# CSE 544
# Principles of Database Management Systems

Lectures 17-18 - Transactions: recovery

# Announcements

- Project presentations tomorrow! 10am – 2pm

- Paper review was due today

# References

- **Concurrency control and recovery.**

  Michael J. Franklin. The handbook of computer science and engineering. A. Tucker ed. 1997


- **Database management systems.**

  Ramakrishnan and Gehrke.

  Third Ed. **Chapters 16 and 18.**

# Outline

- **Review of ACID properties**
  - Today we will cover techniques for ensuring atomicity and durability in face of failures

- **Review of buffer manager and its policies**

- **Write-ahead log + simple UNDO / REDO recovery**

- **ARIES method for failure recovery**

# ACID Properties

- Atomicity: Either all changes performed by transaction occur or none occurs

- Consistency: A transaction as a whole does not violate integrity constraints

- Isolation: Transactions appear to execute one after the other in sequence

- Durability: If a transaction commits, its changes will survive failures

# What Could Go Wrong?

- Concurrent operations
  - That's what we discussed last time (atomicity and isolation properties)


- Failures can occur at any time
  - Today (isolation and durability properties)

# Problem Illustration

Client 1:

    START TRANSACTION
    INSERT INTO SmallProduct(name, price)
        SELECT pname, price
        FROM Product
        WHERE price <= 0.99

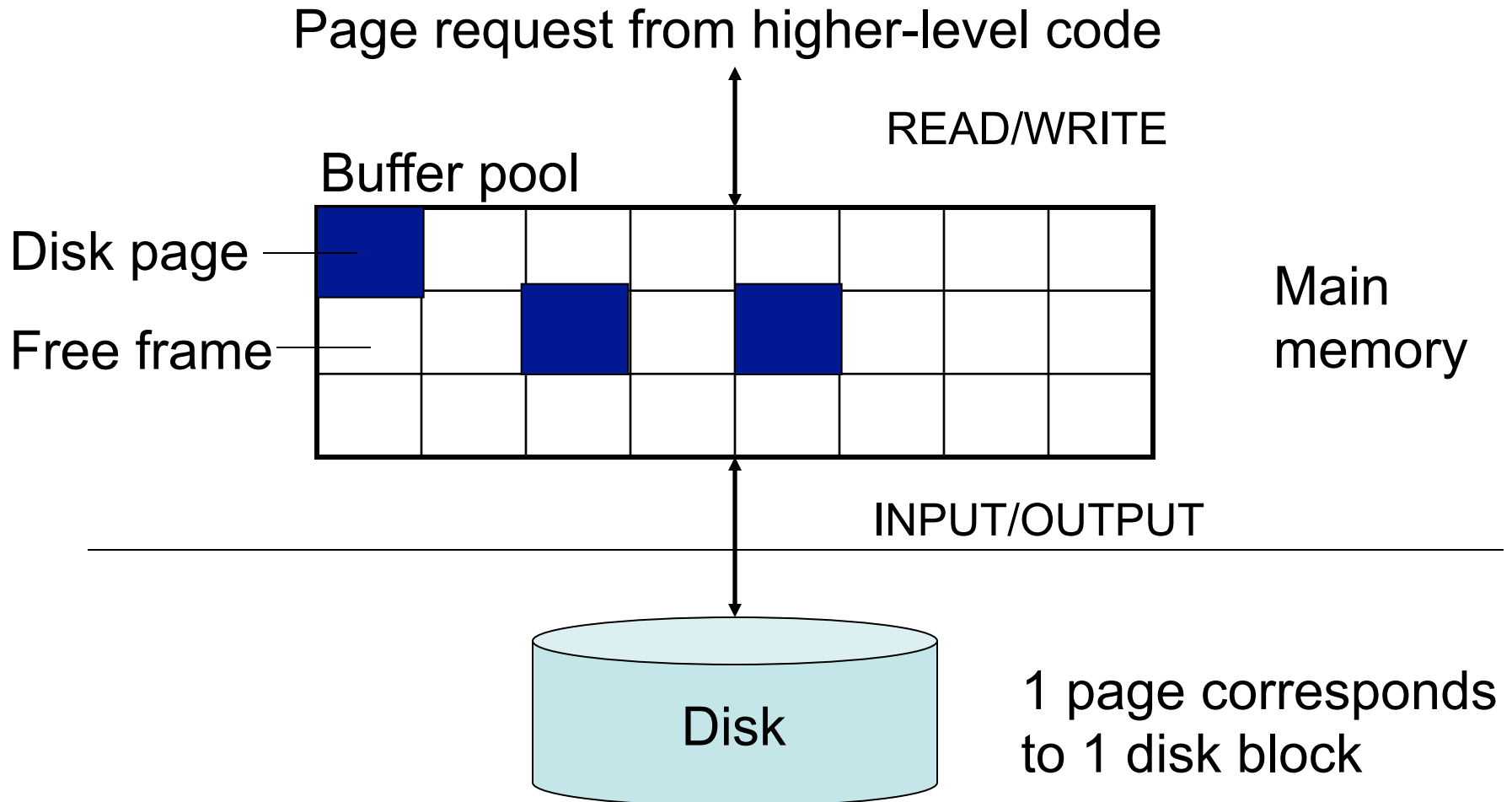Crash !

