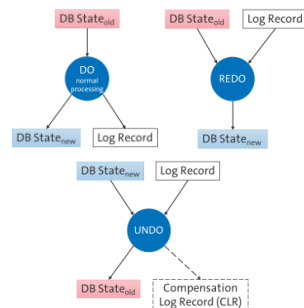# Summary

- Architectures of Database Systems
- **Transaction Management**
- Modern Database Technology
- Data Warehouses and OLAP
- Data Mining
- Big Data Analytics

**Transaction Failure** | **System Failure**
**Device Failure** | **Disaster**

**R1 Recovery** partial undo
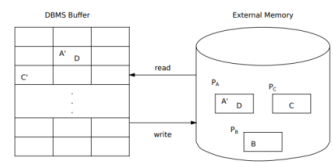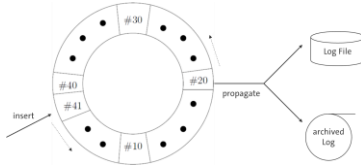
**R4 Recovery** (snapshot +) global redo

**R2 Recovery** Committed transactions (partial redo)
**R3 Recovery** Uncommitted transactions (global undo)

DB State$_{old}$ — DO normal processing — DB State$_{new}$, Log Record

DB State$_{old}$, Log Record — REDO — DB State$_{new}$

DB State$_{new}$, Log Record — UNDO — DB State$_{old}$, Compensation Log Record (CLR)

Example

```
#    T₁              T₂           Log [LSN, TA, PageID, Redo, Undo, PrevLSN]
1.   BOT                          [#1, T₁, BOT, 0]
2.   r(A,a₁)
3.                   BOT          [#2, T₂, BOT, 0]
4.                   r(C,c₂)
5.   a₁:=a₁-50
6.   w(A,a₁)                      [#3, T₁, P_A, A-=50, A+=50, #1]
7.                   c₂:=c₂+100
8.                   w(C,c₂)      [#4, T₂, P_C, C+=100, C-=100, #2]
9.   r(B,b₁)
10.  b₁:=b₁+50
```

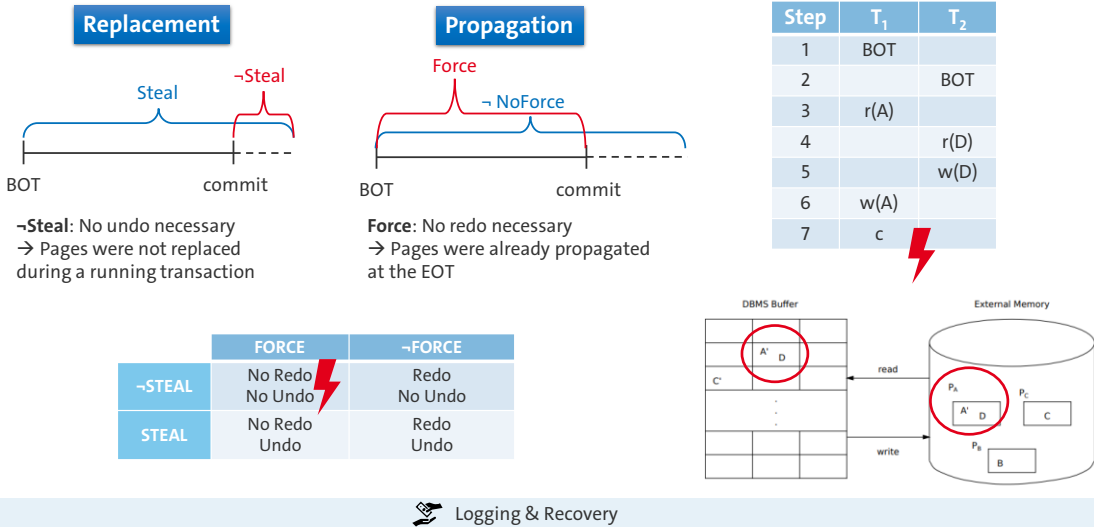DBMS Buffer — External Memory

read / write

Log File / archived Log

propagate / insert

**Replacement, Propagation, and Insertion Strategies**

**Recap replacement, propagation, and insertion strategies**

Replacement — Propagation — Insertion

Direct
$P_A \rightarrow P_A$

Indirect with Twin-Block
$P_A$ state i / state i+1 $\rightarrow P_A$

¬Steal — Steal
BOT — commit

Force — ¬ NoForce
BOT — commit

**¬Steal**: No undo necessary
→ Pages were not replaced during a running transaction

**Force**: No redo necessary
→ Pages were already propagated at the EOT

Logging & Recovery

Force writes changes at commit time at the latest → No indirect insertion strategy possible

