# Recovery

## R&G Chapter 18

(slides adapted from content by J.Gehrke, J.Shanmugasundaram, and/or C.Koch)

1

# A.C.I.D.

- **Atomicity**: <u>All</u> actions in a transaction happen, <u>or none</u> happen.

- **Consistency**: If the transaction maintains consistency, and the DB <u>starts consistent</u>, then the database <u>ends consistent</u>.

- **Isolation**: The execution of one transaction is isolated from all other transactions.

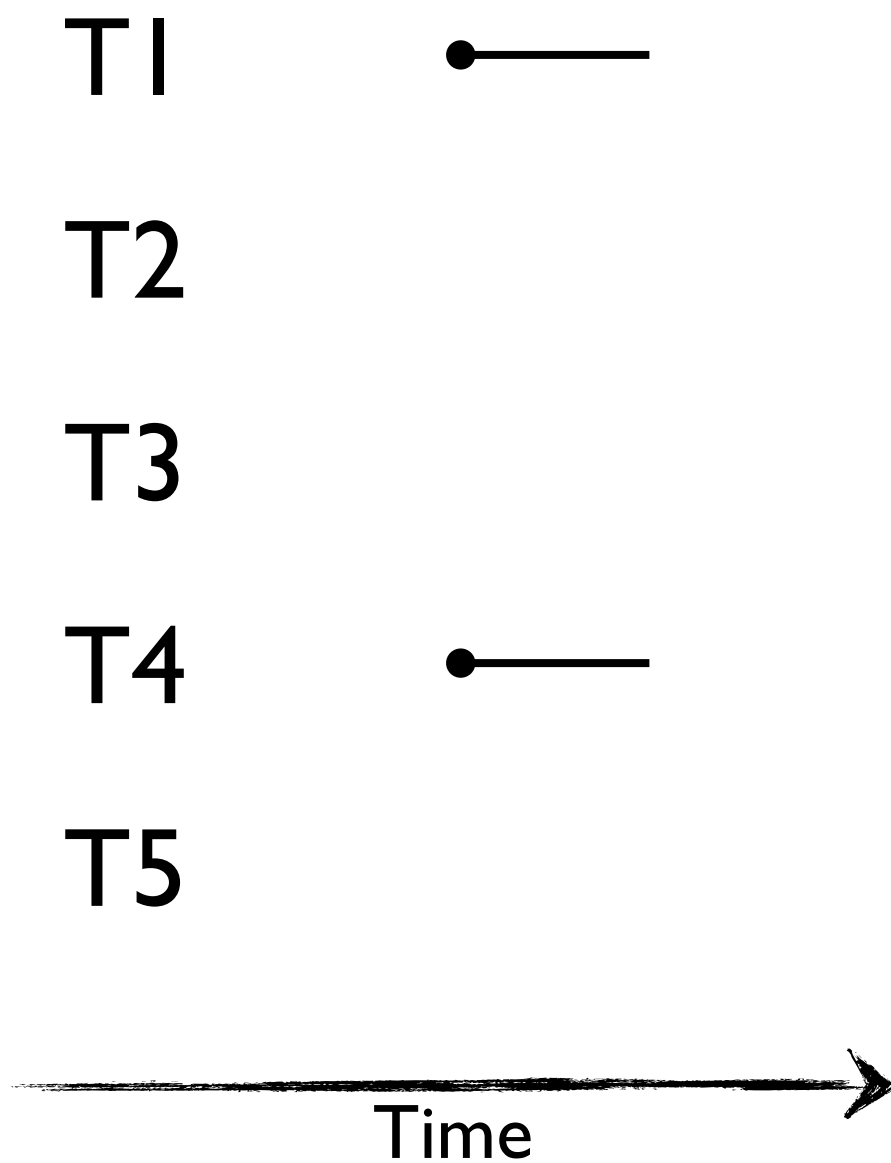- **Durability**: If a transaction commits, its effects persist.

