# Introduction to Database Systems
# CSE 444

Lectures 9-10

Transactions: recovery

# Outline

- We are starting to look at DBMS internals

- Next pair of lectures: transactions & recovery
  - Disks 13.2
  - Undo logging 17.2
  - Redo logging 17.3
  - Redo/undo 17.4

# The Mechanics of Disk

Mechanical characteristics:

- Rotation speed (5400 RPM)
- Number of platters (1-30)
- Number of tracks (<=10000)
- Number of bytes/track($10^5$)

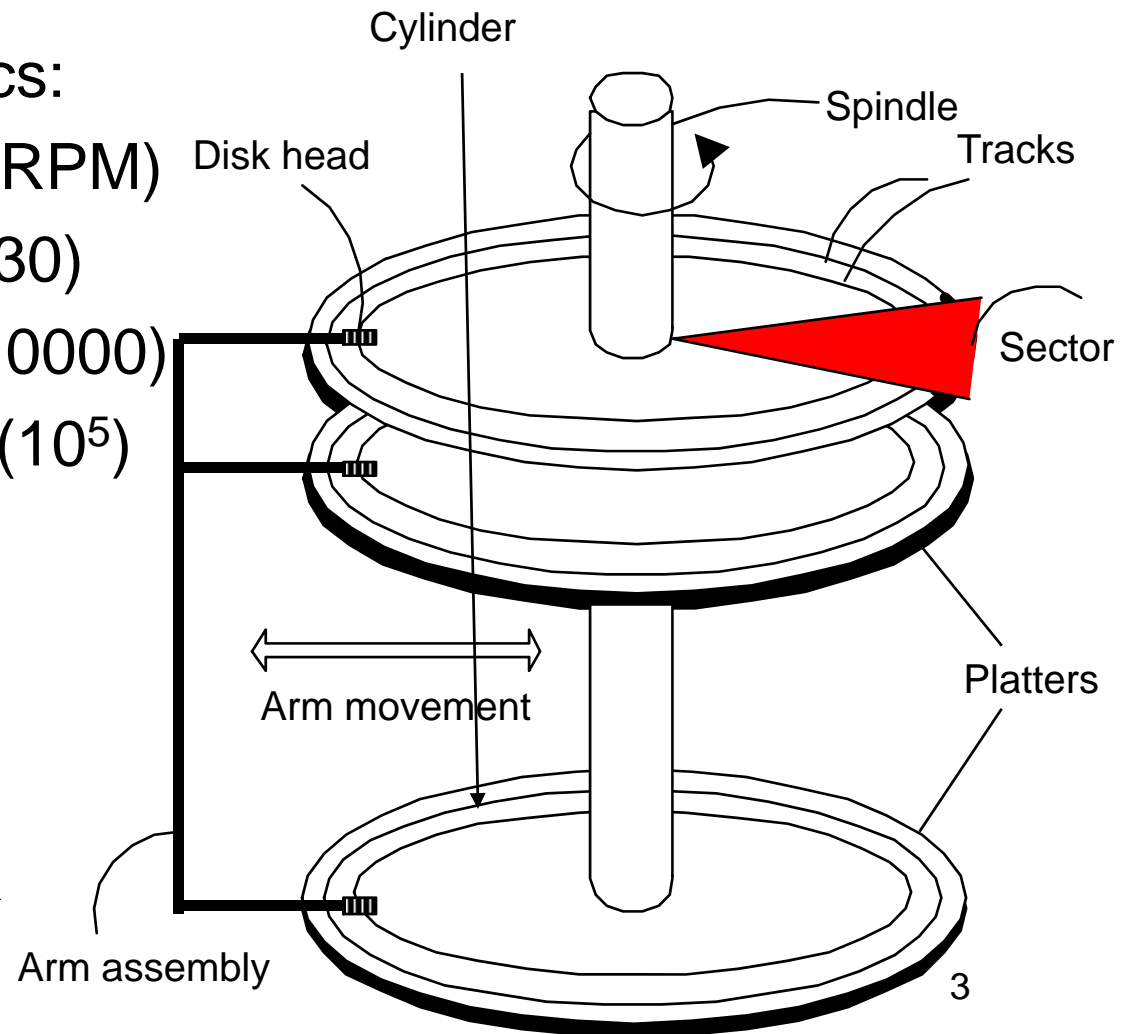Unit of read or write:
**disk block**

Once in memory:
**page**

Typically: 4k or 8k or 16k

Cylinder

Spindle

Tracks

Disk head

Sector

Arm movement

Platters

Arm assembly

3

# RAID

Several disks that work in parallel
• Redundancy: use parity to recover from disk failure
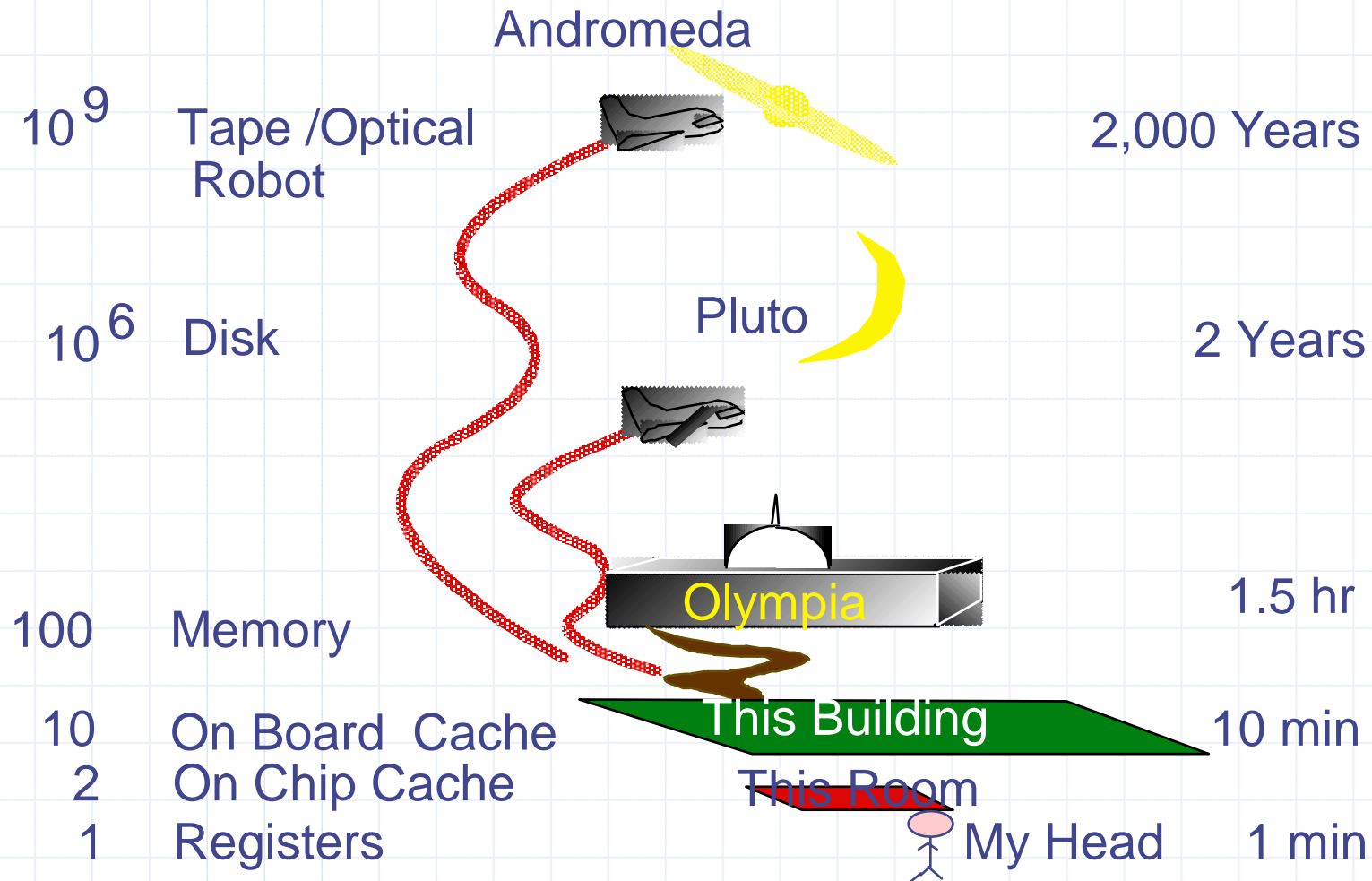• Speed: read from several disks at once

