
we tell that

flushed to disk

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.

↑ At a later point we add a
txn end record

to log

nothing after

this point no

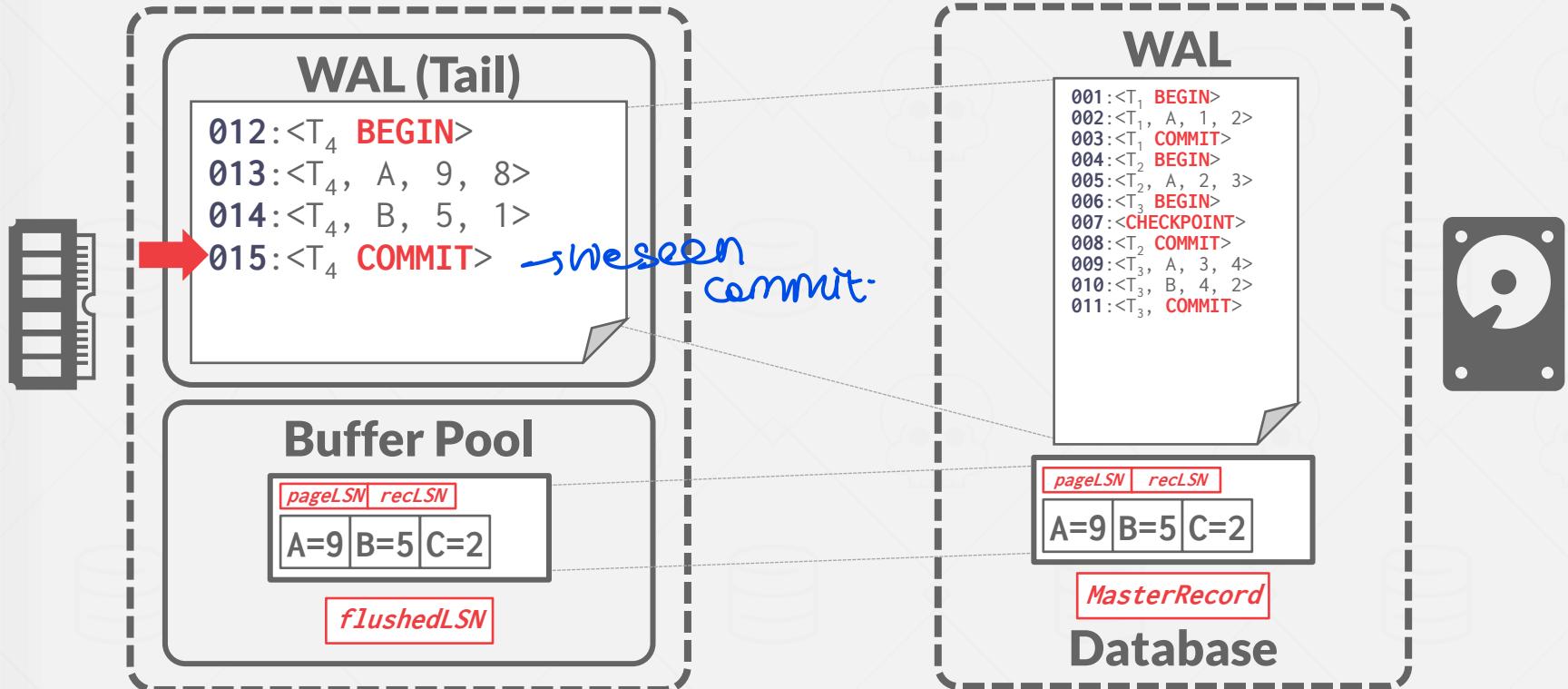
records will happen
for this txn.

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

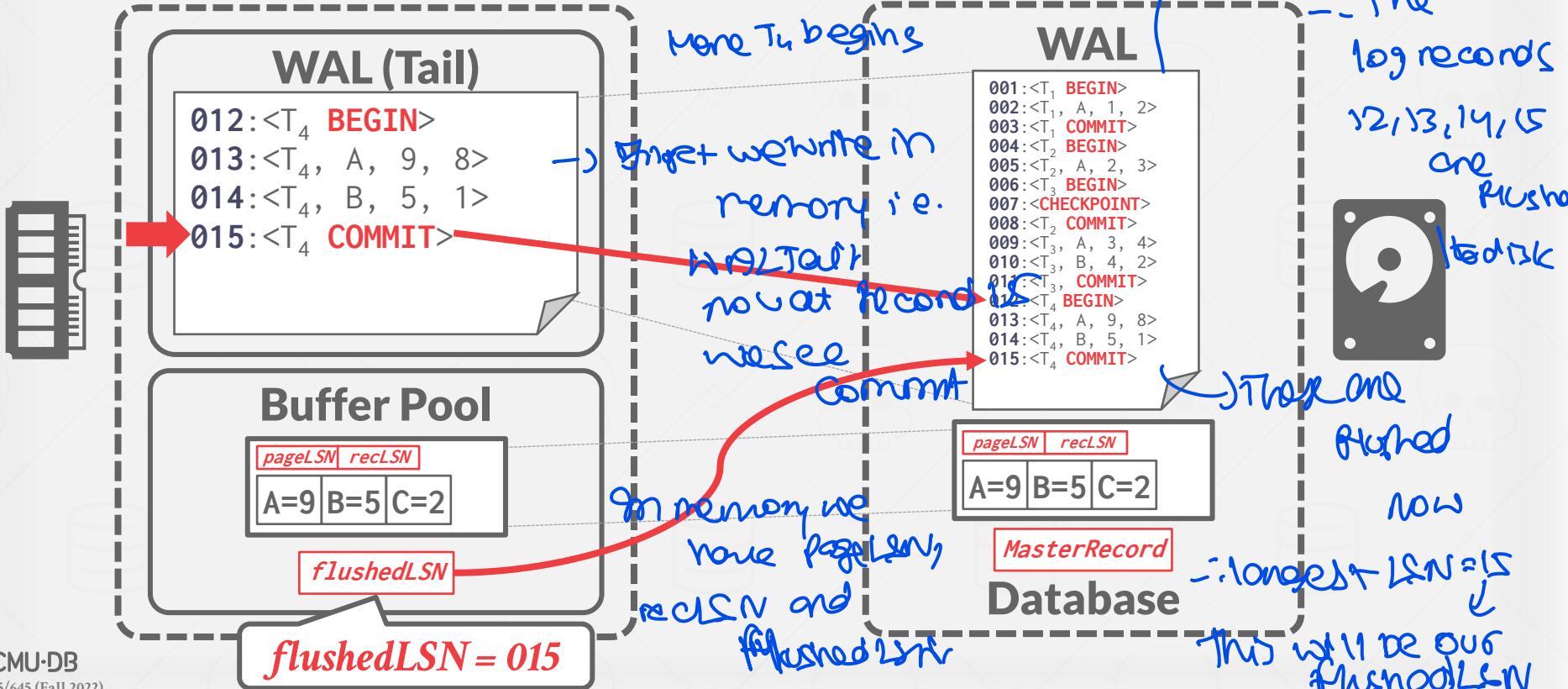
- Indicates that no new log record for a txn will appear in the log ever again.
- This does not need to be flushed immediately.

↳ this need not be flushed immediately.