    DELETE Product
        WHERE price <=0.99
    COMMIT

# Handling Failures

- Types of failures
  - Transaction failure
  - System failure
  - Media failure -> we will not talk about this now

- Required capability: undo and redo

- Challenge: buffer manager
  - Changes performed in memory
  - Changes written to disk only from time to time

# Impact of Buffer Manager

Page request from higher-level code

READ/WRITE

Buffer pool

Disk page

Free frame

Main memory

INPUT/OUTPUT

Disk

1 page corresponds to 1 disk block

# Primitive Operations

- READ(X,t)
  - copy value of data item X to transaction local variable t

- WRITE(X,t)
  - copy transaction local variable t to data item X


- INPUT(X)
  - read page containing data item X to memory buffer

- OUTPUT(X)
  - write page containing data item X to disk

# Running Example

BEGIN TRANSACTION

READ(A,t);

t := t*2;

WRITE(A,t);

READ(B,t);

t := t*2;

WRITE(B,t)

COMMIT;

Initially, A=B=8.

**Atomicity** requires that either
(1) T commits and A=B=16, or
(2) T does not commit and A=B=8.

READ(A,t); t := t*2; WRITE(A,t);
READ(B,t); t := t*2; WRITE(B,t)

| | Transaction | Buffer pool | | Disk | |
| :---: | :---: | :---: | :---: | :---: | :---: |
| **Action** | **t** | **Mem A** | **Mem B** | **Disk A** | **Disk B** |
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |

Is this bad ?

| Action | t | Mem A | Mem B | Disk A | Disk B |
|--------|---|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |

Crash !

**Is this bad ?**   **Yes it's bad: A=16, B=8….**

| Action | t | Mem A | Mem B | Disk A | Disk B |
|---|---|---|---|---|---|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |

**Crash !**

## Is this bad ?

| Action | t | Mem A | Mem B | Disk A | Disk B |
|---|---|---|---|---|---|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |

Crash !

| Action | t | Mem A | Mem B | Disk A | Disk B |
|---|---|---|---|---|---|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |

Crash !

## Is this bad ?

| Action | t | Mem A | Mem B | Disk A | Disk B |
|--------|---|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |

Crash !

| Action | t | Mem A | Mem B | Disk A | Disk B |
|--------|---|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |

**Crash !**

## Typically, OUTPUT is after COMMIT (why?)

| Action | t | Mem A | Mem B | Disk A | Disk B |
|---|---|---|---|---|---|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| COMMIT | | | | | |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |

## Typically, OUTPUT is after COMMIT (why?)

| Action | t | Mem A | Mem B | Disk A | Disk B |
|---|---|---|---|---|---|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| COMMIT | | | | | |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |

Crash !

# Write-Ahead Log

- **Log: append-only file containing log records**
- For every update, commit, or abort operation
  - Write a log record
  - Multiple transactions run concurrently, log records are interleaved
- After a system crash, use log to:
  - Redo transactions that did commit
  - Undo other transactions that didn't commit

# Log Granularity

What is an "element" X:

- Physical log records: X = a block(page)
  - Position on a particular page where update occurred
  - Both before and after image for undo/redo logs
  - Benefits: Idempotent & updates are fast to redo/undo

- Logical log records: X = a record
  - Record only high-level information about the operation
  - Benefit: Smaller log
  - BUT difficult to implement because crashes can occur in the middle of an operation

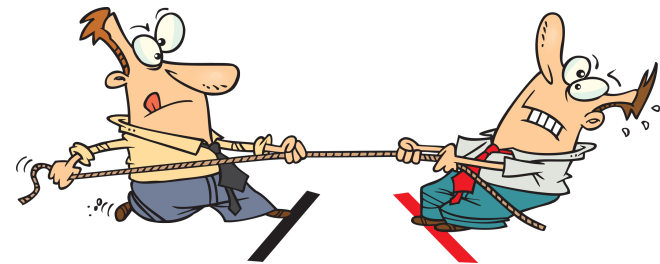# Buffer Manager Policies

- **STEAL or NO-STEAL**
  - Can an update made by an uncommitted transaction overwrite the most recent committed value of a data item on disk?

