# ARIES (& Logging)

*April 2-4, 2018*

What does it mean for a transaction to be committed?

If commit <u>returns</u> <u>successfully</u>, the transaction…

- … is recorded completely (atomicity)

- … left the database in a stable state (consistency)

- …'s effects are independent of other xacts (isolation)

- … will survive failures (durability)

commit
returns
successfully

=

the xact's
effects
are visible
<u>forever</u>

commit
returns
successfully

=

the xact's
effects
are visible
<u>forever</u>

commit
called but
doesn't
return

=

the xact's
effects
<u>may</u> be
visible

# Motivation

T1

T2

T3

T4

T5

Time

# Motivation

T1

T2

T3

T4

T5

Time

# Motivation

T1

T2

T3

T4

T5

Time

# Motivation

T1

T2

T3

T4

T5

Time

# Motivation

T1
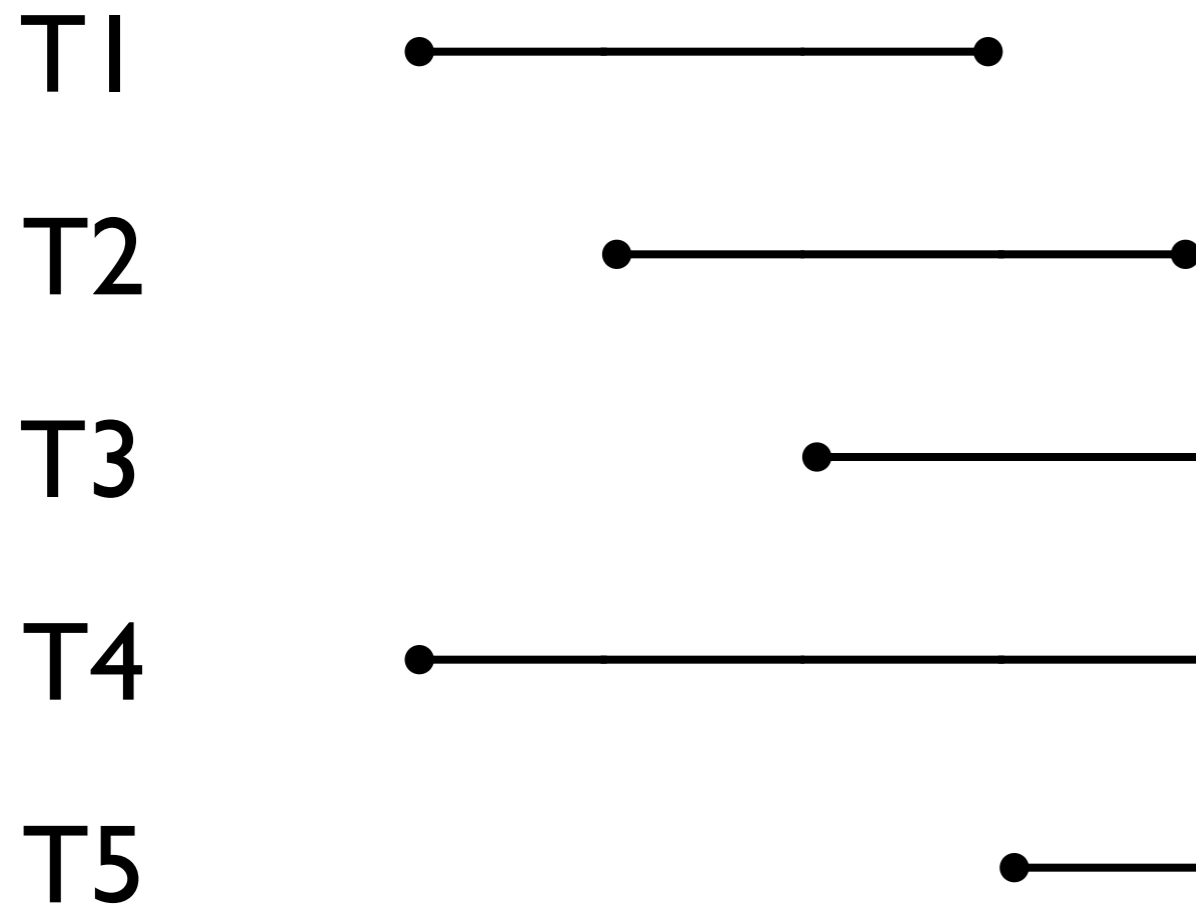
T2

T3

T4

T5

Time

# Motivation



T1

T2

T3

T4

T5

Time

# Motivation

T1

T2

T3

T4

T5

Time

# Motivation

T1

T2

T3

T4

T5

CRASH!