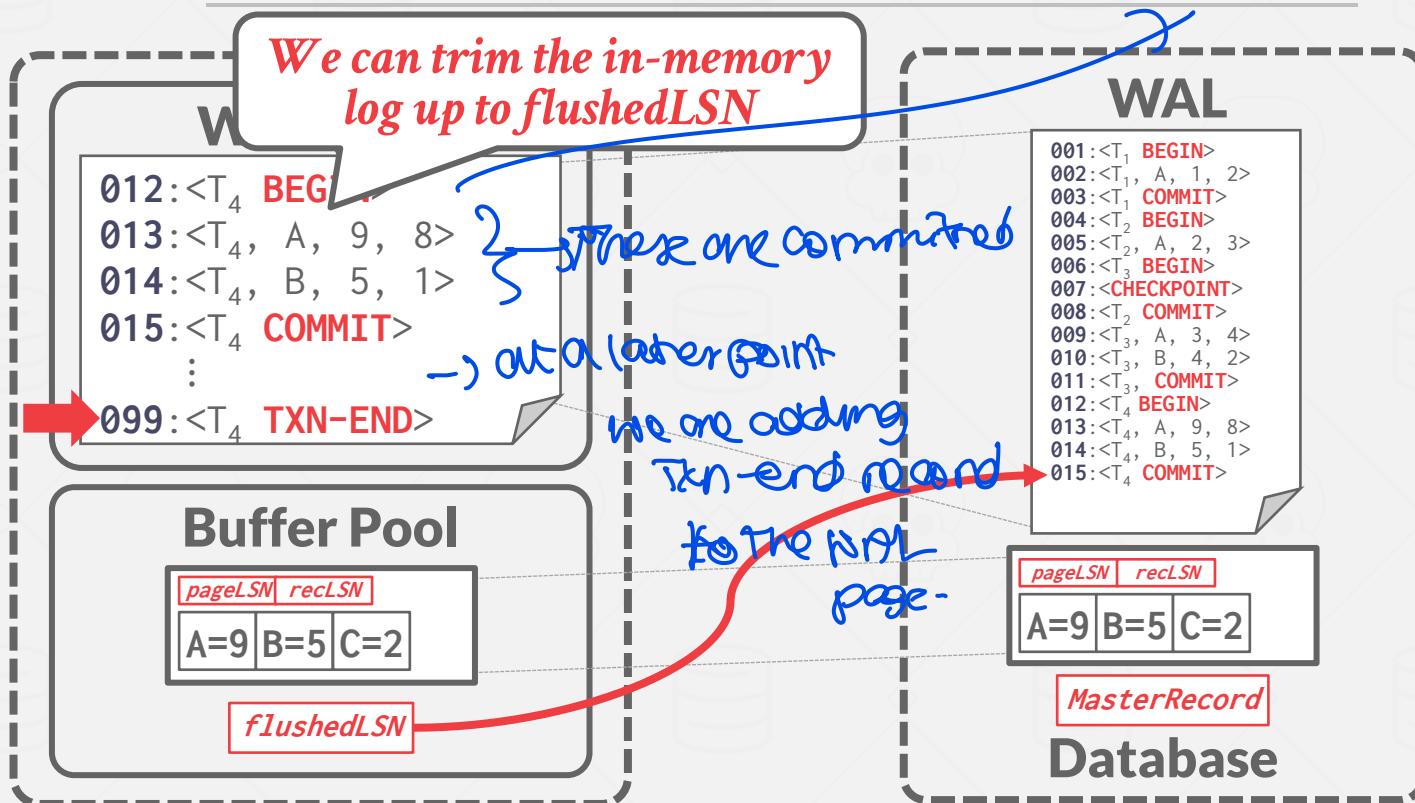
TRANSACTION COMMIT



TRANSACTION COMMIT



TRANSACTION COMMIT



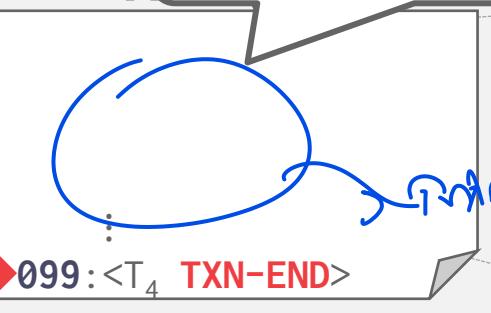
^ In this page whatever deflashed space of txns can be trimmed away.