Friday, March 22, 13

# Motivation

T1

T2

T3

T4

T5

→ Time

3

# Motivation

# Motivation



T1

T2

T3

T4

T5

Time

3

# Motivation



T1

T2

T3

T4

T5

Time

3

# Motivation

T1

T2

T3

T4

T5

Time

# Motivation



T1

T2

T3

T4

T5

Time

3

# Motivation

# Motivation

T1

T2

T3

T4

T5

Time

CRASH!

3

# Motivation



Committed Transactions.
These should be present when the DB restarts.

T1

T2

T3

T4

T5

Time

Image copyright: Wikimedia Commons

# Motivation

Committed Transactions.
These should be present when the DB restarts.

CRASH!

T1 •————————•

T2 •————————•

T3 •————————•

T4 •——————————————•

T5 •——————————•

Time

Uncommitted Transactions.
These should leave no trace

Image copyright: Wikimedia Commons

# Goals

- Support <u>recovery</u> from crashes

- Ensure that the effects of <u>committed</u> transactions are <u>recovered</u> fully.

- Ensure that the effects of <u>uncommitted</u> transactions are <u>discarded</u> fully.

- We want a simple scheme to guarantee atomicity and durability.

4

# Assumptions

- Concurrency control is in effect
  - Strict 2-Phase Locking Specifically
- Updates are happening "in place"
  - Updates are applied directly to the (single) on-disk copy of the data.

5

# The Buffer Pool

- An intermediary between memory and disk.

- Challenges:

  - **Force**: Every write goes directly to disk?

    - Provides durability.

    - … but results in poor response time.

6

# The Buffer Pool

- An intermediary between memory and disk.

- Challenges:

  - **Steal**: Allow buffer-pool frames to be stolen from uncommitted xacts?

    - If stealing not allowed, poor throughput.

    - If stealing allowed, atomicity not preserved.

7

# The Buffer Pool

|  | No Steal | Steal |
|---|---|---|
| **Force** | Trivial! | |
| **No Force** | | Desired |

# Steal

- Enforcing **atomicity** is hard

  - To steal frame F, the current page P in F is written to disk.

  - Some transaction T has written to P.

  - What if T aborts?

  - Need to be able to UNDO T's writes to P.

9

# No Force

- Enforcing **durability** is hard

  - What if the system crashes before a modified page is written to disk?

  - What if the system crashes while writing several pages to disk?

  - Need to be able to REDO writes that haven't been finalized yet.

10

# Basic Idea: Logging

- Record all information needed to perform REDO and UNDO operations in a <u>log</u>.

  - Log accesses are all <u>sequential</u> writes.

  - Can use a separate disk for logs.

  - Minimal info (delta) is written to log, so can fit many updates in one page.

11

# Example Log Record

- Transaction ID

- Page ID

- Offset

- Length

- Old Data

- New Data

- (Some extra control info that we'll see shortly)

12

# Write-Ahead Logging

- The write-ahead logging protocol

  - Force the log record for an update <u>before</u> the corresponding data page is <u>written</u>.

    - Guarantees Atomicity (REDO).

  - Force all log records for a transaction <u>before</u> the transaction <u>commits</u>.

    - Guarantees Durability (UNDO).

13

# Write-Ahead Logging

- The write-ahead logging protocol

  - Force the log record for an update <u>before</u> the corresponding data page is <u>written</u>.

    - Guarantees Atomicity (REDO).

  - Force all log records for a transaction <u>before</u> the transaction <u>commits</u>.
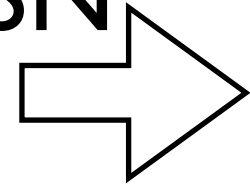
    - Guarantees Durability (UNDO).

  How would recovery be implemented?

13

# WAL & The Log

- Each log record has a unique Log Sequence Number (LSN)

- Each data page contains a Page LSN.