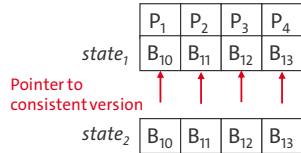## Recap replacement, propagation, and insertion strategies

**Replacement**

Steal
¬Steal

BOT — commit

**¬Steal**: No undo necessary
→ Pages were not replaced during a running transaction

**Propagation**

Force
¬ NoForce

BOT — commit

**Force**: No redo necessary
→ Pages were already propagated at the EOT

| | FORCE | ¬FORCE |
|---|---|---|
| **¬STEAL** | No Redo No Undo | Redo No Undo |
| **STEAL** | No Redo Undo | Redo Undo |

| Step | $T_1$ | $T_2$ |
|---|---|---|
| 1 | BOT | |
| 2 | | BOT |
| 3 | r(A) | |
| 4 | | r(D) |
| 5 | | w(D) |
| 6 | w(A) | |
| 7 | c | |

DBMS Buffer

A' D
C'

read
write

External Memory

$P_A$  A' D   $P_C$  C
$P_B$  B

Logging & Recovery

- Changes of $T_1$ on page $P_A$ are forced to propagate at commit in step 7 because of FORCE
- $T_2$ has not committed, yet → Because of ¬STEAL, page $P_A$ (which includes element D that $T_2$ is working on) must not be replaced while $T_2$ is still running
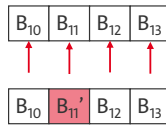
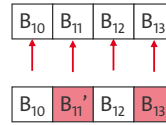# Recap Indirect Insertion with Redundancy (Twin-Block)
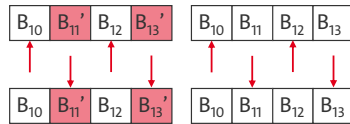
Mapping of pages to blocks:

Operations:

$T_1$: $w(P_2)$

$T_1$: $w(P_4)$

| $P_1$ | $P_2$ | $P_3$ | $P_4$ |
|---|---|---|---|
| $B_{10}$ | $B_{11}$ | $B_{12}$ | $B_{13}$ |

Pointer to consistent version

$state_1$

$state_2$

| $B_{10}$ | $B_{11}$ | $B_{12}$ | $B_{13}$ |
|---|---|---|---|

$T_1$: Commit

$T_1$: Abort

Current state: Like initial state
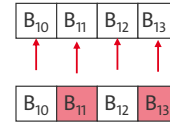
- Copy $B_{11}'$
- Switch pointer

Current state: Like initial state but with switched pointer

- Copy $B_{11}$
- Do NOT switch pointer

UHH Universitä
DER FORSCHUNG | DER LEHRI

5

**Recap Indirect Insertion with Redundancy (Twin-Block)**

Operations:

$T_1: w(P_2)$

| $B_{10}$ | $B_{11}$ | $B_{12}$ | $B_{13}$ |
|---|---|---|---|

| $B_{10}$ | $B_{11}'$ | $B_{12}$ | $B_{13}$ |
|---|---|---|---|

$T_1: w(P_4)$

| $B_{10}$ | $B_{11}$ | $B_{12}$ | $B_{13}$ |
|---|---|---|---|

| $B_{10}$ | $B_{11}'$ | $B_{12}$ | $B_{13}'$ |
|---|---|---|---|

$T_2: w(P_4)$

| $B_{10}$ | $B_{11}$ | $B_{12}$ | $B_{13}$ |
|---|---|---|---|

| $B_{10}$ | $B_{11}'$ | $B_{12}$ | $B_{13}''$ |
|---|---|---|---|

Can cause Lost Update
→ Wait for $T_1$ to commit first

$T_2: r(P_4)$

| $B_{10}$ | $B_{11}$ | $B_{12}$ | $B_{13}$ |
|---|---|---|---|

| $B_{10}$ | $B_{11}'$ | $B_{12}$ | $B_{13}'$ |
|---|---|---|---|

$T_2$ reads the consistent version

**Exercise**

| $P_1$ | $P_2$ | $P_3$ | $P_4$ |
|---|---|---|---|
| $B_{10}$ | $B_{11}$ | $B_{12}$ | $B_{13}$ |

Assuming each page can only hold a single integer number and $P_1 = P_2 = P_3 = P_4 = 1$. What would the Twin Block look like after the following operations?