TRANSACTION COMMIT

We can trim the in-memory log up to flushedLSN



W



Trimmed space
until LSN
it's been
flushed

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>
  
```



Buffer Pool

pageLSN	recLSN
A=9	B=5 C=2

flushedLSN

pageLSN	recLSN
A=9	B=5 C=2

MasterRecord

Database

when a
txn abort
w/o open type.

TRANSACTION ABORT

- To keep track of

what we
aborted last

so we
keep new field

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

i.e. $\text{prevLSN} \rightarrow \text{NULL}$ for 1st record of a txn.

We need to add another field to our log records:

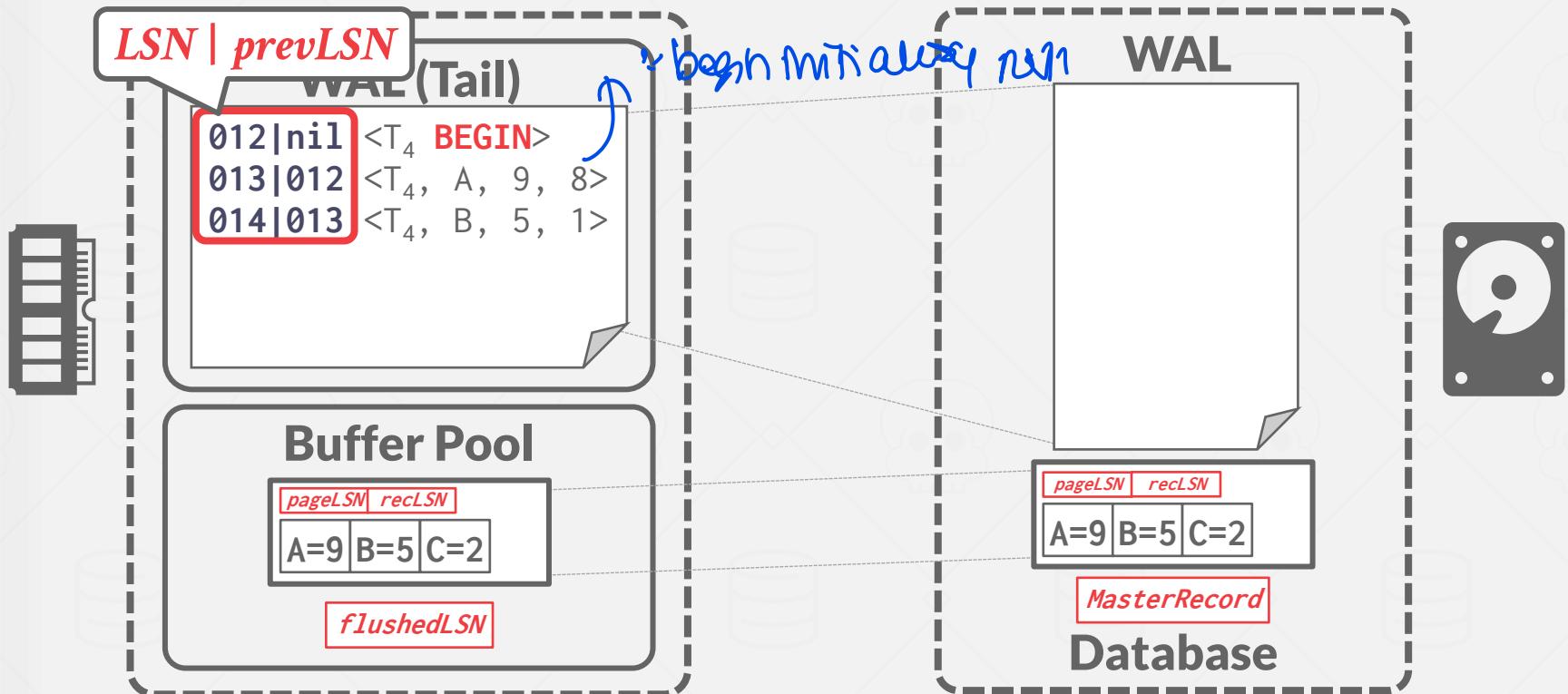
→ **prevLSN**: The previous *LSN* for the txn.

→ This maintains a linked-list for each txn that makes it easy to walk through its records.

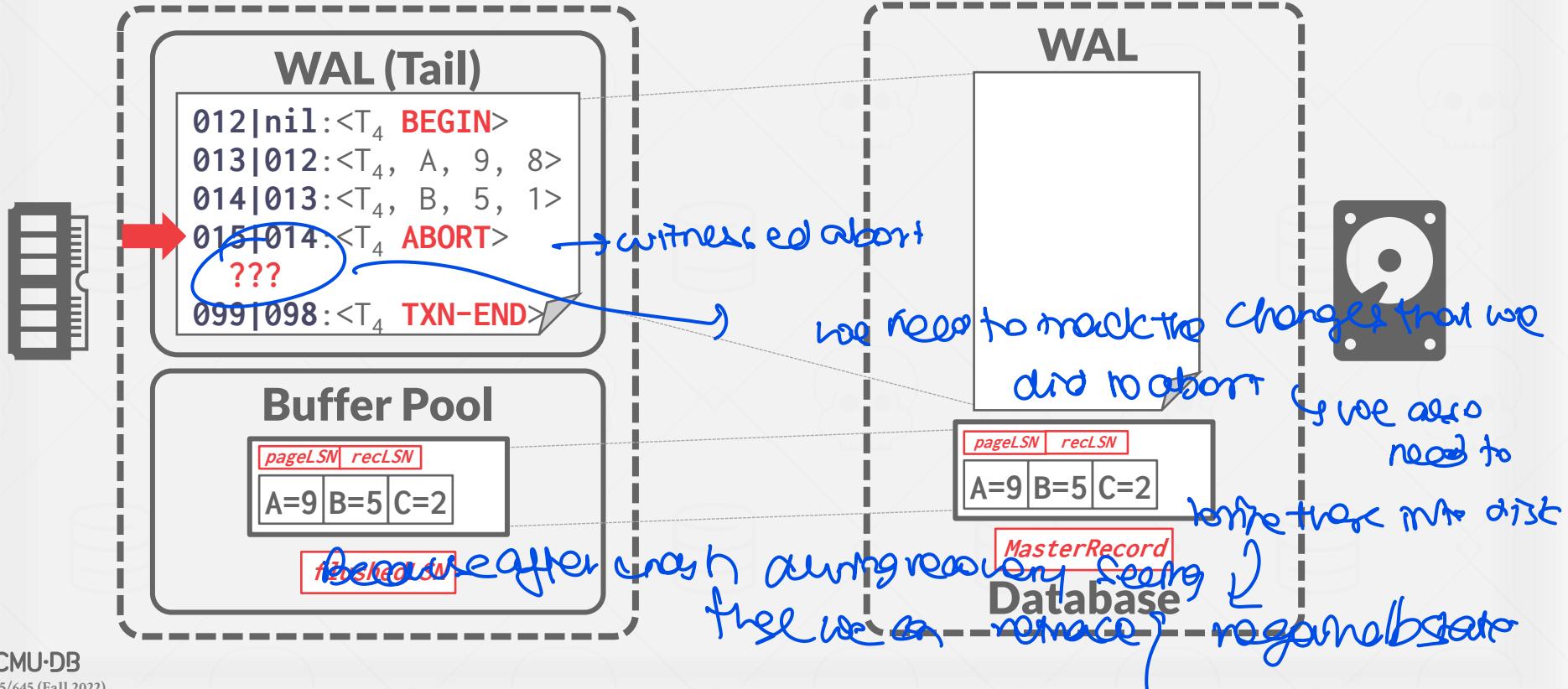
- LSN | prevLSN

how we
maintain

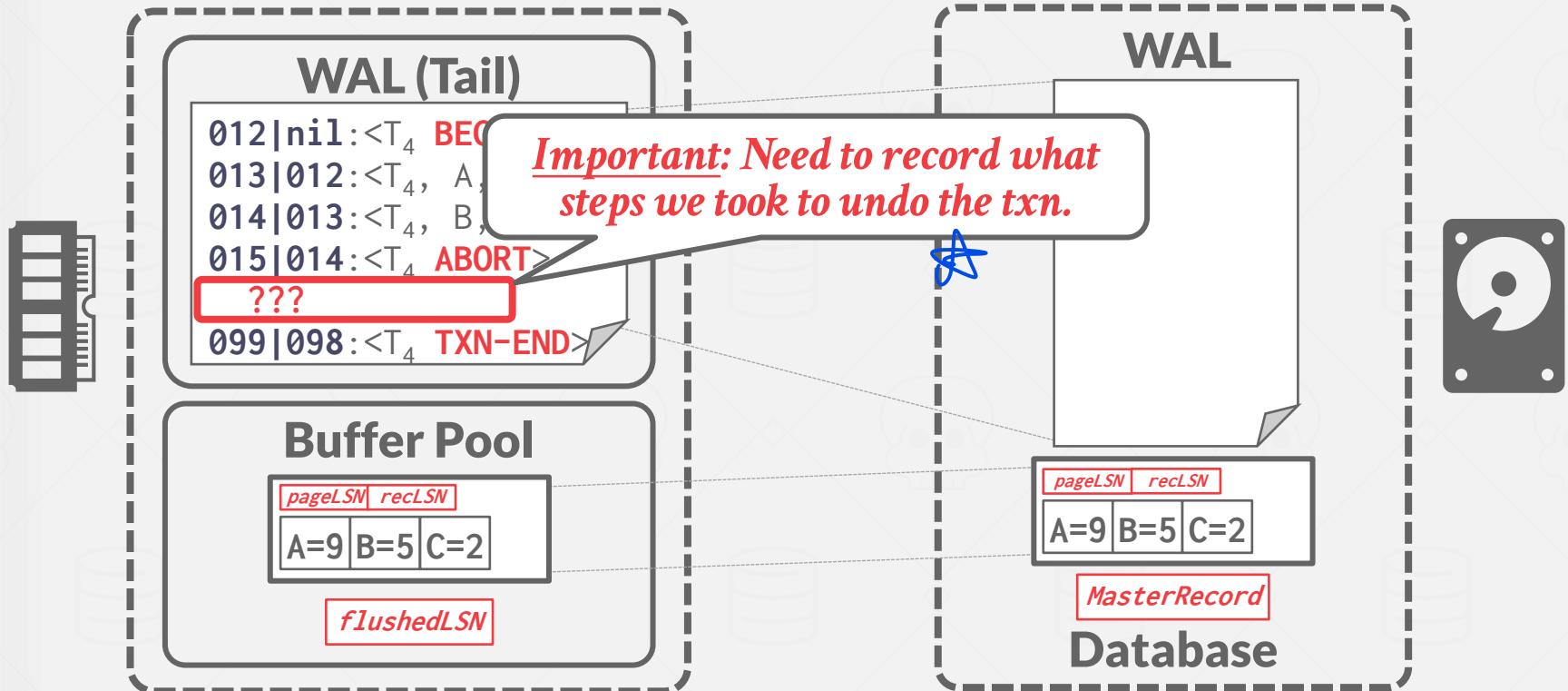
TRANSACTION ABORT



TRANSACTION ABORT



TRANSACTION ABORT



for now
purposive
use CLR

Compensation Log Record (CLR)

COMPENSATION LOG RECORDS

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

~~already affected
or-pending record~~

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

(specifying which one to undo next)

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

CLRs are added to log records but DBMS does

*actions taken
to undo action
of a previous
update
record*

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	-	-	-	-

↴ Seen Abort ↴ Undoing
 ↴ : Create a CR to undo previous updates ↴ Previous update

TRANSACTION ABORT - CLR EXAMPLE

newly introduced

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

↳ A record of CLR
created to
undo this

↳ undone
changes there is
CLR.

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

setting
redo
log

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

Letting
undo
next

The LSN of the next log record to be undone.

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
027	026	T ₁	TXN-END	-	-	-	nil

In this log, we see the end of transaction
records written into the log from the
transaction.

ABORT ALGORITHM

- ① First write an **ABORT** record to log for the txn.
- ② Then analyze the txn's updates in reverse order.
For each update record:
 - Write a **CLR** entry to the log. ✓
 - Restore old value. ↴
- ③ Lastly, write a **TXN-END** record and release locks.

Notice: CLRs never need to be undone.

They will stay in the log / wal record only

TODAY'S AGENDA

Log Sequence Numbers

Normal Commit & Abort Operations

Fuzzy Checkpointing

Recovery Algorithm

Logs can grow forever
to understand until
where we want to
look back we use
Checkpointing .

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 easy.

Bad for perf Good for recovery

Then
flush

dirty
pages

then take
snapshot

We halt everything
when we begin a
txn. If there are
txns that are
updating we
wait until
they're
finished

Slightly
better
way.

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.

Here

and checkpoint

writes all 3 pages

(P1, P2, and P3) to disk.

Then starts to update page 3

but before it updates page 3 and completes the whole update check

point does see

join and transaction half

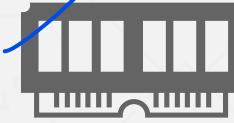
transaction unlock

Transaction

Pause the modification from
writing back to a checkpoint

white marks

Checkpoint



Page #1

Page #2

Page #3

Page #4

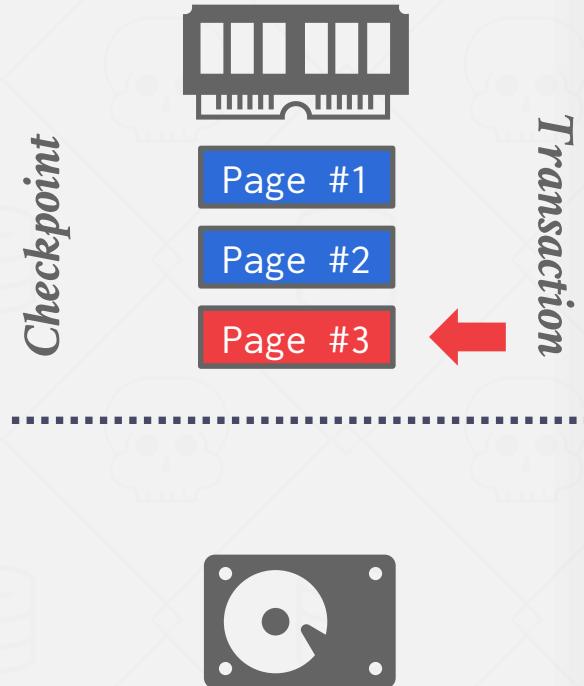


update
page

SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the DBMS takes the checkpoint.

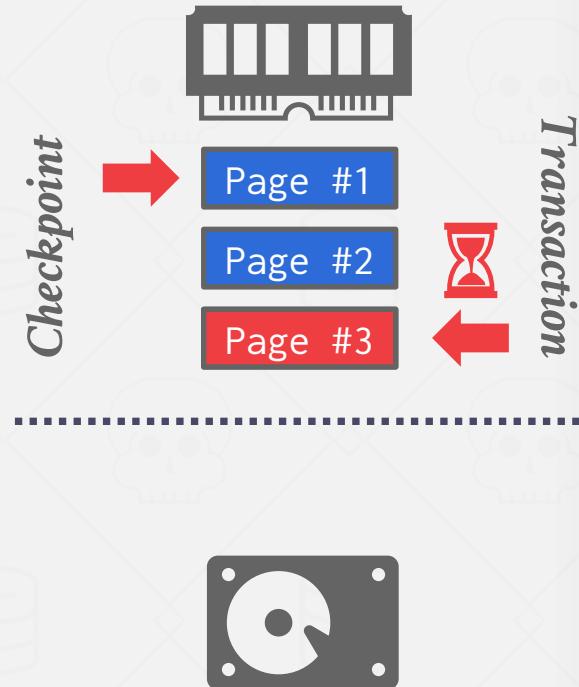