Various configurations (called *levels*):
• RAID 1 = mirror
• RAID 4 = n disks + 1 parity disk
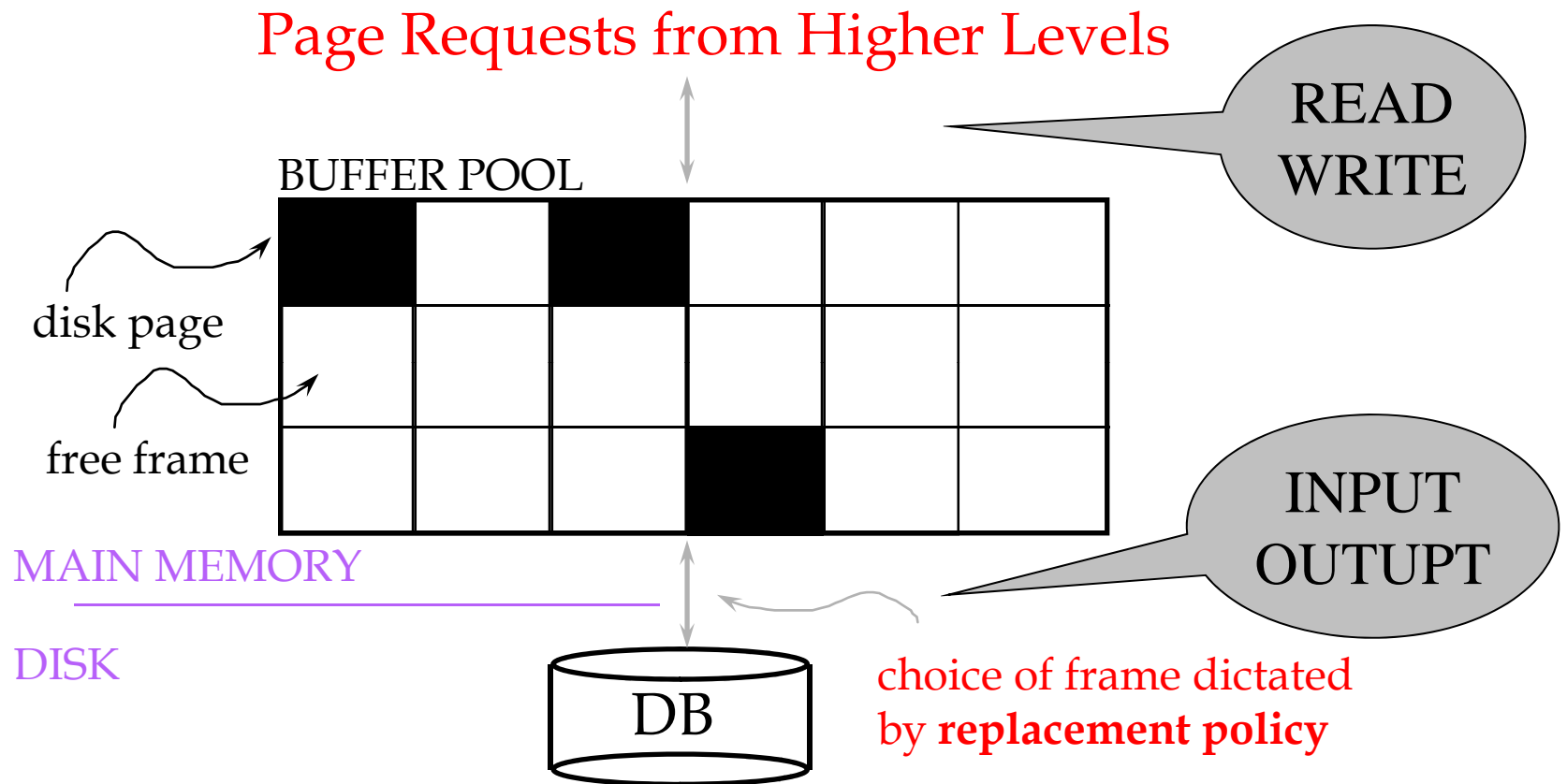• RAID 5 = n+1 disks, assign parity blocks round robin
• RAID 6 = "Hamming codes"

# Disk Access Characteristics

- Disk latency = time between when command is issued and when data is in memory

- Disk latency = seek time + rotational latency
  - Seek time = time for the head to reach cylinder
    - 10ms – 40ms
  - Rotational latency = time for the sector to rotate
    - Rotation time = 10ms
    - Average latency = 10ms/2
- Transfer time = typically 40MB/s
- Disks read/write one block at a time

# Storage Latency:
# How Far Away is the Data?

Andromeda

$10^9$    Tape /Optical Robot                                 2,000 Years

Pluto

$10^6$   Disk                                         2 Years

Olympia            1.5 hr

100    Memory

10    On Board  Cache    This Building    10 min

2    On Chip Cache    This Room

1    Registers                      My Head    1 min

# Buffer Management in a DBMS

Page Requests from Higher Levels

BUFFER POOL

READ
WRITE

disk page

free frame

MAIN MEMORY

DISK

DB

INPUT
OUTUPT

choice of frame dictated
by **replacement policy**

- Data must be in RAM for DBMS to operate on it!