- **FORCE or NO-FORCE**
  - Should all updates of a transaction be forced to disk before the transaction commits?

- Easiest for recovery: **NO-STEAL**/**FORCE**

- Highest performance: **STEAL**/**NO-FORCE**

# Outline

- **Review of ACID properties**
  - Today we will cover techniques for ensuring atomicity and durability in face of failures

- **Review of buffer manager and its policies**

- **Write-ahead log + simple UNDO / REDO recovery**

- **ARIES method for failure recovery**

# UNDO Log

FORCE and STEAL

# Undo Logging

Log records

- <START T>
  - transaction T has begun
- <COMMIT T>
  - T has committed
- <ABORT T>
  - T has aborted
- <T,X,v>
  - T has updated element X, and its _old_ value was v

| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|--------|---|-------|-------|--------|--------|----------|
|  |  |  |  |  |  | \<START T\> |
| INPUT(A) |  | 8 |  | 8 | 8 |  |
| READ(A,t) | 8 | 8 |  | 8 | 8 |  |
| t:=t*2 | 16 | 8 |  | 8 | 8 |  |
| WRITE(A,t) | 16 | 16 |  | 8 | 8 | \<T,A,8\> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |  |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |  |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |  |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | \<T,B,8\> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |  |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |  |
| COMMIT |  |  |  |  |  | \<COMMIT T\> |

| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|--------|---|-------|-------|--------|--------|----------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | Crash ! |
| COMMIT | | | | | | <COMMIT T> |

WHAT DO WE DO ?

| Action | t | Mem A | Mem B | Disk A | Disk B | **UNDO** Log |
|---|---|---|---|---|---|---|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | Crash ! |
| COMMIT | | | | | | <COMMIT T> |

WHAT DO WE DO ?   We UNDO by setting B=8 and A=8

| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|--------|---|-------|-------|--------|--------|----------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |

What do we do now ?

Crash !

| Action | t | Mem A | Mem B | Disk A | Disk B | **UNDO** Log |
|---|---|---|---|---|---|---|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |

What do we do now ?

Nothing: log contains COMMIT

# Recovery with Undo Log

```
…
…
<T6,X6,v6>

…
…
<START T5>
<START T4>
<T1,X1,v1>
<T5,X5,v5>
<T4,X4,v4>
<COMMIT T5>
<T3,X3,v3>
<T2,X2,v2>
```

Crash !

Question1: Which updates are undone ?

Question 2:
How far back do we need to read in the log ?

Question 3:
What happens if there is a second crash, during recovery ?

| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|--------|---|-------|-------|--------|--------|----------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | | 8 | |
| READ(A,t) | 8 | | | | 8 | |
| t:=t*2 | 16 | 8 | | | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |

When must we force pages to disk ?

| Action | t | Mem A | Mem B | Disk A | Disk B | **UNDO** Log |
|--------|-----|-------|-------|--------|--------|--------------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | FORCE | <COMMIT T> |

RULES: log entry _before_ OUTPUT _before_ COMMIT

# Undo-Logging Rules

U1: If T modifies X, then <T,X,v> must be written to disk before OUTPUT(X)

U2: If T commits, then OUTPUT(X) must be written to disk before <COMMIT T>

FORCE

• Hence: OUTPUTs are done _early_, before the transaction commits

# REDO Log

NO-FORCE and NO-STEAL

# Redo Logging

One minor change to the undo log:

- <T,X,v>= T has updated element X, and its _new_ value is v

| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|--------|---|-------|-------|--------|--------|----------|
| | | | | | | <START T> |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT | | | | | | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |

| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|--------|---|-------|-------|--------|--------|----------|
|        |   |       |       |        |        | <START T> |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT | | | | | | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | Crash ! |

How do we recover ?

| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|---|---|---|---|---|---|---|
| | | | | | | <START T> |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT | | | | | | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | Crash ! |