Time

Image copyright: Wikimedia Commons

# Motivation

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



T1

T2

T3

T4

T5

CRASH!

Time

Image copyright: Wikimedia Commons

# Motivation



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

T1

T2

T3

T4

T5

Time

Uncommitted Transactions.
These should leave no trace

CRASH!

6

Image copyright: Wikimedia Commons

- How do we guarantee durability under failures?

- How do aborted transactions get rolled back?

- How do we guarantee atomicity under failures?

**Problem 1**: Providing durability under failures.

# Simplified Model

When a write succeeds, the data is completely written

# Problems

- A crash occurs part-way through the write.

- A crash occurs before buffered data is written.

# Write-Ahead Logging

Before writing to the database, first write what you plan to write to a log file…

**Log**

`W(A:10)`

| A | 8 |
|---|---|
| B | 12 |
| C | 5 |
| D | 18 |
| E | 16 |

# Write-Ahead Logging

Once the log is safely on disk you can write the database

**Log**

`W(A:10)`

| | |
|---|---|
| **A** | ~~8~~ 10 |
| **B** | 12 |
| **C** | 5 |
| **D** | 18 |
| **E** | 16 |

# Write-Ahead Logging

Log is append-only, so writes are always efficient

**Log**

```
W(A:10)
W(C:8)
W(E:9)
```

| | |
|---|---|
| **A** | ~~8~~ 10 |
| **B** | 12 |
| **C** | 5 |
| **D** | 18 |
| **E** | 16 |

# Write-Ahead Logging

…allowing random writes
to be safely batched

**Log**

```
W(A:10)
W(C:8)
W(E:9)
```

| A | 8 10 |
|---|---|
| B | 12 |
| C | 5 8 |
| D | 18 |
| E | 16 9 |

**Problem 2**: Providing rollback.

# Single DB Model

**Txn 1**

A = 20
B = 14
COMMIT

**Txn 2**

E = 19
B = 15
ABORT

| A | 8 |
|---|---|
| B | 12 |
| C | 5 |
| D | 18 |
| E | 16 |

# Single DB Model

**Txn 1**

A = 20
B = 14
COMMIT

**Txn 2**

E = 19
B = 15
ABORT

| A | 8  20 |
|---|-------|
| B | 12 |
| C | 5 |
| D | 18 |
| E | 16 |

Image copyright: OpenClipart (rg1024)

# Single DB Model

**Txn 1**

A = 20
B = 14
COMMIT

**Txn 2**

E = 19
B = 15
ABORT

| A | 8   20 |
|---|--------|
| B | 17 2 |
| C | 5 |
| D | 18 |
| E | 16 19 |

# Single DB Model

**Txn 1**

```
A = 20
B = 14
COMMIT
```

**Txn 2**

```
E = 19
B = 15
ABORT
```

| A | 8 20 |
|---|------|
| B | 12 14 |
| C | 5 |
| D | 18 |
| E | 16 19 |

Image copyright: OpenClipart (rg1024)

# Single DB Model



**Txn 1**

A = 20
B = 14
COMMIT

**Txn 2**

E = 19
B = 15
ABORT

| | | | |
|---|---|---|---|
| A | ~~8~~ | 20 | |
| B | ~~12~~ | ~~14~~ | 15 |
| C | 5 | | |
| D | 18 | | |
| E | ~~16~~ | 19 | |

# Staged DB Model

**Txn 1**

➤ A = 20
B = 14
COMMIT

**Txn 2**

➤ E = 19
B = 15
ABORT

| | |
|---|---|
| A | 8 |
| B | 12 |
| C | 5 |
| D | 18 |
| E | 16 |

| | |
|---|---|
| A | 8 |
| B | 12 |
| C | 5 |
| D | 18 |
| E | 16 |

# Staged DB Model

**Txn 1**

```
A = 20
B = 14
```
➡ `COMMIT`

**Txn 2**

```
E = 19
B = 15
```
➡ `ABORT`

| | | |
|---|---|---|
| **A** | ~~8~~ | 20 |
| **B** | ~~12~~ | 14 |
| **C** | 5 | |
| **D** | 18 | |
| **E** | 16 | |

| | | |
|---|---|---|
| **A** | 8 | |
| **B** | ~~12~~ | 15 |
| **C** | 5 | |
| **D** | 18 | |
| **E** | ~~16~~ | 19 |

# Staged DB Model

**Txn 1**

```
A = 20
B = 14
```
→ `COMMIT`

**Txn 2**

```
E = 19
B = 15
```
→ `ABORT`



| | | |
|---|---|---|
| **A** | ~~8~~ | 20 |
| **B** | ~~12~~ | 14 |
| **C** | 5 | |
| **D** | 18 | |
| **E** | 16 | |

Is staging always possible?

