# CSE 562
# Database Systems

## Failure Recovery

Some slides are based or modified from originals by
*Database Systems: The Complete Book,*
*Pearson Prentice Hall 2nd Edition*
*©2008 Garcia-Molina, Ullman, and Widom*

*cse@buffalo*

---

## Integrity or Correctness of Data

- Would like data to be "accurate" or "correct" at all times

EMP

| Name | Age |
|------|-----|
| White | 52 |
| Green | 3421 |
| Gray | 1 |

---

## Integrity or Consistency Constraints

- Predicates data must satisfy
- Examples:
  - x is key of relation R
  - x → y holds in R
  - Domain(x) = {Red, Blue, Green}
  - $\alpha$ is valid index for attribute x of R
  - no employee should make more than twice the average salary

---

## Definition

- <u>Consistent state:</u> satisfies all constraints
- <u>Consistent DB:</u> DB in consistent state

1

## Constraints (as we use here) may not capture "full correctness"

**Example 1:** Transaction constraints
- When salary is updated,
    new salary > old salary
- When account record is deleted,
    balance = 0

---

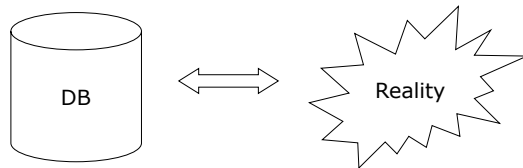Note: could be "emulated" by simple
constraints, e.g.,

account

| Acct # | …. | balance | deleted? |
|--------|-----|---------|----------|

---

## Constraints (as we use here) may not capture "full correctness"

**Example 2:** Database should reflect real world

DB $\longleftrightarrow$ Reality

---

In any case, continue with constraints…

**Observation:** DB cannot be consistent always!
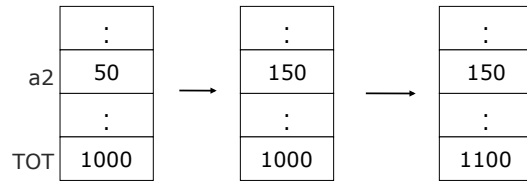
**Example:** a1 + a2 +…. an = TOT (constraint)

Deposit $100 in a2:
$$\begin{cases} a2 \leftarrow a2 + 100 \\ TOT \leftarrow TOT + 100 \end{cases}$$

**Example:** a1 + a2 +…. an = TOT (constraint)

Deposit $100 in a2:
$$a2 \leftarrow a2 + 100$$
$$TOT \leftarrow TOT + 100$$

| | a2 | | | | |
|---|---|---|---|---|---|
| : | 50 | → | 150 | → | 150 |
| : | : | | : | | : |
| TOT | 1000 | | 1000 | | 1100 |

**Transaction:** collection of actions
that preserve consistency

Consistent DB — T — Consistent DB'

## Big Assumption

If T starts with consistent state +
  T executes in isolation
⇒ T leaves consistent state

## Correctness (informally)

- If we stop running transactions,
  DB left consistent
- Each transaction sees a consistent DB

## How Can Constraints Be Violated?

- Transaction bug
- DBMS bug
- Hardware failure
    - e.g., disk crash alters balance of account
- Data sharing, e.g.:

  T1: give 10% raise to programmers
  T2: change programmers $\Rightarrow$ systems analysts

## How Can We Prevent/Fix Violations?

- Chapter 17: Due to failures only
- Chapter 18: Due to data sharing only
- Chapter 19: Due to failures and sharing

## Will Not Consider

- How to write correct transactions
- How to write correct DBMS
- Constraint checking & repair
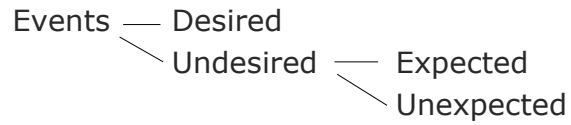    - That is, solutions studied here do not need to know constraints

## Chapter 17:  Recovery

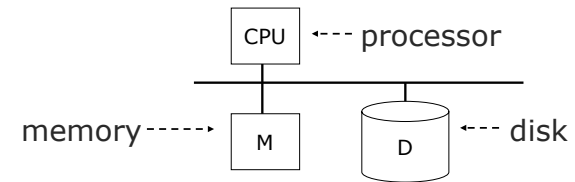- First order of business:
    - Failure Model

Events — Desired

Undesired — Expected

Unexpected

## Our Failure Model

CPU ◄--- processor

memory ----► M      D ◄--- disk

Desired events: see product manuals….

Undesired expected events:

System crash

- memory lost
- cpu halts, resets

═══════ that's it!! ═══════

Undesired Unexpected: Everything else!

Undesired Unexpected: Everything else!