How do we recover ?   We REDO by setting A=16 and B=16

# Recovery with Redo Log

<START T1>
<T1,X1,v1>
<START T2>
<T2, X2, v2>
<START T3>
<T1,X3,v3>
<COMMIT T2>
<T3,X4,v4>
<T1,X5,v5>

**Crash !**

Show actions during recovery

| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|--------|---|-------|-------|--------|--------|----------|
|        |   |       |       |        |        | <START T> |
| READ(A,t) | 8 | 8 |  |  | 8 |  |
| t:=t*2 | 16 | 8 |  |  | 8 |  |
| WRITE(A,t) | 16 | 16 |  | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |  |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |  |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT |  |  |  |  |  | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |  |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |  |

When must we force pages to disk ?

| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|--------|---|-------|-------|--------|--------|----------|
|  |  |  |  |  |  | <START T> |
| READ(A,t) | 8 | 8 |  | 8 | 8 |  |
| t:=t*2 | 16 | 8 |  | 8 | 8 |  |
| WRITE(A,t) | 16 | 16 |  | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |  |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |  |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT | NO-STEAL |  |  |  |  | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |  |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |  |

RULE: OUTPUT _after_ COMMIT

43

# Redo-Logging Rules

R1: If T modifies X, then both <T,X,$v$> and <COMMIT T> must be written to disk before OUTPUT(X)

<span style="color:blue">NO-STEAL</span>

- Hence: OUTPUTs are done *late*

# Comparison Undo/Redo

- Undo logging: OUTPUT must be done early:
  - Inefficient

- Redo logging: OUTPUT must be done late:
  - Inflexible

# Outline

- **Review of ACID properties**
  - Today we will cover techniques for ensuring atomicity and durability in face of failures

- **Review of buffer manager and its policies**

- **Write-ahead log + simple UNDO / REDO recovery**

- **ARIES method for failure recovery**

# ARIES

# Aries

- ARIES pieces together several techniques into a comprehensive algorithm

- Developed at IBM Almaden, by C. Mohan

- IBM botched the patent, so everyone uses it now

- Several variations, e.g. for distributed transactions

# Granularity in ARIES

- *Physiological logging*
  - Log records refer to a single page
  - But record logical operation within the page

- Page-oriented logging for REDO
  - Necessary since can crash in middle of complex operation

- Logical logging for UNDO
  - Enables **tuple-level locking**!
  - Why physical logging for REDO and logical logging for UNDO? (answer at the end of the lecture)

# ARIES Method

**Recovery from a system crash is done in 3 passes:**

1. **Analysis pass**
   – Figure out what was going on at time of crash
   – List of dirty pages and active transactions

2. **Redo pass (repeating history principle)**
   – Redo all operations, even for transactions that will not commit
   – Get back to state at the moment of the crash

3. **Undo pass**
   – Remove effects of all uncommitted transactions
   – Log changes during undo in case of another crash during undo

# ARIES Recovery Manager

- A redo/undo log
- Physiological logging
  - Physical logging for REDO
  - Logical logging for UNDO
- Efficient checkpointing

Why do we
do checkpointing?

# ARIES Recovery Manager

Log entries:

- <START T>   -- when T begins

- Update: <T,X,u,v>

  - T updates X, _old_ value=u, _new_ value=v

  - In practice: _undo only_ and _redo only_ entries

- <COMMIT T> or <ABORT T>

- CLR's – we'll talk about them later.

# ARIES Recovery Manager

Rule:

- If T modifies X, then <T,X,u,v> must be written to disk before OUTPUT(X)

We are free to OUTPUT early or late

# LSN = Log Sequence Number

- **<u>LSN</u>** = identifier of a log entry
  - Log entries belonging to the same TXN are linked

- Each page contains a **pageLSN**:
  - LSN of log record for latest update to that page

# ARIES Data Structures

- **Active Transactions Table**
  - Lists all active TXN's
  - For each TXN: lastLSN = its most recent update LSN

- **Dirty Page Table**
  - Lists all dirty pages
  - For each dirty page: recoveryLSN (recLSN)= first LSN that caused page to become dirty

- **Write Ahead Log**
  - LSN, prevLSN = previous LSN for same txn

$W_{T100}(P7)$
$W_{T200}(P5)$
$W_{T200}(P6)$
$W_{T100}(P5)$

# ARIES Data Structures

## Dirty pages

| pageID | recLSN |
|--------|--------|
| P5 | 102 |
| P6 | 103 |
| P7 | 101 |

## Log (WAL)

| LSN | prevLSN | transID | pageID | Log entry |
|-----|---------|---------|--------|-----------|
| 101 | - | T100 | P7 | |
| 102 | - | T200 | P5 | |
| 103 | 102 | T200 | P6 | |
| 104 | 101 | T100 | P5 | |

## Active transactions

| transID | lastLSN |
|---------|---------|
| T100 | 104 |
| T200 | 103 |

## Buffer Pool

| P8 | P2 | . . . |
|----|----|-------|
| | . . . | |
| **P5** PageLSN=104 | **P6** PageLSN=103 | **P7** PageLSN=101 |
| | | |

# ARIES Normal Operation

T writes page P

- What do we do ?