- Staging takes up more memory.

- Merging after-the-fact can be harder.

- Merging after-the-fact introduces more latency!

for the single database model

**Problem 2**: Providing rollback.
∧

# UNDO Logging

Store both the "old" and the "new" values of the record being replaced

**Log**

```
W(A:8→10)
W(C:5→8)
W(E:16→9)
```

| A | 8 10 |
| B | 12 |
| C | 5 8 |
| D | 18 |
| E | 16 9 |

Image copyright: OpenClipart (rg1024)

# UNDO Logging



| | | |
|---|---|---|
| **A** | ~~8~~ | 10 |
| **B** | 12 | |
| **C** | ~~5~~ | 8 |
| **D** | 18 | |
| **E** | ~~16~~ | 9 |

**Active Xacts**

Xact:1, Log: 45

Xact:2, Log: 32

**Log**

```
43:W(A:8→10)
44:W(C:5→8)
45:W(E:16→9)
```

# UNDO Logging

**Active Xacts**

Xact:1, Log: 45 **ABORT**

Xact:2, Log: 32

**Log**

➡️ 
```
43:W(A:8→10)
44:W(C:5→8)
45:W(E:16→9)
```

| | |
|---|---|
| **A** | 8̶ 10 |
| **B** | 12 |
| **C** | 5̶ 8 |
| **D** | 18 |
| **E** | 16̶ 9 |

# UNDO Logging

**Active Xacts**

Xact:1, Log: 45 **ABORT**

Xact:2, Log: 32

**Log**

➡ 43:W(A:8→10)
44:W(C:5→8)
45:W(E:16→9)

| | |
|---|---|
| **A** | 8̶ 10 |
| **B** | 12 |
| **C** | 5̶ 8 |
| **D** | 18 |
| **E** | 16 |

# UNDO Logging

**Active Xacts**

Xact:1, Log: 45 **ABORT**

Xact:2, Log: 32

**Log**

43: W(A:8→10)

➡ 44: W(C:5→8)

45: W(E:16→9)

| | |
|---|---|
| **A** | ~~8~~ 10 |
| **B** | 12 |
| **C** | 5 |
| **D** | 18 |
| **E** | 16 |

# UNDO Logging

**Active Xacts**

Xact:1, Log: 45 **ABORT**

Xact:2, Log: 32

**Log**

➡ `43: W(A:8→10)`
`44: W(C:5→8)`
`45: W(E:16→9)`

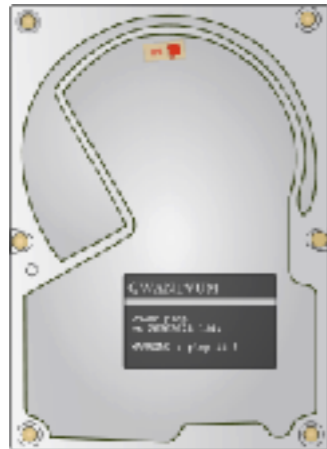| | |
|---|---|
| A | 8 |
| B | 12 |
| C | 5 |
| D | 18 |
| E | 16 |