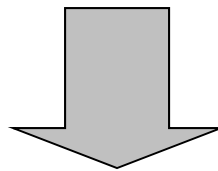- Table of <frame#, pageid> pairs is maintained

# Buffer Manager

- Enables higher layers of the DBMS to assume that needed data is in main memory

- Needs to decide on page replacement policy
  - LRU, clock algorithm, or other

- Both work well in OS, but not always in DB

# Least Recently Used (LRU)

- Order pages by the time of last accessed
- Always replace the least recently accessed

P5, P2, P8, P4, P1, P9, P6, P3, P7

Access P6

P6, P5, P2, P8, P4, P1, P9, P3, P7

LRU is expensive (why ?); the clock algorithm is good approx

# Buffer Manager

- Why not use the OS for the task??

- Reason 1: Correctness
  - DBMS needs fine grained control for transactions
  - Needs to force pages to disk for recovery purposes

- Reason 2: Performance
  - DBMS may be able to anticipate access patterns
  - Hence, may also be able to perform prefetching
  - May select better page replacement policy
  - May want to pin pages in the buffer

# Transaction Management and the Buffer Manager

Transaction manager operates on buffer pool

- **Recovery**: 'log-file write-ahead', then careful policy about which pages to force to disk

- **Concurrency control**: locks at the page level, multiversion concurrency control

# Transaction Management

Two parts:

- Recovery from crashes:  <u>A</u>CI<u>D</u>
- Concurrency control:      AC<u>I</u>D

Both operate on the buffer pool

# 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

## What do we do now?

# Recovery

| Type of Crash | Prevention |
|---|---|
| Wrong data entry | Constraints and Data cleaning |
| Disk crashes | Redundancy: e.g. RAID, archive |
| Fire, theft, bankruptcy… | Buy insurance, Change jobs… |
| System failures | DATABASE RECOVERY |

# Main Idea for Recovery

- Each transaction has *internal state*
- When system crashes, internal state is lost
  - Don't know which parts executed and which didn't
  - Need ability to *undo* and *redo*
- Remedy: use a **log**

  - File that records every single action of all running transactions
  - After a crash, transaction manager reads the log to find out exactly what each transaction did or did not do

# Transactions

- Assumption: db composed of **_elements_**
  - Usually 1 element = 1 block
  - Can be smaller (=1 record) or larger (=1 relation)

- Assumption: each transaction reads/writes some elements

# Primitive Operations of Transactions

- ## READ(X,t)
  - copy element X to transaction local variable t
- ## WRITE(X,t)
  - copy transaction local variable t to element X

- ## INPUT(X)
  - read element X to memory buffer
- ## OUTPUT(X)
  - write element X to disk

# Example

START TRANSACTION

READ(A,t);

t := t*2;

WRITE(A,t);

READ(B,t);

t := t*2;

WRITE(B,t);

COMMIT;

Atomicity:
BOTH A and B
are multiplied by 2

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 |
| READ(A,t) | | | | | |
| t:=t*2 | | | | | |
| WRITE(A,t) | | | | | |
| INPUT(B) | | | | | |
| READ(B,t) | | | | | |
| t:=t*2 | | | | | |
| WRITE(B,t) | | | | | |
| OUTPUT(A) | | | | | |
| OUTPUT(B) | | | | | |

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) | | | | | |
| t:=t*2 | | | | | |
| WRITE(A,t) | | | | | |
| INPUT(B) | | | | | |
| READ(B,t) | | | | | |
| t:=t*2 | | | | | |
| WRITE(B,t) | | | | | |
| OUTPUT(A) | | | | | |
| OUTPUT(B) | | | | | |

```
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) |  |  |  |  |  |
| INPUT(B) |  |  |  |  |  |
| READ(B,t) |  |  |  |  |  |
| t:=t*2 |  |  |  |  |  |
| WRITE(B,t) |  |  |  |  |  |
| OUTPUT(A) |  |  |  |  |  |
| OUTPUT(B) |  |  |  |  |  |

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

```
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) |  |  |  |  |  |
| t:=t*2 |  |  |  |  |  |
| WRITE(B,t) |  |  |  |  |  |
| OUTPUT(A) |  |  |  |  |  |
| OUTPUT(B) |  |  |  |  |  |

```
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) | | | | | |
| OUTPUT(A) | | | | | |
| OUTPUT(B) | | | | | |

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) |  |  |  |  |  |
| OUTPUT(B) |  |  |  |  |  |

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) | | | | | |

```
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 |

| 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 | Crash ! |
| OUTPUT(B) | 16 | 16 | 16 | 16 | |

Crash occurs after OUTPUT(A), before OUTPUT(B)
We lose atomicity

28

# 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

# The Log

- Log = append-only file containing log records
- Multiple transactions run concurrently, log records are interleaved
- After a system crash, use log to:
  - Redo some transactions that did commit
  - Undo other transactions that did not commit