# ARIES Normal Operation

T writes page P

- What do we do ?

<br>

- Write **<T,P,u,v>** in the **Log**
- **prevLSN**=**lastLSN**
- **pageLSN**=**LSN**
- **lastLSN**=**LSN**
- **recLSN**=if isNull then **LSN**

# ARIES Normal Operation

Buffer manager wants to OUTPUT(P)

- What do we do ?

Buffer manager wants INPUT(P)

- What do we do ?

# ARIES Normal Operation

Buffer manager wants to OUTPUT(P)

- Flush log up to **pageLSN**

- Remove P from **Dirty Pages** table

Buffer manager wants INPUT(P)

- Create entry in **Dirty Pages** table
  **recLSN** = NULL

# ARIES Normal Operation

Transaction T starts

- What do we do ?

Transaction T commits/aborts

- What do we do ?

# ARIES Normal Operation

Transaction T starts

- Write **<START T>** in the log

- New entry T in Active TXN;
  lastLSN = null

Transaction T commits/aborts

- Write **<COMMIT T>** in the log

- Flush log up to this entry

# Checkpoints

Write into the log

- Entire active transactions table
- Entire dirty pages table

Recovery always starts by analyzing latest checkpoint

Background process periodically flushes dirty pages to disk

# ARIES Recovery

1. **Analysis pass**
   - Figure out what was going on at time of crash
   - List of dirty pages and active transactions

2. **Redo pass (repeating history principle)**
   - Redo all operations, even for transactions that will not commit
   - Get back to state at the moment of the crash

3. **Undo pass**
   - Remove effects of all uncommitted transactions
   - Log changes during undo in case of another crash during undo

# ARIES Method Illustration



Figure 3: The Three Passes of ARIES Restart

First undo and first redo log entry might be in reverse order

[Figure 3 from Franklin97]

# 1. Analysis Phase

- Goal
  - Determine point in log (**firstLSN**) where to start REDO
  - Determine set of dirty pages when crashed
    - Conservative estimate of dirty pages
  - Identify active transactions when crashed

- Approach
  - Rebuild active transactions table and dirty pages table
  - Reprocess the log from the checkpoint
    - Only update the two data structures
  - Compute: **firstLSN** = smallest of all recoveryLSN

# 1. Analysis Phase

**Log**     Checkpoint                    (crash)

**firstLSN**= ???

*Where do we start the REDO phase ?*

**Dirty pages**

| pageID | recLSN | pageID |
|--------|--------|--------|
|        |        |        |
|        |        |        |

**Active txn**

| transID | lastLSN | transID |
|---------|---------|---------|
|         |         |         |
|         |         |         |

# 1. Analysis Phase

**Log**      Checkpoint      (crash)

**firstLSN**=min(**recLSN**)

**Dirty pages**

| pageID | recLSN | pageID |
|--------|--------|--------|
|        |        |        |
|        |        |        |

**Active txn**

| transID | lastLSN | transID |
|---------|---------|---------|
|         |         |         |
|         |         |         |

# 1. Analysis Phase

**Log**    Checkpoint        (crash)

**firstLSN**

**Dirty pages**

| pageID | recLSN | pageID |
|---|---|---|
|  |  |  |
|  |  |  |

Replay history

| pageID | recLSN | pageID |
|---|---|---|
|  |  |  |
|  |  |  |

**Active txn**

| transID | lastLSN | transID |
|---|---|---|
|  |  |  |
|  |  |  |

| transID | lastLSN | transID |
|---|---|---|
|  |  |  |
|  |  |  |

# 2. Redo Phase

Main principle: replay history

- Process Log forward, starting from **firstLSN**

- Read every log record, sequentially

- Redo actions are not recorded in the log

- Needs the Dirty Page Table

# 2. Redo Phase: Details

For each Log entry record LSN: **&lt;T,P,u,v&gt;**

- Re-do the action P=u and WRITE(P)

- But which actions can we skip, for efficiency ?