- Prevent queries from acquiring write latch on table/index pages.
- Don't have to wait until all txns finish before taking the checkpoint.



SLIGHTLY BETTER CHECKPOINTS

Pause modifying txns while the DBMS takes the checkpoint.

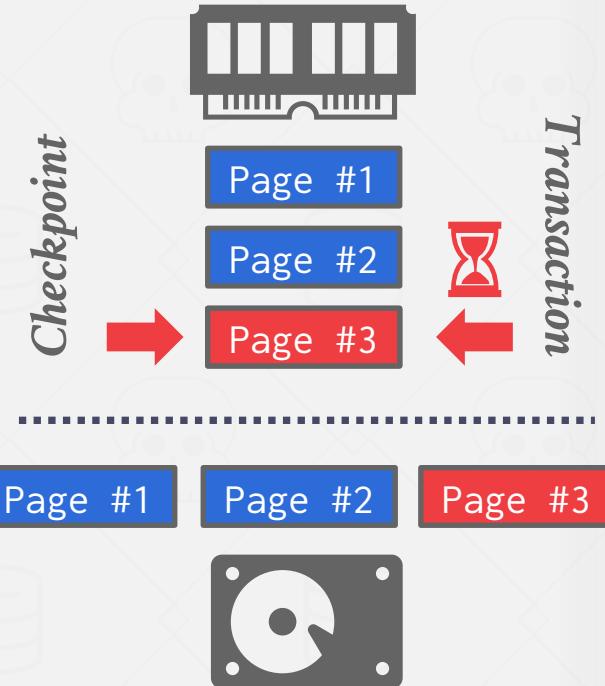
- Prevent queries from acquiring write latch on table/index pages.
- Don't have to wait until all txns finish before taking the checkpoint.



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.

- A
- **Dirty Page Table (DPT)**
- (ATT)



all dirty pages that are present during taking a check point

- This causes
unatomic
update
inconsistent
checkpoint -

update
inconsistent
checkpoint



Page #1

Page #2

Page #3



Transaction

Page #1

Page #2

Page #3



by nonone error

↑ maintains list of active transactions while taking checkpoint

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.

Remove entry after the **TXN-END** record.

Upon
Txn-end
record
remove
this entry

Txn Status Codes:

- **R** → Running
- **C** → Committing
- **U** → Candidate for Undo

Each txn stored as

$\langle \text{txnId}, \text{status}, \text{lastLSN} \rangle$

Can be 3 ways:

R → Running

C → Committing

U → Candidate for Undo

DIRTY PAGE TABLE

Keep track of which pages in the buffer pool contain changes that have not been flushed to disk.

One entry per dirty page in the buffer pool:

→ **recLSN**: The *LSN* of the log record that first caused the page to be dirty.

Each dirty page → recLSN → The LSN of the log record that first caused the page to be dirty.

↳ all dirty pages in buffer pool which haven't been flushed to disk.

SLIGHTLY BETTER CHECKPOINT

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}). *More T_3 is running at the time of checkpoint*

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

*T_2 may well be still running
for T_2 has not been seen yet.*

WAL	
< T_1	BEGIN>
< T_2	BEGIN>
< T_1 ,	$A \rightarrow P_{11}$, 100, 120>
< T_1	COMMIT>
< T_2 ,	$C \rightarrow P_{22}$, 100, 120>
< T_1	TXN-END >
<CHECKPOINT	
→ ATT={ T_2 },	ATT
DPT={ P_{22} } >	ATT
< T_3	BEGIN>
< T_2 ,	$A \rightarrow P_{11}$, 120, 130>
< T_2	COMMIT>
< T_3 ,	$B \rightarrow P_{33}$, 200, 400>
<CHECKPOINT	
→ ATT={ T_2, T_3 },	
DPT={ P_{11}, P_{33} } >	
< T_3 ,	$B \rightarrow P_{33}$, 400, 600>

FUZZY CHECKPOINTS

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

→ No attempt to force dirty pages to disk.

Start of checkpoint

New log records ↑ to track checkpoint boundaries:

- **CHECKPOINT-BEGIN**: Indicates start of checkpoint
- **CHECKPOINT-END**: Contains ATT + DPT.

Till now checkpointing was forcing dirty pages to be written back

↑
To keep track of
the trans &
dirty pages.

Allows active txns to continue running while system writes log records for checkpoint

↓ Our wal page contains these 2 records

FUZZY CHECKPOINT

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.

T_2 was active previously

P_{22} is dirty.

After checkpoint-end
LSN of checkpoint
begin

not written to Master record

WAL
< T_1 BEGIN>
< T_2 BEGIN>
< T_1 , $A \rightarrow P_{11}$, 100, 120>
< T_1 COMMIT>
< T_2 , $C \rightarrow P_{22}$, 100, 120>
< T_1 TXN-END >
<CHECKPOINT-BEGIN> $\xrightarrow{\text{log}}$
< T_3 BEGIN>
< T_2 , $A \rightarrow P_{11}$, 120, 130>
<CHECKPOINT-END
ATT={ T_2 },
DPT={ P_{22} }>
< T_2 COMMIT> $\xrightarrow{\text{T2 comm}}$
< T_3 , $B \rightarrow P_{33}$, 200, 400>
<CHECKPOINT-BEGIN>
< T_3 , $B \rightarrow P_{33}$, 10, 12> $\xrightarrow{\text{T3 txn-end}}$
<CHECKPOINT-END
ATT={ T_2, T_3 },
DPT={ P_{11}, P_{33} }>

FUZZY CHECKPOINT

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.

WAL