Examples:

• Disk data is lost

• Memory lost without CPU halt

• CPU implodes wiping out universe…

## Is This Model Reasonable?

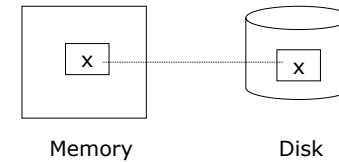Approach: Add low level checks + redundancy to increase probability model holds

E.g., Replicate disk storage (stable store)
Memory parity
CPU checks

## Second Order of Business:

Storage hierarchy



Memory                    Disk

## Operations

- Input (x): block containing x → memory
- Output (x): block containing x → disk

- Read (x,t): do input(x) if necessary
  t ← value of x in block
- Write (x,t): do input(x) if necessary
  value of x in block ← t

## Key Problem: Unfinished Transaction

**Example**    Constraint: A=B
T1:  A ← A × 2
      B ← B × 2

6

T1:  Read (A,t);  t ← t×2
     Write (A,t);
     Read (B,t);  t ← t×2
     Write (B,t);
     Output (A);
     Output (B);                    failure!

A: 8̶ 16
B: 8̶ 16

memory

A: 8̶ 16
B: 8

disk

---

- Need atomicity:    execute all actions of
                     a transaction or none
                     at all

---

One solution: undo logging
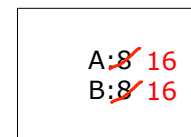
(immediate modification)


due to: Hansel and Gretel, 782 AD

---

## Undo Logging (Immediate modification)
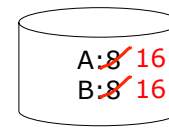
T1:  Read (A,t);  t ← t×2        A=B
     Write (A,t);
     Read (B,t);  t ← t×2
     Write (B,t);
     Output (A);
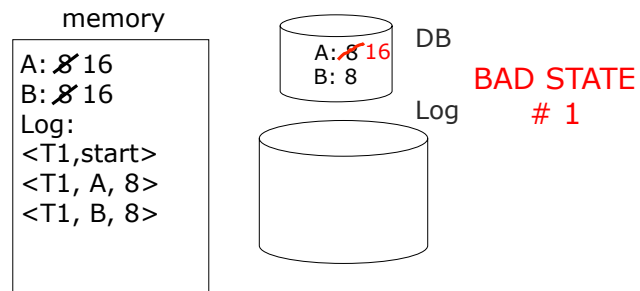     Output (B);

A: 8̶ 16
B: 8̶ 16

memory

A: 8̶ 16
B: 8̶ 16

disk

<T1, start>
<T1, A, 8>
<T1, B, 8>
<T1, commit>

log

7

## One "Complication"

- Log is first written in memory
- Not written to disk on every action

memory

A: ~~8~~ 16
B: ~~8~~ 16
Log:
<T1,start>
<T1, A, 8>
<T1, B, 8>

DB

A: ~~8~~ 16
B: 8

Log

**BAD STATE # 1**

## One "Complication"

- Log is first written in memory
- Not written to disk on every action

memory

A: ~~8~~ 16
B: ~~8~~ 16
Log:
<T1,start>
<T1, A, 8>
<T1, B, 8>
<T1, commit>

DB

A: ~~8~~ 16
B: 8

Log

⋮
<T1, B, 8>
<T1, commit>

**BAD STATE # 2**

## Undo Logging Rules

(1) For every action generate undo log
     record (containing old value)
(2) Before *x* is modified on disk, log
     records pertaining to *x* must be
     on disk (write ahead logging: WAL)
(3) Before commit is flushed to log, all
     writes of transaction must be
     reflected on disk

## Recovery Rules: Undo Logging

- For every Ti with <Ti, start> in log:
    - If <Ti,commit> or <Ti,abort> in log:
        Do nothing
    - Else $\left\{ \begin{array}{l} \text{For all <Ti, } X, v\text{> in log:} \\ \text{write } (X, v) \\ \text{output } (X) \\ \text{Write <Ti, abort> to log} \end{array} \right.$

IS THIS CORRECT??