- Three kinds of logs: undo, redo, undo/redo
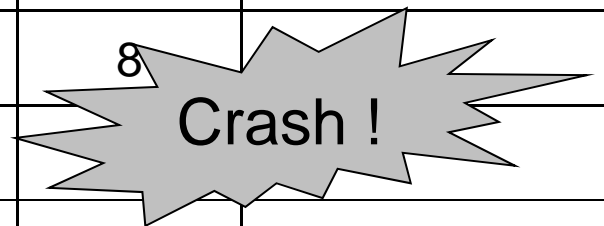
# Undo Logging

Log records

- **<START T>**
  - Transaction T has begun

- **<COMMIT T>**
  - T has committed

- **<ABORT T>**
  - T has aborted

- **<T,X,v>** -- Update record
  - T has updated element X, and its *old* value was v

| Action | T | Mem A | Mem B | Disk A | Disk B | 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 | 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 | Crash ! | |
| COMMIT | | | | | | <COMMIT T> |

WHAT DO WE DO ?

| Action | T | Mem A | Mem B | Disk A | Disk B | 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 ?

Crash !

# After Crash

- In the first example:
  - We UNDO both changes: A=8, B=8
  - The transaction is atomic, since none of its actions has been executed

- In the second example
  - We don't undo anything
  - The transaction is atomic, since both it's actions have been executed

# 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>

- Hence: OUTPUTs are done *early*, before the transaction commits

| Action | T | Mem A | Mem B | Disk A | Disk B | 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> |

# Recovery with Undo Log

After system's crash, run recovery manager

- Idea 1. Decide for each transaction T whether it is completed or not
  - <START T>….<COMMIT T>….    = yes
  - <START T>….<ABORT T>…….  = yes
  - <START T>……………………  = no

- Idea 2. Undo all modifications by incomplete transactions

# Recovery with Undo Log

Recovery manager:

- Read log <u>from the end</u>; cases:

  &lt;COMMIT T&gt;:  mark T as completed

  &lt;ABORT T&gt;: mark T as completed

  &lt;T,X,v&gt;: if T is not completed
                      then write X=v to disk
            else ignore

  &lt;START T&gt;: ignore

# 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

Question 1: 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?

40

# Recovery with Undo Log

- Note: all undo commands are *idempotent*
  - If we perform them a second time, no harm done
  - E.g. if there is a system crash during recovery, simply restart recovery from scratch

# Recovery with Undo Log

When do we stop reading the log ?

- We cannot stop until we reach the beginning of the log file
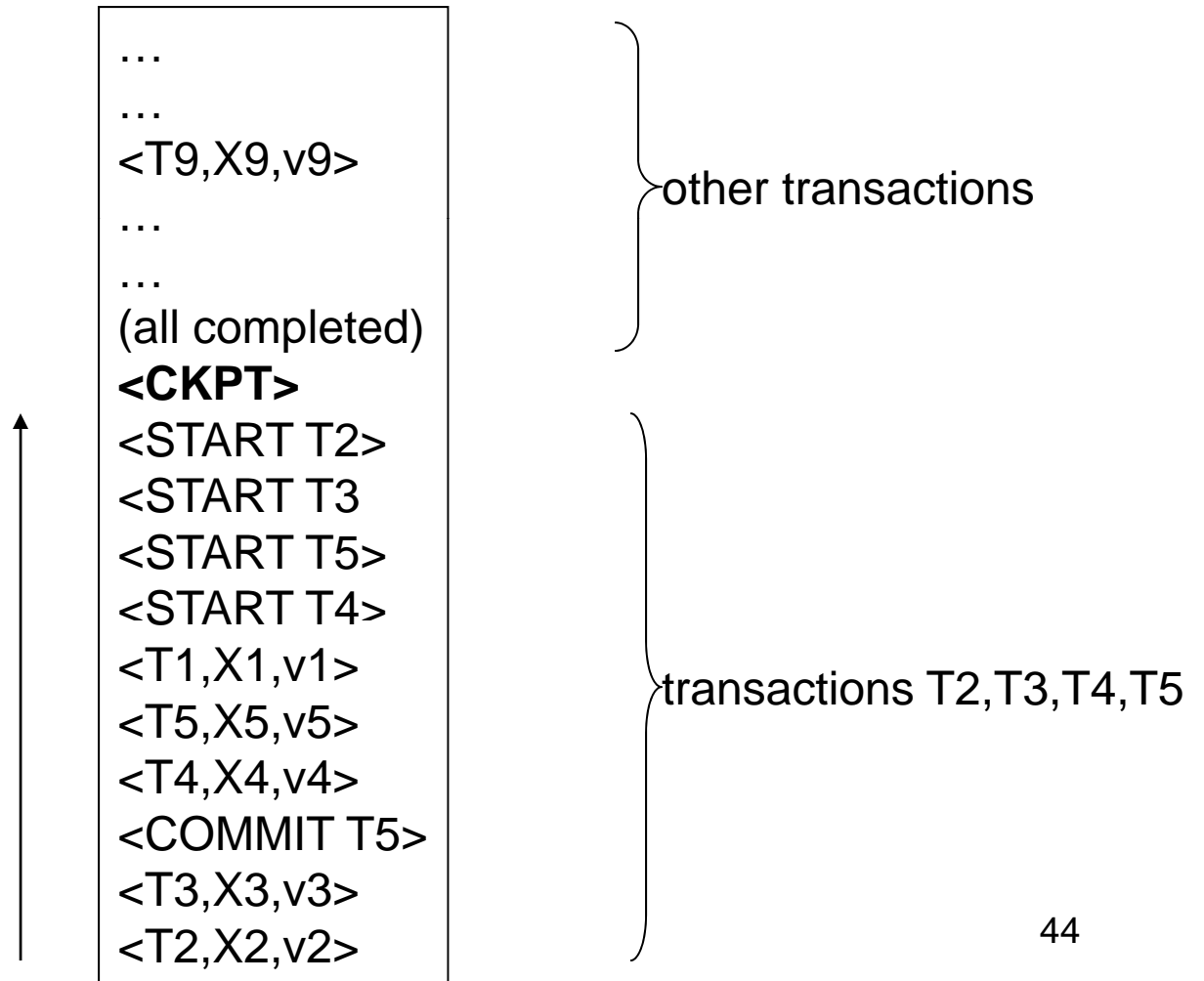
- This is impractical