```

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

```

ARIES - RECOVERY PHASES

↗ Start at latest checkpoint (from WAL records)

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.

After the
3 phases

↓
consistent
db state

create the ATT & DPT

containing all
active txns and
dirty pages.

Phase #2 – Redo

- Repeat all actions starting from an appropriate point in the log (even txns that will abort).

→ Redo all the txns that are committed
before the crash
Rewrite them into disk

Phase #3 – Undo

- Undo changes of txns that committed before the crash.

→ Undo changes of txns that committed before the crash.

have not

before the
crash.

ARIES - OVERVIEW

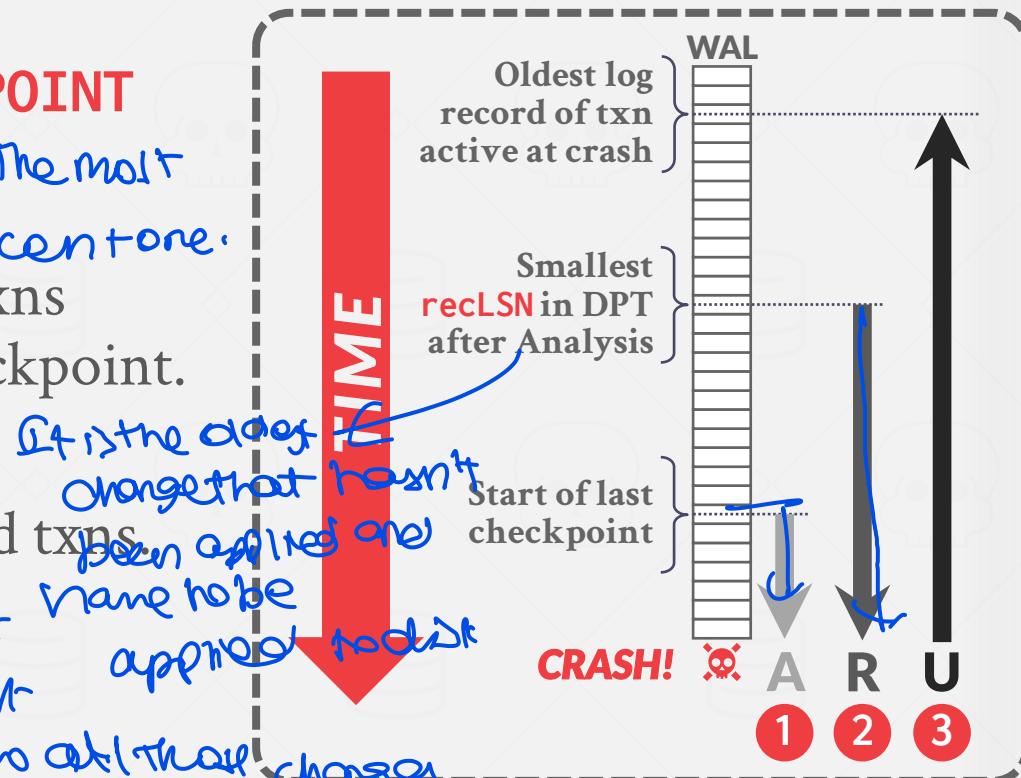
Start from last **BEGIN-CHECKPOINT**
found via **MasterRecord**. ↗ The most
↑ Start of last checkpoint from recent one.

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

Redo: Repeat all actions.

Undo: Reverse effects of failed txns

↓
Grant at the end find the oldest log record of txn that is active at crash Undo all those changes



If seen after end then

remove corr. txn from ATT.

*If commit,
then change*

*status to
commit*

*If a txn
not in ATT*

*add it to
undo
one.*

ANALYSIS PHASE

from last

*succ. check
point*

*to end of
final log
records.*

Scan log forward from last successful checkpoint.

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 status **UNDO**.
- On commit, change txn status to **COMMIT**.

For update log records:

- If page **P** not in DPT, add **P** to DPT, set its **recLSN=LSN**.

↳ if a page not in DPT

*Add it to DPT
Add LSN
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.

At the end ATT → all the txns that are active at the time of crash

DPT, all the dirty pages that haven't been written back to disk

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		
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!

```

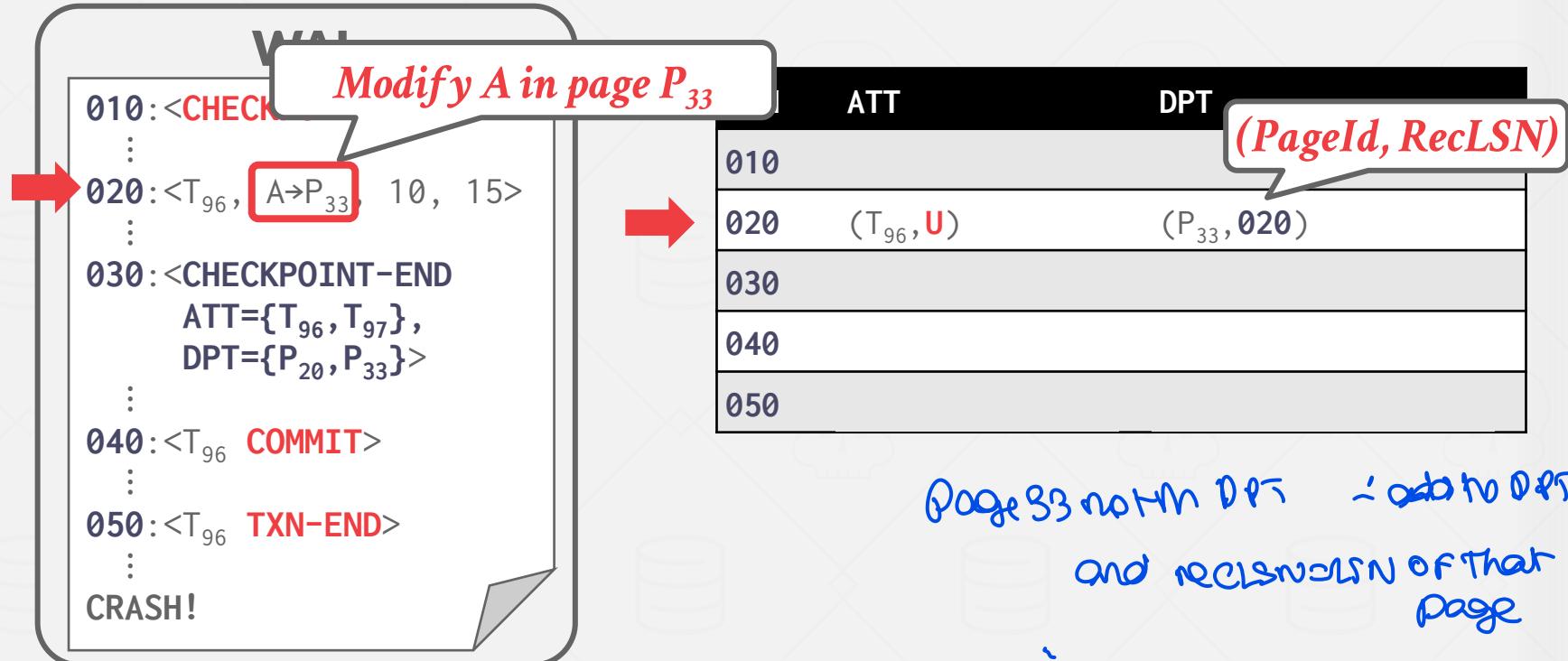


LSN	ATT	DPT
010		
020	(T ₉₆ , U)	
030		
040		
050		

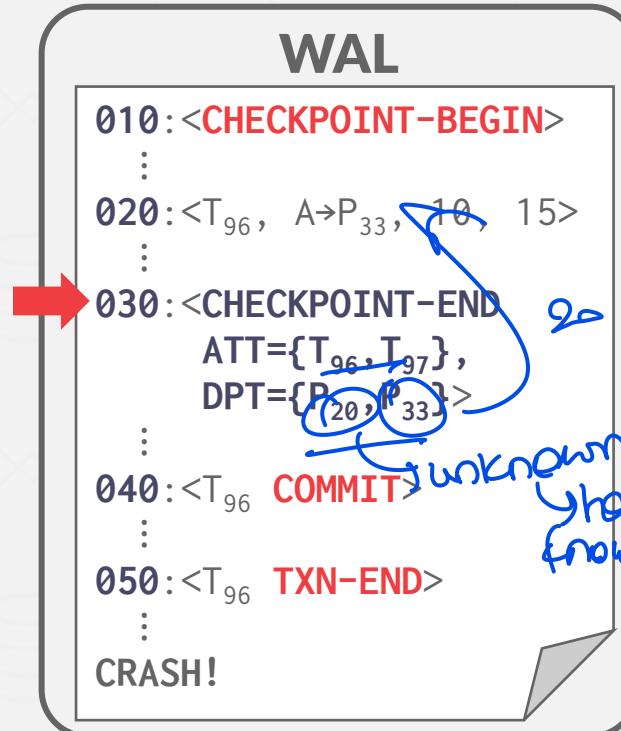
(TxnId, Status)