# Log Sequence Number Linked Lists
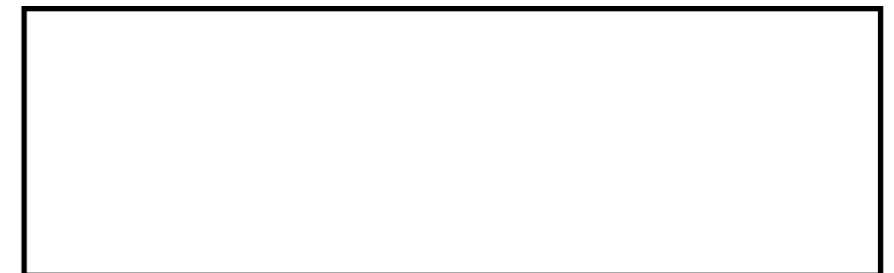
Transaction Table

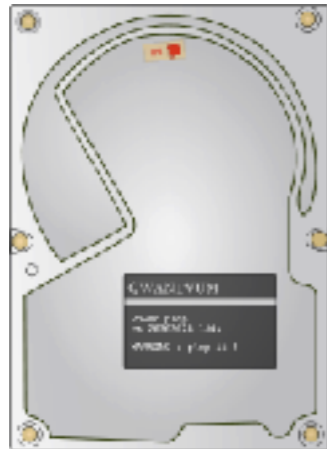Log

# Log Sequence Number Linked Lists

**Transaction Table**

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

**Log**

**(necessary for crash recovery)**

# Log Sequence Number Linked Lists

Transaction Table

| |
|---|
| XID, LastLSN |

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

Log

(necessary for crash recovery)

# Log Sequence Number
# Linked Lists

Transaction Table
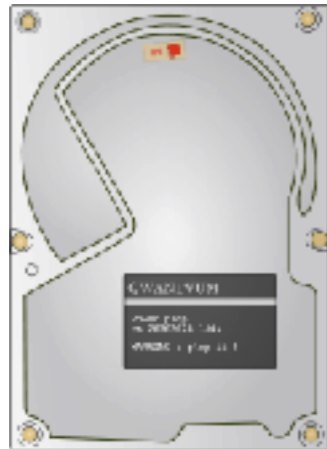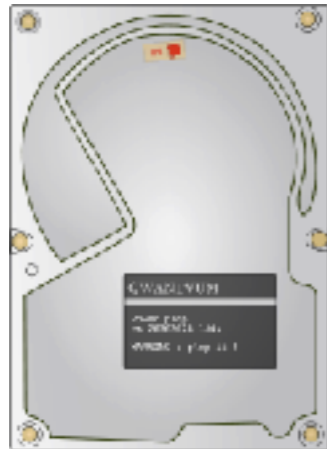
| XID, LastLSN |
| --- |

LSN, Prev LSN,
Prev Image, …

ABORT
[XID]

Log

(necessary for crash recovery)

# Log Sequence Number Linked Lists

Transaction Table

| XID, LastLSN |

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

Log

(necessary for crash recovery)

# Log Sequence Number
# Linked Lists

Transaction Table

| XID, LastLSN |
| --- |

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

Log

(necessary for crash recovery)

**Problem 3**: Providing atomicity.

**Goal**: Be able to reconstruct all state at the time of the DB's crash (minus all running xacts)

# What state is relevant?

# DB State



**Active Xacts**

Xact:1, Log: 45

Xact:2, Log: 32

**Log**

```
43:W(A:8→10)
44:W(C:5→8)
45:W(E:16→9)
```

| | |
|---|---|
| **A** | ~~8~~  10 |
| **B** | 12 |
| **C** | ~~5~~  8 |
| **D** | 18 |
| **E** | ~~16~~  9 |

# DB State

**On-Disk
(or rebuildable)**

**In-Memory
Only!**

**On-Disk**

**Active Xacts**

**Log**

Xact:1, Log: 45

Xact:2, Log: 32

```
43:W(A:8→10)
44:W(C:5→8)
45:W(E:16→9)
```

| | | |
|---|---|---|
| **A** | ~~8~~ | 10 |
| **B** | 12 | |
| **C** | ~~5~~ | 8 |
| **D** | 18 | |
| **E** | ~~16~~ | 9 |

# Rebuilding the Xact Table

Log every COMMIT
(replay triggers commit process)

Log every ABORT
(replay triggers abort process)

New message: END
(replay removes Xact from Xact Table)

# Rebuilding the Xact Table

Log every COMMIT
(replay triggers commit process)

Log every ABORT
(replay triggers abort process)

New message: END
(replay removes Xact from Xact Table)

What about BEGIN?
(when does an Xact get added to the Table?)

# Transaction Commit

- Write **Commit** Record to Log

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

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

- Commit() returns.

- Write **End** record to log.

# 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?)

# Rebuilding the Xact Table

**Optimization**: Write the Xact Table to the log periodically.
(checkpointing)

# ARIES Crash Recovery

- Start from checkpoint stored in master record.

- Analysis: Rebuild the Xact Table

- Redo: Replay operations from all live Xacts (even uncommitted ones).

- Undo: Revert operations from all uncommitted/aborted Xacts.

Oldest log record
of transaction
active at crash

Smallest recLSN
in dirty page table
after Analysis

Last Checkpoint

CRASH

A  R  U