8

## Recovery Rules: Undo Logging

(1) Let S = set of transactions with
    &lt;Ti, start&gt; in log, but no
    &lt;Ti, commit&gt; (or &lt;Ti, abort&gt;) record
                       in log

(2) For each &lt;Ti, X, v&gt; in log,
    in reverse order (latest → earliest) do:
    - if Ti ∈ S then { - write (X, v)
                         - output (X)

(3) For each Ti ∈ S do
    - write &lt;Ti, abort&gt; to log

## Question

- Can writes of &lt;Ti, abort&gt; records be done in any order (in Step 3)?
  - **Example:** T1 and T2 both write A
  - T1 executed before T2
  - T1 and T2 both rolled-back
  - &lt;T1, abort&gt; written but NOT &lt;T2, abort&gt;

         ↑               ↑      →
    T1 write A      T2 write A    time/log

---

What if failure during recovery?
    No problem!       Undo **idempotent!**

## To Discuss:

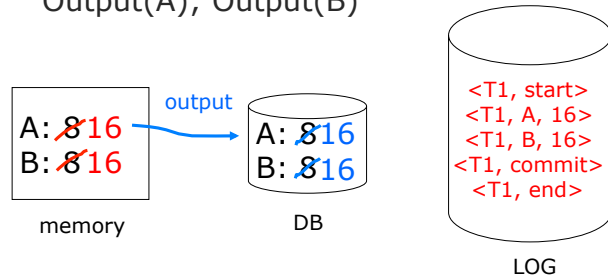- Redo logging
- Undo/redo logging, why both?
- Real world actions
- Checkpoints
- Media failures

## Redo Logging (deferred modification)

T1: Read(A,t); t←t×2; write (A,t);
    Read(B,t); t←t×2; write (B,t);
    Output(A); Output(B)

A: 8̶ 16    output    A: 8̶ 16
B: 8̶ 16              B: 8̶ 16

memory          DB

LOG:
<T1, start>
<T1, A, 16>
<T1, B, 16>
<T1, commit>
<T1, end>

## Redo Logging Rules

(1) For every action, generate redo log
     record (containing new value)
(2) Before X is modified on disk (DB),
     all log records for transaction that
     modified X (including commit) must
     be on disk
(3) Flush log at commit
(4) Write END record after DB updates
     flushed to disk

## Recovery Rules: Redo Logging

- For every Ti with <Ti, commit> in log:
  – For all <Ti, X, v> in log:
        Write(X, v)
        Output(X)

        IS THIS CORRECT??

## Recovery Rules: Redo Logging

(1) Let S = set of transactions with
     <Ti, commit> (and no <Ti, end>) in log
(2) For each <Ti, X, v> in log, in forward
     order (earliest → latest) do:
     - If Ti ∈ S then  Write(X, v)
                       Output(X)
(3) For each Ti ∈ S, write <Ti, end>

10

## Combining <Ti, end> Records

• Want to delay DB flushes for hot objects

Actions:

Say X is branch balance:
T1: ... update X...
T2: ... update X...
T3: ... update X...
T4: ... update X...

write X
~~output X~~
write X
~~output X~~
write X
~~output X~~
write X
output X
combined <end> **(checkpoint)**

## Solution: Checkpoint

• no <ti, end> actions
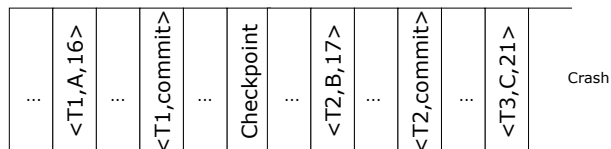• simple checkpoint

Periodically:

(1) Do not accept new transactions
(2) Wait until all transactions finish
(3) Flush all log records to disk (log)
(4) Flush all buffers to disk (DB) (do not discard buffers)
(5) Write "checkpoint" record on disk (log)
(6) Resume transaction processing

## Example: What To Do At Recovery?

Redo log (disk):

| ... | <T1,A,16> | ... | <T1,commit> | ... | Checkpoint | ... | <T2,B,17> | ... | <T2,commit> | ... | <T3,C,21> |
|-----|-----------|-----|-------------|-----|------------|-----|-----------|-----|-------------|-----|-----------|

Crash

## Key Drawbacks

• *Undo logging:* cannot bring backup DB
                copies up to date
• *Redo logging:* need to keep all modified
                blocks in memory
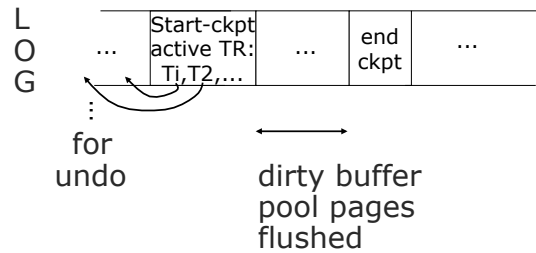                until commit

11

## Solution: Undo/Redo Logging!
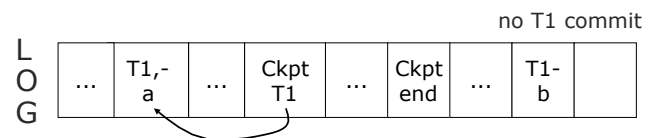
Update ⇒  <Ti, Xid, New X val, Old X val>
page X