T₉₆ is not in ATT ∴ add it to
ATT and undo

ANALYSIS PHASE EXAMPLE



ANALYSIS PHASE EXAMPLE



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

Checkpoint end: odd contents in ATT and DPT

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!

```

at LSN 40

→ Take commit

Therefore change status to Commit

LSN	ATT	DPT
010		
020	(T ₉₆ , U)	(P ₃₃ , 020)
030	(T ₉₆ , U), (T ₉₇ , U)	(P ₃₃ , 020), (P ₂₀ , 008)
040	(T ₉₆ , C), (T ₉₇ , U)	(P ₃₃ , 020), (P ₂₀ , 008)
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!

```



Dirty page names.

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

→ Tag seen Txn end at LSN 50 - remove its entry from ATT

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 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, or
- Affected page is in DPT but that record's **LSN** is less than the page's **recLSN**.

*if aff not in DPT → changes have been applied
record's LSN > page recLSN'*

REDO PHASE

To redo an action:

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

:To redo an action.

Set pageLSN to record's LSN.

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

at the end of

*redo
phase*

*write TXN-END for
actions with status C*

*? remove them
from 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.

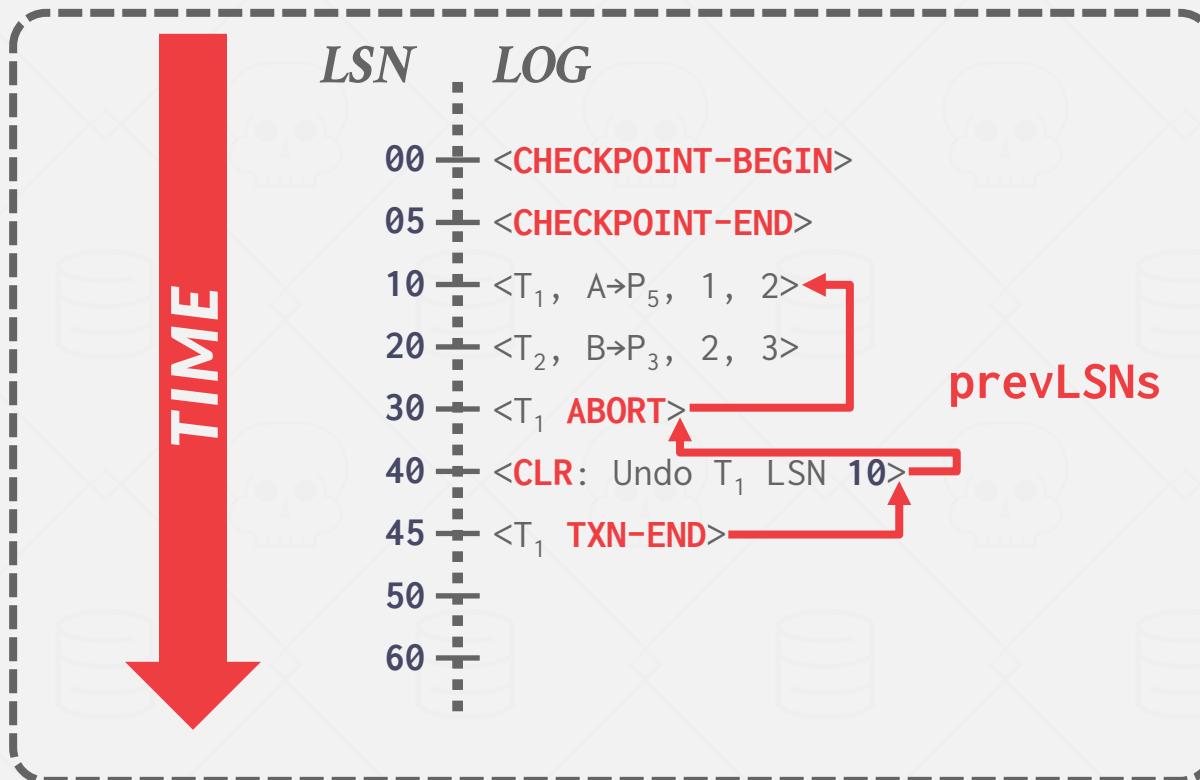
Process them in reverse *LSN* order using the **lastLSN** to speed up traversal.

Write a **CLR** for every modification.

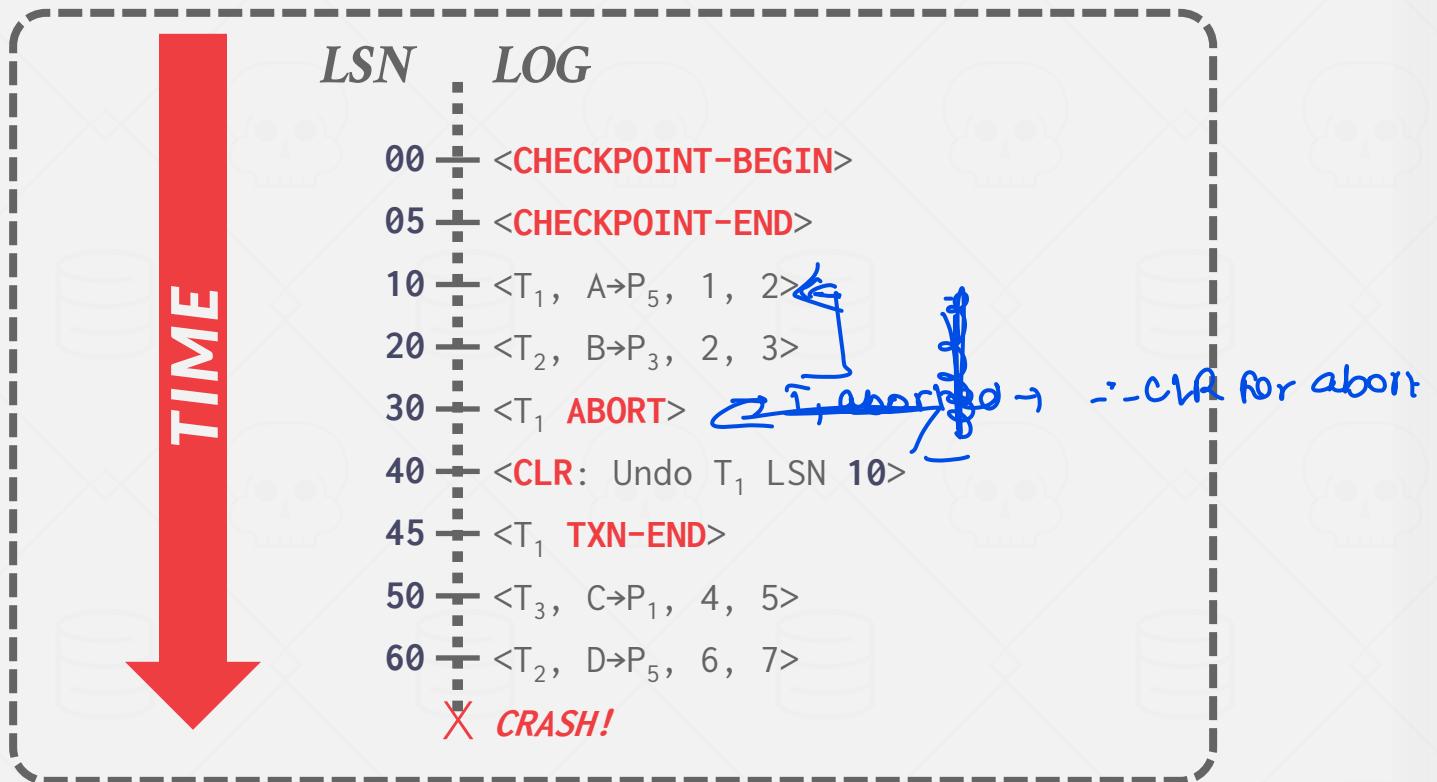
Read it in
reverse order
op2n to
Speed Up