Instead: use checkpointing

# Checkpointing

Checkpoint the database periodically

- Stop accepting new transactions
- Wait until all current transactions complete
- Flush log to disk
- Write a <CKPT> log record, flush
- Resume transactions

# Undo Recovery with Checkpointing

During recovery,
Can stop at first
<CKPT>

```
…
…
<T9,X9,v9>
…
…
(all completed)
**<CKPT>**
<START T2>
<START T3
<START T5>
<START T4>
<T1,X1,v1>
<T5,X5,v5>
<T4,X4,v4>
<COMMIT T5>
<T3,X3,v3>
<T2,X2,v2>
```

other transactions

transactions T2,T3,T4,T5

# Nonquiescent Checkpointing

- Problem with checkpointing: database freezes during checkpoint
- Would like to checkpoint while database is operational
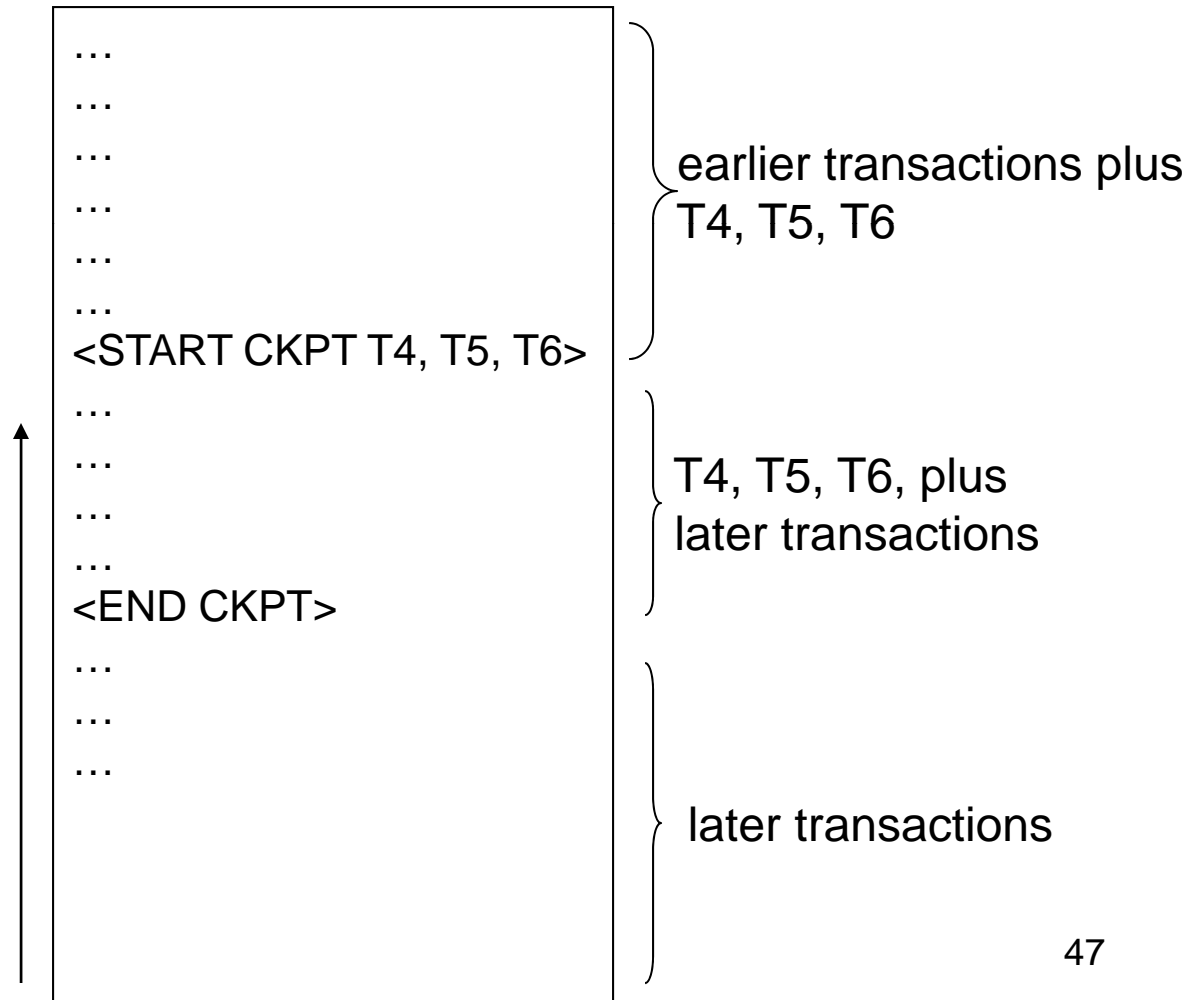- Idea: nonquiescent checkpointing