## Rules

- Page X can be flushed before or after Ti commit
- Log record flushed before corresponding updated page (WAL)
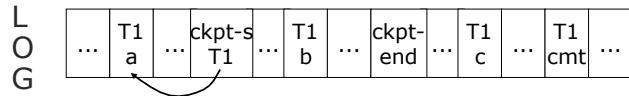- Flush at commit (log only)

## Non-Quiesce Checkpoint

```
L       ┌──────────┬─────┬──────┬─────┐
O  ...  │Start-ckpt│ ... │ end  │ ... │
G       │active TR:│     │ ckpt │     │
        │Ti,T2,... │     │      │     │
        └──────────┴─────┴──────┴─────┘
    ⋮
   for              ←──────→
  undo          dirty buffer
                pool pages
                flushed
```

## Example: What To Do At Recovery Time?

```
                                            no T1 commit
L     ┌─────┬──────┬─────┬──────┬─────┬──────┬─────┬──────┬─────┐
O ... │ ... │ T1,- │ ... │ Ckpt │ ... │ Ckpt │ ... │ T1-  │     │
G     │     │  a   │     │  T1  │     │ end  │     │  b   │     │
      └─────┴──────┴─────┴──────┴─────┴──────┴─────┴──────┴─────┘
```

Undo T1 (undo a,b)

12

## Example

L
O
G

| ... | T1 a | ... | ckpt-s T1 | ... | T1 b | ... | ckpt-end | ... | T1 c | ... | T1 cmt | ... |

Redo T1 (redo b,c)

## Recover From <u>Valid</u> Checkpoint

L
O
G

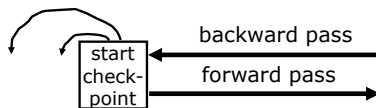| ... | ckpt start | ... | ckpt end | ... | T1 b | ... | ckpt-start | ... | T1 c | ... | |

start
of latest
valid
checkpoint

## Recovery Process

- **Backwards pass**
  (end of log ➜ latest valid checkpoint start)
  – construct set S of committed transactions
  – undo actions of transactions not in S
- **Undo pending transactions**
  – follow undo chains for transactions
    in (checkpoint active list) - S
- **Forward pass** (latest checkpoint start ➜ end of log)
  – redo actions of S transactions

| start check-point |
backward pass
forward pass

## Real World Actions

E.g., dispense cash at ATM
   Ti = a1 a2 … … aj … … an

$

13

## Solution

(1) execute real-world actions after commit
(2) try to make idempotent

---

ATM

Give$$
(amt, Tid, time)

lastTid: [      ]
time: [      ]
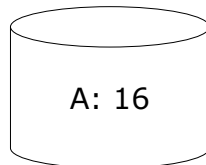
↓ give(amt)

$

---

## Media Failure (loss of non-volatile storage)

A: 16

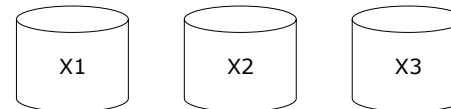Solution: Make copies of data!

---

## Example 1: Triple Modular Redundancy

- Keep 3 copies on separate disks
- Output(X) --> three outputs
- Input(X) --> three inputs + vote

X1          X2          X3

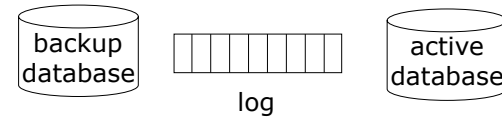## Example 2: Redundant Writes, Single Reads

- Keep N copies on separate disks
- Output(X) --> N outputs
- Input(X) --> Input one copy
  - if ok, done
  - else try another one

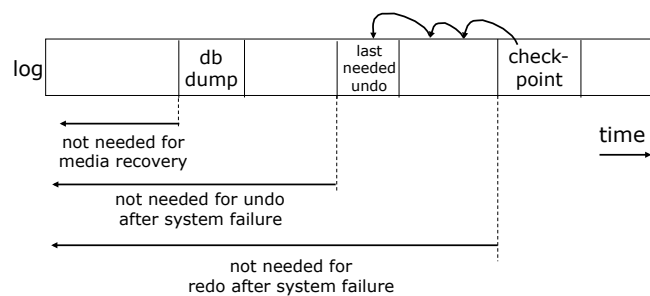Assumes bad data can be detected

## Example 3: DB Dump + Log



backup database          log          active database

- If active database is lost,
  – restore active database from backup
  – bring up-to-date using redo entries in log

## When Can Log Be Discarded?



log | | db dump | | last needed undo | | check-point | |

time

not needed for media recovery

not needed for undo after system failure

not needed for redo after system failure

## Summary

- Consistency of data
- One source of problems: failures
  - Logging
  - Redundancy
- Another source of problems:
  Data Sharing… next