- System keeps (in memory) the max Flushed LSN (max LSN forced to disk)

- Before a page is written, ensure that **Page LSN ≤ Flushed LSN**

Flushed LSN

Tail (Ram Only)

14

# Log Records

- Record Types

  - **Update**

  - **Commit**

  - **Abort**

  - **End** (of commit/abort)

  - CLRs

## Log Record Fields

Prev LSN

XID

Record Type

Page ID

Length

Offset

Before-Image

After-Image

Only needed for **update**

15

# Other Log-Related State

- **Transaction Table**

  - One entry per active transaction

  - Contains <u>XID</u>, Transaction <u>Status</u>, <u>LastLSN</u>

- **Dirty Page Table**

  - One entry per dirty page in the buffer pool.

  - Contains **RecLSN**: LSN of the log record that first caused the page to be dirty.

16

# Transaction Execution

- Normal execution:

  - A Series of <u>reads</u> and <u>writes</u> followed by a <u>commit</u> or an <u>abort</u>.

  - Strict 2-Phase Locking

  - <u>Steal</u>, <u>No-Force</u> buffer management with <u>write-ahead logging</u>.

# The Big Picture

## The Log

**Log Records**
```
Prev LSN
XID
Record Type
Page ID
Length
Offset
Before-Image
After-Image
```

## The DB (Disk)

**Data Pages**
(each with a PageLSN)

**Master Record**

## The DB (Ram)

**Transaction Table**
```
XID
LastLSN
Status
```

**Dirty Page Table**
```
RecLSN
```

**Flushed LSN**

18

# Simple Transaction Abort
# (no crash involved)

"Play Back" the log in reverse order, UNDOing records

Transaction Table

Log

image credit: openclipart.org

# Simple Transaction Abort (no crash involved)

"Play Back" the log in reverse order, UNDOing records



Transaction Table

| | |
|---|---|
| | ABORT [XID] |

Log

(necessary for crash recovery)

Friday, March 22, 13

# Simple Transaction Abort
# (no crash involved)

## "Play Back" the log in reverse order, UNDOing records

Transaction Table

XID, LastLSN

| | ABORT [XID] |
|---|---|

Log

(necessary for crash recovery)

19

# Simple Transaction Abort (no crash involved)

"Play Back" the log in reverse order, UNDOing records

Transaction Table

| XID, LastLSN |
|---|

| LSN, Prev LSN, Prev Image, … | | ABORT [XID] |
|---|---|---|

Log

(necessary for crash recovery)

19

image credit: openclipart.org

# Simple Transaction Abort
# (no crash involved)

"Play Back" the log in reverse order, UNDOing records

Transaction Table

| XID, LastLSN |
| --- |

| LSN, Prev LSN, Prev Image, … | | ABORT [XID] |
| --- | --- | --- |

Log

(necessary for crash recovery)

19

image credit: openclipart.org

# Simple Transaction Abort
# (no crash involved)

"Play Back" the log in reverse order, UNDOing records

Transaction Table

XID, LastLSN

| | LSN, Prev LSN, Prev Image, … | | LSN, Prev LSN, Prev Image, … | | ABORT [XID] |
|---|---|---|---|---|---|

Log

(necessary for crash recovery)

19

image credit: openclipart.org

# Simple Transaction Abort (supporting crash recovery)

- Before restoring the old value of a page, write a Compensation Log Record (CLR).
  - Logging continues <u>during</u> UNDO processing.
  - CLR has an extra field: UndoNextLSN
    - Points to the next LSN to undo (the PrevLSN of the record currently being undone)
  - CLRs are never UNDOne.
    - But might be REDOne when repeating history.
    - (Why?)

20

# Transaction Commit

- Write **Commit** Record to Log

- All Log records up to the transaction's LastLSN are flushed.

  - Guarantees that FlushedLSN $\geq$ LastLSN

  - Note that Log Flushes are Sequential, Synchronous Writes to Disk

- Commit() returns.

- Write **End** record to log.

21

# Next week: Crash Recovery