Quiescent = being quiet, still, or at rest; inactive
Non-quiescent = allowing transactions to be active

# Nonquiescent Checkpointing

- Write a <START CKPT(T1,…,Tk)>
  where T1,…,Tk are all active transactions.

- Continue normal operation

- When all of T1,…,Tk have completed, write
  <END CKPT>.

# Undo Recovery with Nonquiescent Checkpointing

During recovery,
Can stop at first
<CKPT>

```
...
...
...
...
...
...
<START CKPT T4, T5, T6>
...
...
...
...
<END CKPT>
...
...
...
```

earlier transactions plus
T4, T5, T6

T4, T5, T6, plus
later transactions

later transactions

Q: Do we need
<END CKPT> ?

# Implementing ROLLBACK

- Recall: a transaction can end in COMMIT or ROLLBACK

- Idea: use the undo-log to implement ROLLBCACK

- How ?
  - LSN = Log Sequence Number
  - Log entries for the same transaction are linked, using the LSN's
  - Read log in reverse, using LSN pointers

# Undo Logging Critique

- Works!
- But….
  - Requires physical OUTPUT before transaction can commit
    - Can cause unnecessary I/O ops if more updates will be done on the same buffer page soon
    - What if two transactions share the same buffer page and only one is ready to commit?  (this one is subtle – more later…)

# Redo 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 _new_ value is v

| Action | T | Mem A | Mem B | Disk A | Disk B | 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 T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |

# Redo-Logging Rules

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

- Hence: OUTPUTs are done *late*

| Action | T | Mem A | Mem B | Disk A | Disk B | 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 T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |

# Recovery with Redo Log

After system's crash, run recovery manager

- Step 1. Decide for each transaction T whether it is completed or not
  - \<START T\>....\<COMMIT T\>....   = yes
  - \<START T\>....\<ABORT T\>.......  = yes
  - \<START T\>.........................  = no

- Step 2. Read log from the beginning, redo all updates of _committed_ transactions