So LSN with greater
number is
first rolled back

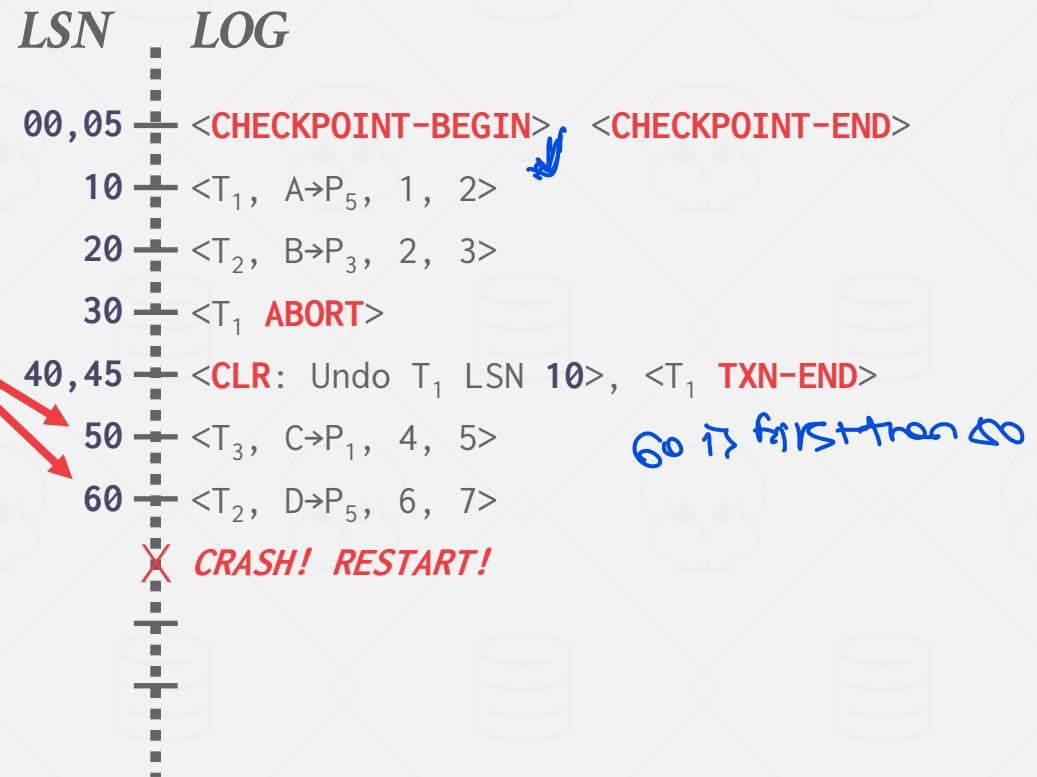
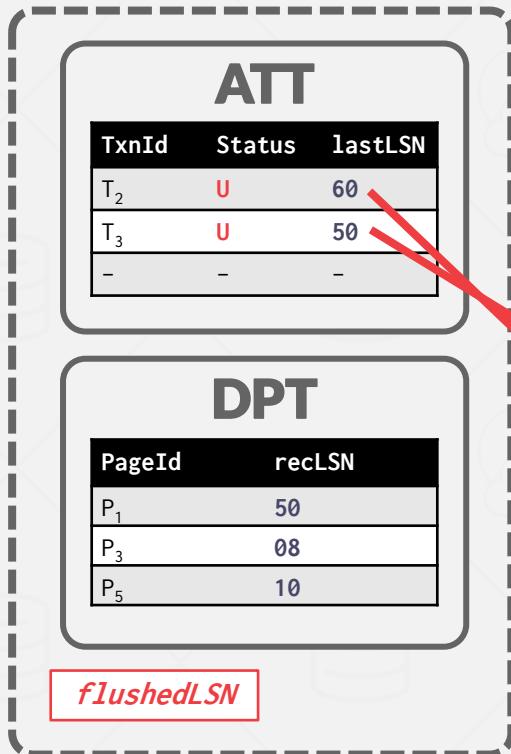
FULL EXAMPLE



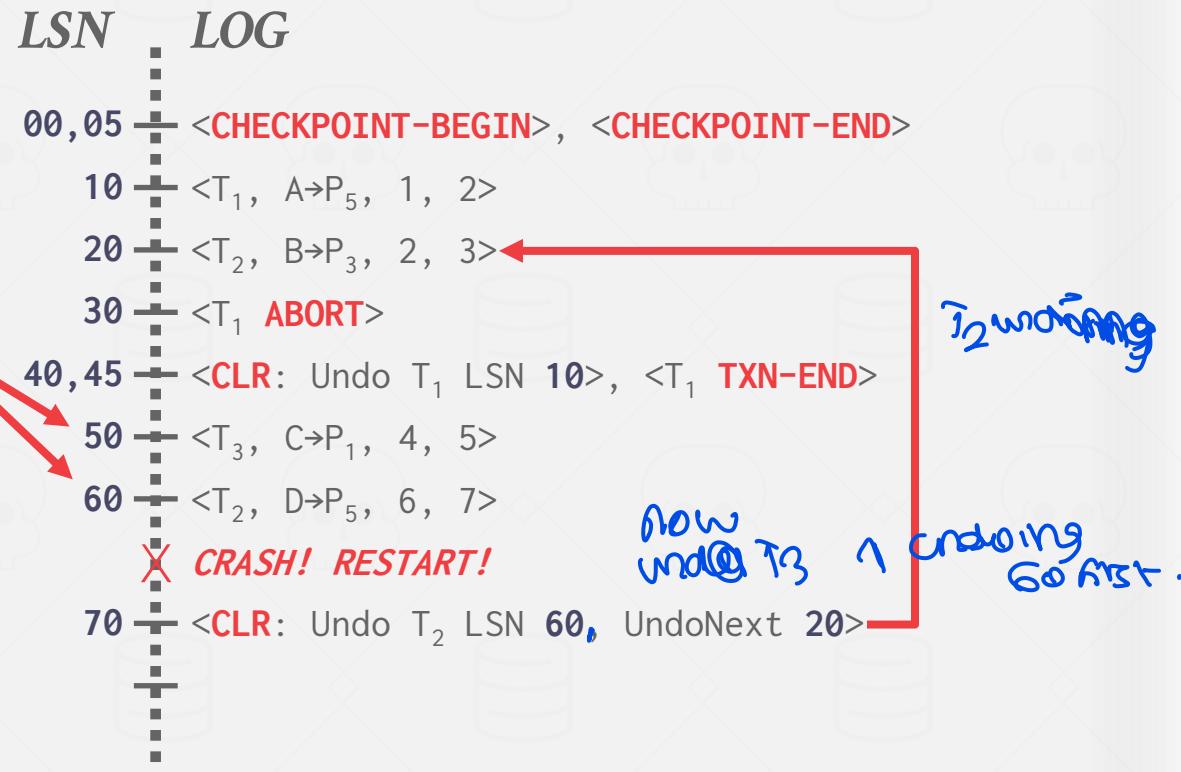
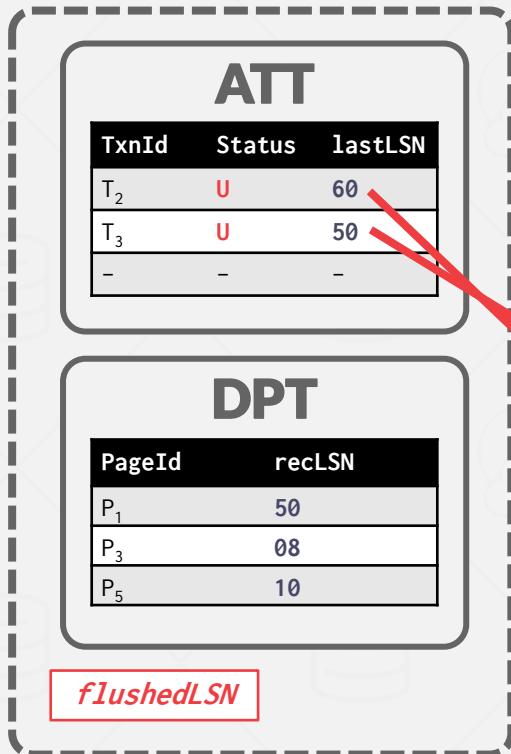
FULL EXAMPLE