T1: w($P_2$ =5)
T1: r($P_1$)
T2: r($P_2$)
T1: A= $P_1$+3
T1: w($P_3$=A)
T2: w($P_4 = P_2$)
T1: Commit
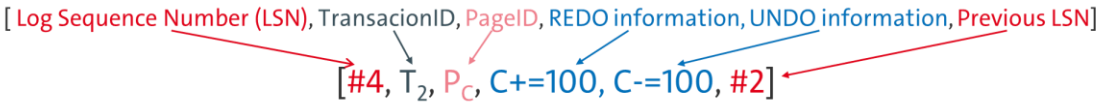T2: Commit

6

Logging & Recovery

**Difference to Shadow Storage**
- Shadow storage only creates duplicate if data is changed, i.e. here only for $B_{11}$ and $B_{13}$
- Advantage: Less memory consumption
- Disadvantage: Fragmentation (addresses are scattered)

## Log Data

[ Log Sequence Number (LSN), TransacionID, PageID, REDO information, UNDO information, Previous LSN]

$$[\#4, T_2, P_C, C\mathrel{+}=100, C\mathrel{-}=100, \#2]$$

**Rules for writing Log data**

- **Write Ahead Log (WAL)** principle for UNDO information
  → Log entries must be written to file before an affected page is replaced
- **Force-Log-at-Commit**
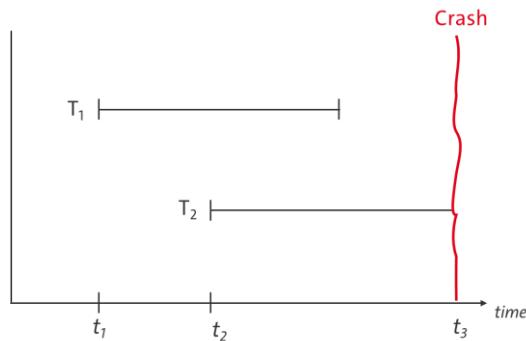  → Redo information must be written before the commit

**WAL**
- Only relevant for STEAL → NOSTEAL does not replace pages during transaction
- Important for direct insert strategies

**Force-Log-at-Commit**
- Required for Crash-Recovery with NOFORCE (i.e. when changes might not have been propagated)
- Required for R4 recovery (with FORCE and NOFORCE) → even if changes were propagated, they were lost
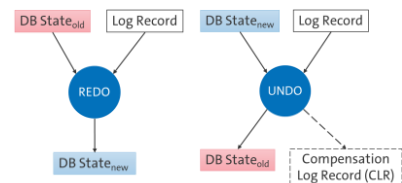- For direct and indirect insertion strategies

Recovery

Crash

$T_1$

$T_2$

$t_1$ $t_2$ $t_3$ time

Transaction $T_1$ is a winner → Redo
Transaction $T_2$ is a loser → Undo

**Phases of Recovery**

1. **Analysis**
   → Identify winners and losers
2. **Redo**
   → Repetition of history
   → Reads log file forwards
3. **Undo** of losers
   → Reads log file backwards
   → Write CLR

DB State_old | Log Record | DB State_new | Log Record
REDO | UNDO
DB State_new | DB State_old | Compensation Log Record (CLR)

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Logging & Recovery

10

**Analysis**
- Started transactions marked by BOT
- WAL important to not lose this log information in case of a crash, e.g. to not lose a commit that happened but was not entered in the log file
- Winners: Transactions with commit entry
  → T1 is a winner
- Losers: Transactions without commit entry
  → T2 is a loser

**REDO**
- For winners and losers
- Go forward in log file
- Fetch referenced page from secondary storage into buffer
- Compare LSN of page with LSN in log (see slide 22 in last lecture)
  - $LSN_{page} >= LSN_{log}$: Ok
  - $LSN_{page} < LSN_{log}$: Redo and update LSN of page
- Propagate page

**UNDO**

- For all loser transactions irrespective of the LSN
- Go backwards through log file
- Skip all entries of winner transactions
- For all loser transactions:
  - Fetch referenced page from secondary storage into buffer
  - Execute undo operation
  - Write CLR
- Propagate Page

**Fault Tolerance of Recovery**

**How to make sure that another crash during the recovery does not change the result of the recovery?**
→ REDO and UNDO phase must be idempotent (result must always be the same irrespective of how often the operation has been applied)

Idempotence of REDO
LSN prevents repeated REDO

Idempotence of UNDO
Compensation Log Record for every UNDO operation

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Logging & Recovery

11

For each operation (that changes data), the following must be fulfilled:
- undo(undo(...undo(a))...)) = undo(a)
- redo(redo(...redo(a))...)) = redo(a)

An object *a* with an operator ○ that fulfills the property *a* ○ *a* = *a* is called idempotent in respect of operator ○.