# 2. Redo Phase: Details

For each Log entry record LSN: **<T,P,u,v>**

- If P is not in Dirty Page then **no update**

- If recLSN > LSN, then **no update**

- Read page from disk:
  If **pageLSN** > LSN, then **no update**

- Otherwise perform update

# 2. Redo Phase: Details

What happens if system crashes during REDO ?

# 2. Redo Phase: Details

What happens if system crashes during REDO ?

We REDO again !  Each REDO operation is _idempotent_:
  doing it twice is the as as doing it once.

# 3. Undo Phase

- Cannot "unplay" history, in the same way as we "replay" history

- WHY NOT ? Time to answer this question☺

# 3. Undo Phase

- Cannot "unplay" history, in the same way as we "replay" history

- WHY NOT ? Time to answer this question☺

- Need to support ROLLBACK!
  Selective undo, for one transaction only
  - Cannot simply undo physical actions
  - E.g. Txn updates a record on a page, another Txn updates another record on the same page: don't undo the latter
  - E.g. Txn updates a B$^+$-tree, causing rebalancing, other Txn do other update: don't undo the latter!

- Hence, *logical* undo v.s. *physical* redo

# 3. Undo Phase

Main principle: "logical" undo

- Start from end of Log, move backwards

- Read only affected log entries

- Undo actions *are* written in the Log as special entries: CLR (Compensating Log Records)

- CLRs are redone, but never undone

# 3. Undo Phase: Details

- "Loser transactions" =
  - Uncommitted transactions in Active Transactions Table
  - Or transactions to be rolled back

- **ToUndo** = set of lastLSN of loser transactions

# 3. Undo Phase: Details

While **ToUndo** not empty:

- Choose most recent (largest) LSN in **ToUndo**
- If LSN = regular record **<T,P,u,v>**:
  - Undo v
  - Write a CLR where CLR.undoNextLSN = LSN.prevLSN
- If LSN = CLR record:
  - Don't undo !
- if CLR.**undoNextLSN** not null, insert in **ToUndo** otherwise, write **<END TRANSACTION>** in log

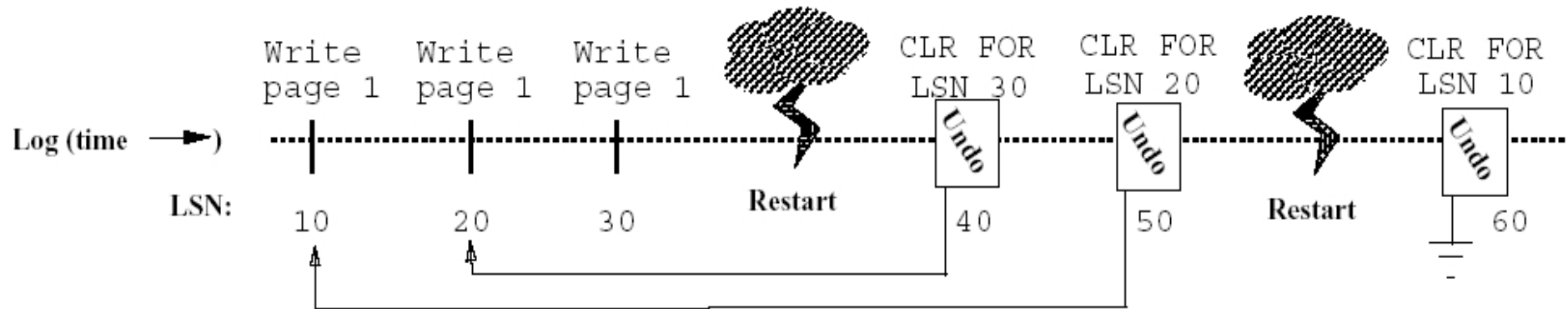# 3. Undo Phase: Details



Figure 4: The Use of CLRs for UNDO

[Figure 4 from Franklin97]

# 3. Undo Phase: Details

What happens if system crashes during UNDO ?

# 3. Undo Phase: Details

What happens if system crashes during UNDO ?

We do not UNDO again !  Instead, each CLR is a REDO record: we simply redo the undo

# Physical v.s. Logical Loging

Why are redo records *physical* ?


Why are undo records *logical* ?

# Physical v.s. Logical Loging

Why are redo records *physical* ?

*   Simplicity: replaying history is easy, and idempotent

Why are undo records *logical* ?

*   Required for transaction rollback: this not "undoing history", but selective undo