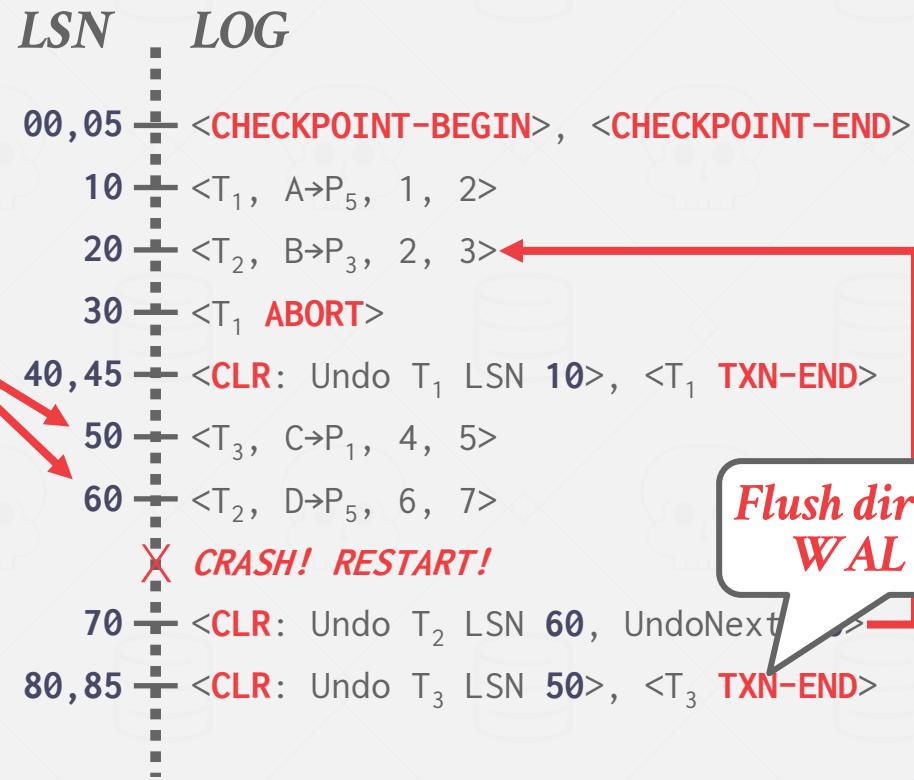
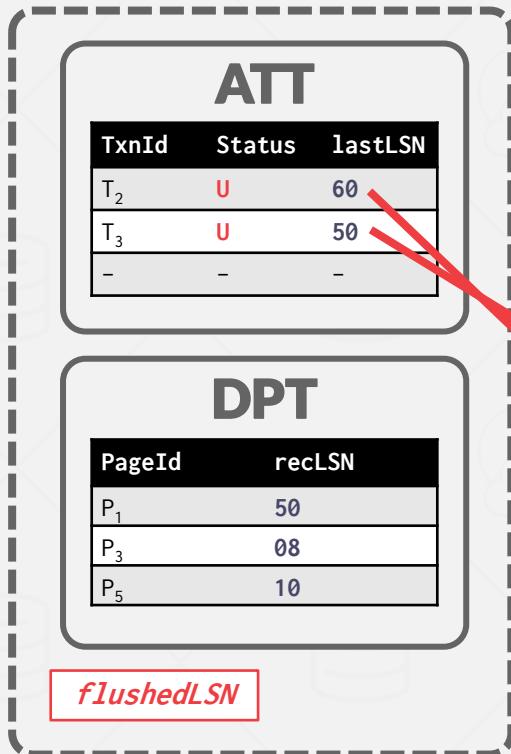
FULL EXAMPLE



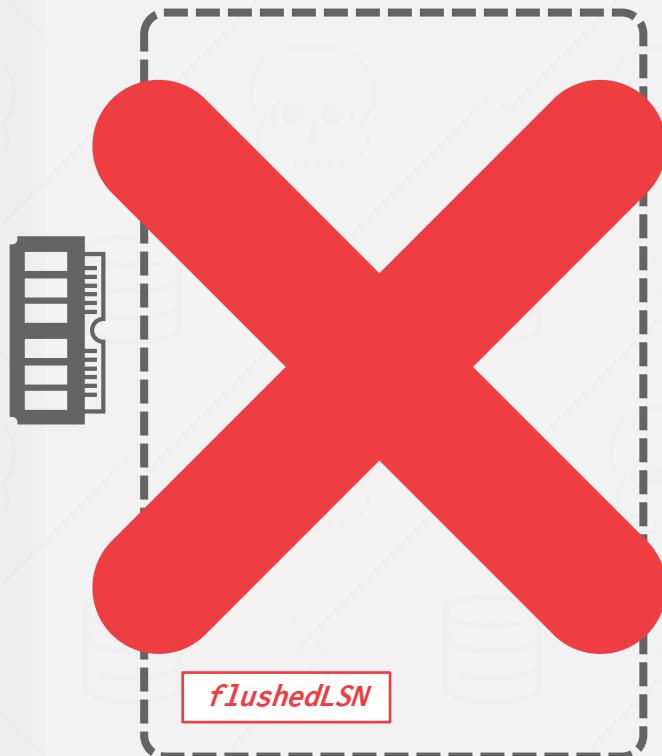
FULL EXAMPLE



FULL EXAMPLE



FULL EXAMPLE



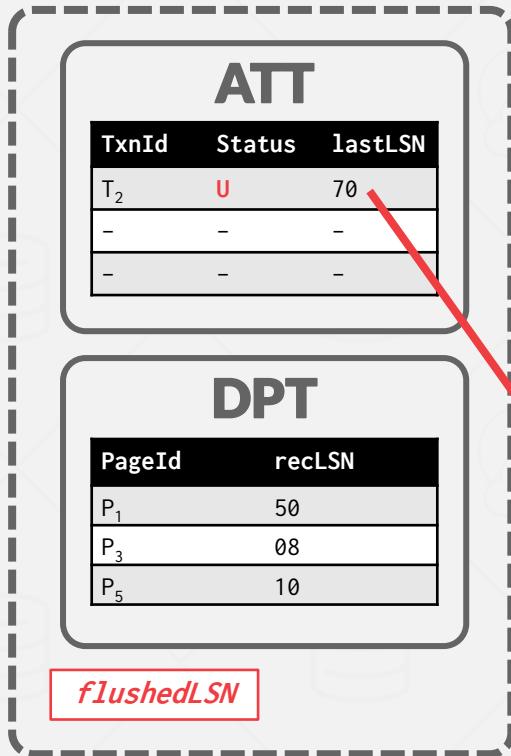
<i>LSN</i>	<i>LOG</i>
00, 05	<CHECKPOINT-BEGIN>, <CHECKPOINT-END>
10	<T ₁ , A→P ₅ , 1, 2>
20	<T ₂ , B→P ₃ , 2, 3>
30	<T ₁ ABORT>
40, 45	<CLR: Undo T ₁ LSN 10>, <T ₁ TXN-END>
50	<T ₃ , C→P ₁ , 4, 5>
60	<T ₂ , D→P ₅ , 6, 7>
X	CRASH! RESTART!
70	<CLR: Undo T ₂ LSN 60, UndoNext>
80, 85	<CLR: Undo T ₃ LSN 50>, <T ₃ TXN-END>
X	CRASH! RESTART!

Flush dirty pages + WAL to disk!

pay

at this point been Txn flushed

FULL EXAMPLE



LOG

LSN	Log Entry
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>
70	X CRASH! RESTART!
70	<CLR: Undo T ₂ LSN 60, UndoNext 20> ←
80, 85	<CLR: Undo T ₃ LSN 50>, <T ₃ TXN-END>
90, 95	X CRASH! RESTART!
90, 95	<CLR: Undo T ₂ LSN 20>, <T ₂ TXN-END>

Another crash.

ADDITIONAL CRASH ISSUES (1)

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

- Nothing. Just run recovery again.

*Start from
beginning 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!