# Compensation Log Records

No undo information

<LSN, TransactionID, PageID, Redo information, PrevLSN, UndoNxtLSN>

<#7', $T_2$, PA, A+=100, #7, #4>

> **LSN**: Modified version of according log entry used for undo
> **Redo information**: Successful undo during recovery
> **PrevLSN**: Previous log entry, can be a CLRs
> **UndoNxtLSN**: Referes to the next change that has to be undone

12

Logging & Recovery

**Recovery Example**

[#1, T$_1$, **BOT**, 0]
[#2, T$_2$, **BOT**, 0]
[#3, T$_1$, A, A-=50, A+=50, #1]
[#4, T$_2$, C, C+=100, C-=100, #2]
[#5, T$_1$, B, B+=50, B-=50, #3]
[#6, T$_1$, **commit**, #5]
[#7, T$_2$, A, A-=100, A+=100, #4]

**Crash**

<#7', T2, A, A+=100, #7, #4>
<#4', T2, C, C-=100, #7', #2>
<#2', T2, -, -, #4', 0>

**Analysis**
Winner: T$_1$
Loser: T$_2$

**Redo**
IF LSN(A) < 3 THEN A-=50 (and replace LSN(A))
IF LSN(C) < 4 THEN C+=100 (and replace LSN(C))
IF LSN(B) < 5 THEN B+=50  (and replace LSN(B))
IF LSN(A) < 7 THEN A-=100 (and replace LSN(A))

**Undo** → Only for losers (T$_2$)
Entry #7
  • A+=100
  • Write CLR
  • LSN(A) = #7'
Entry #4
  • C-=100
  • Write CLR
  • LSN(C) = #4'
Entry #2
  • Write CLR

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Logging & Recovery

13

- No operation done in BOT and EOT → Only CLR entry required in undo phase
- If checkpoint exists: Redo can start at state of checkpoint, else Redo starts at start of DB

# Recovery Exercise

[#1, T2, BOT, 0]

[#2, T2, B, B=B-1, B=B+1, #1]

[#3, T1, BOT, 0]

[#4, T1, A, A=A+2, A=A-2, #3]

[#5, T1, A, A=A*5, A=A/5, #4]

**Crash**

**Analysis**
Winners: ?
Losers: ?

**Redo**
IF LSN(?) < ? THEN ?
…

**Undo**
For all losers:
    Undo operation
    CLR entry:

        <LSN, TransactionID, PageID, Redo information, PrevLSN, UndoNxtLSN>

            <#7', T$_2$, PA, A+=100, #7, #4>

    New LSN

Universität Hamburg
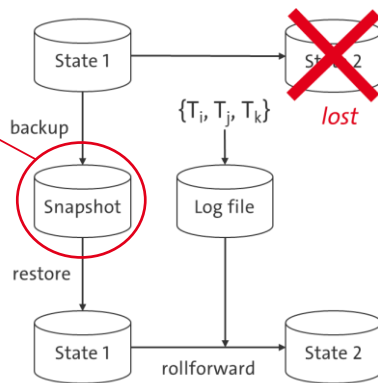DER FORSCHUNG | DER LEHRE | DER BILDUNG

## Checkpoints

**Reduces the amount of REDO operations during recovery after a system failure (or worse)**

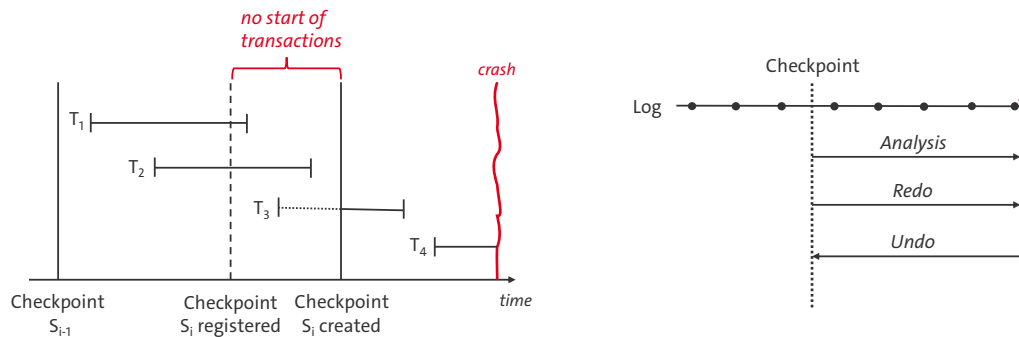**Types of checkpoints**
- (Global) transaction consistent save points
- Action consistent save points
- Fuzzy save points