# 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>
…

…

# Nonquiescent Checkpointing

- Write a <START CKPT(T1,…,Tk)> where T1,…,Tk are all active transactions

- Flush to disk all blocks of committed transactions (*dirty blocks*), while continuing normal operation
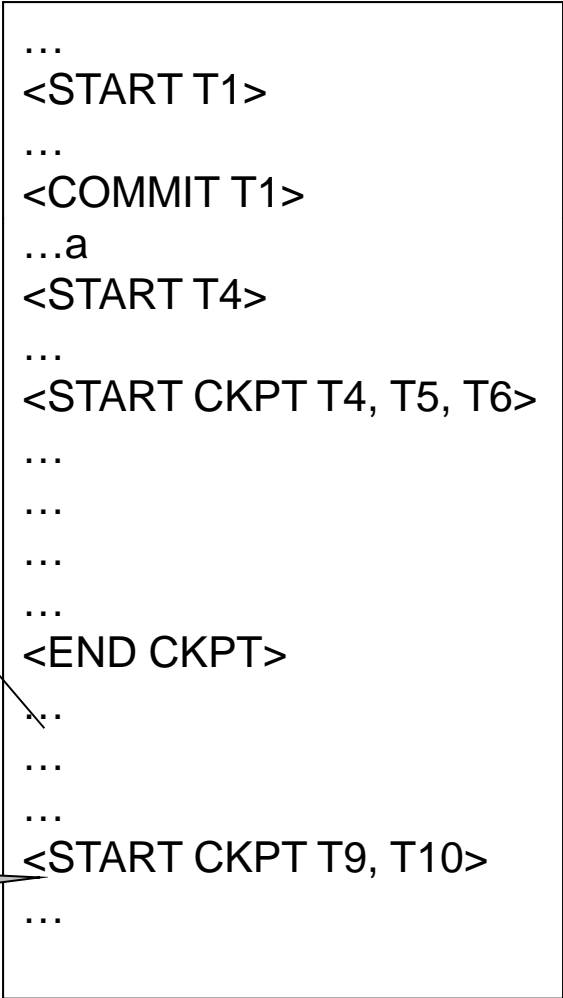
- When all blocks have been written, write <END CKPT>

# Redo Recovery with Nonquiescent Checkpointing

**Step 1: look for
The last
<END CKPT>**

All OUTPUTs
of T1 are
known to be on disk

Cannot
use

```
…
<START T1>
…
<COMMIT T1>
…a
<START T4>
…
<START CKPT T4, T5, T6>
…
…
…
…
<END CKPT>
…
…
…
<START CKPT T9, T10>
…
```

**Step 2: redo
from the
earliest
start of
T4, T5, T6
ignoring
transactions
committed
earlier**

# Comparison Undo/Redo

- Undo logging:
  - OUTPUT must be done early
  - If <COMMIT T> is seen, T definitely has written all its data to disk (hence, don't need to redo) – inefficient

  Steal/Force

- Redo logging
  - OUTPUT must be done late
  - If <COMMIT T> is not seen, T definitely has not written any of its data to disk (hence there is not dirty data on disk, no need to undo) – inflexible

  No-Steal/No-Force

- Would like more flexibility on when to OUTPUT: undo/redo logging (next)

  Steal/No-Force

# Undo/Redo Logging

Log records, only one change

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

# Undo/Redo-Logging Rule

UR1: If T modifies X, then <T,X,u,v> must be written to disk before OUTPUT(X)

Note: we are free to OUTPUT early or late relative to <COMMIT T>

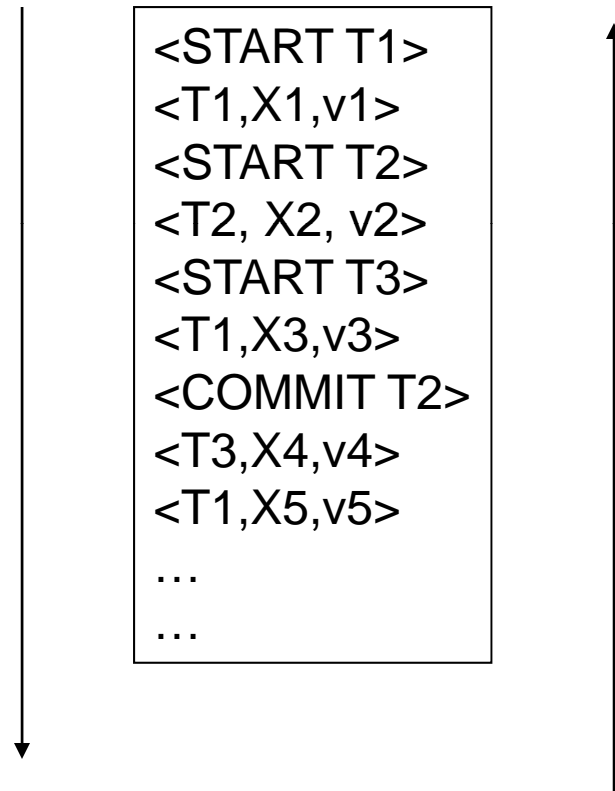| Action | T | Mem A | Mem B | Disk A | Disk B | Log |
|---|---|---|---|---|---|---|
| | | | | | | <START T> |
| REAT(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t)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

Can OUTPUT whenever we want: before/after COMMIT

# Recovery with Undo/Redo Log

After system's crash, run recovery manager

- Redo all committed transaction, top-down
- Undo all uncommitted transactions, bottom-up

# Recovery with Undo/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>
…
…

# Granularity of the Log

- Physical logging: element = physical page
- Logical logging: element = data record


- What are the pros and cons ?

# Granularity of the Log

- Modern DBMS:


- Physical logging for the REDO part
  - Efficiency
- Logical logging for the UNDO part
  - For ROLLBACKs