State 1

backup

{$T_i$, $T_j$, $T_k$}    State 2 *lost*

Snapshot    Log file

restore

State 1    rollforward    State 2
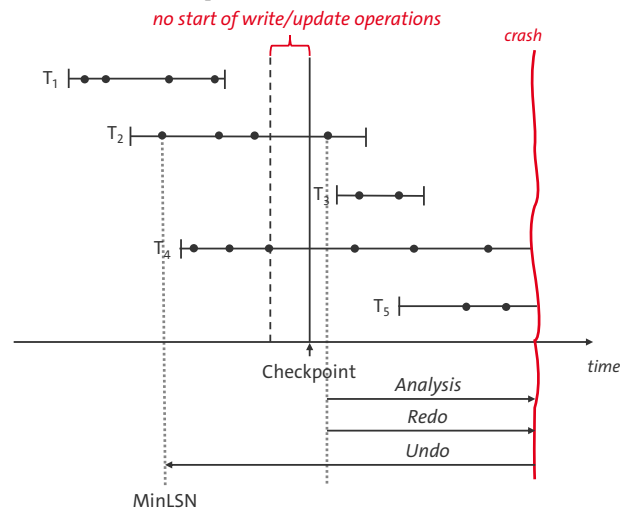
16

Logging & Recovery

- Without checkpoints, all changes since the start of the database must be redone in case of system failure
- Critical: Changes in hot spot pages (pages which are always held in the buffer because there are many operations using that page)

**Transaction Consistent Checkpoints**

*no start of transactions*

*crash*

Checkpoint

Log

*Analysis*

*Redo*

*Undo*

$T_1$

$T_2$

$T_3$

$T_4$

*time*

Checkpoint $S_{i-1}$

Checkpoint $S_i$ registered

Checkpoint $S_i$ created

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

17

Logging & Recovery

- Also called "commit consistent"
- Active transactions (when checkpoint is registered) are finished, new transactions must wait until checkpoint was created
- All modified pages and log information are written to secondary memory → Checkpoint contains transaction consistent state (i.e. no half-finished transactions)
- Expensive because of wait times for transactions (=system down time)→ Not suited for frequent creation
- Log file requires entries indicating when the checkpoint was registered and created:
  - BEGIN_CHKPT Record
  - END_CHKPT Record
- Log Address of last Checkpoint Record is kept in special system file


- "Transaction Oriented Checkpoint" → Just a different way of stating that FORCE was used
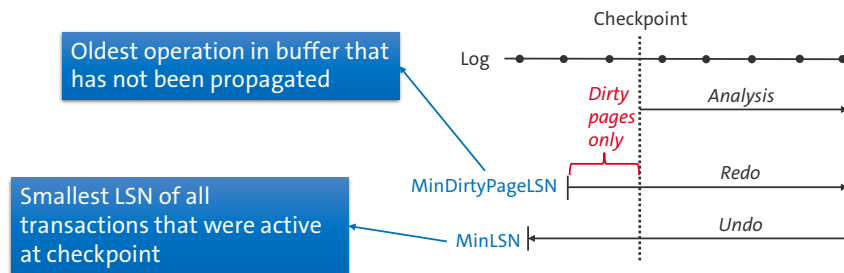  - → Changes are always propagated at the end of transaction
  - → No redo necessary

17

- Before creating a checkpoint, all operations changing any data must be finished (i.e. any write operations)
- All modified pages are written to secondary memory
- During recovery:
    - No redo information older than checkpoint is required
    - Undo information older than checkpoint still required → back to MinLSN
    - MinLSN = LSN of oldest operation of a transaction that was running during the creation of the checkpoint (here: first operation of $T_2$)
- Shorter "down time" but potentially more Undo operations in case of recovery

**Fuzzy Checkpoints**

- Modified pages are not written, only PageID is written

Checkpoint

Log

Oldest operation in buffer that has not been propagated

*Dirty pages only*

*Analysis*

MinDirtyPageLSN

*Redo*

Smallest LSN of all transactions that were active at checkpoint

*Undo*

MinLSN

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

19

Logging & Recovery

Efficiency of recovery depends on buffer management
→ If hot spot pages remain in the buffer and are never propagated, the redo phase must consider the whole log file since fetching these pages
→ Possible solution: Force propagation if a dirty page is part of two consecutive fuzzy checkpoints and not propagated between these two checkpoints