

CSE232A: Database System Principles

Notes 08: Failure Recovery

1

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

2

Integrity or consistency constraints

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

3

Definition:

- Consistent state: satisfies all constraints
- Consistent DB: DB in consistent state

4

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

5

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

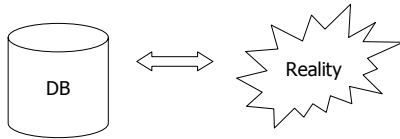
account

Acct #	balance	deleted?
--------	------	---------	----------

6

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

Example 2 Database should reflect real world



7

☞ in any case, continue with constraints...

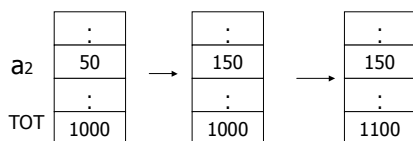
Observation: DB cannot be consistent always!

Example: $a_1 + a_2 + \dots + a_n = \text{TOT}$ (constraint)
Deposit \$100 in a_2 : $\left\{ \begin{array}{l} a_2 \leftarrow a_2 + 100 \\ \text{TOT} \leftarrow \text{TOT} + 100 \end{array} \right.$

8

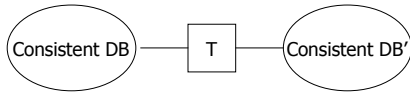
Example: $a_1 + a_2 + \dots + a_n = \text{TOT}$ (constraint)

Deposit \$100 in a_2 : $\left\{ \begin{array}{l} a_2 \leftarrow a_2 + 100 \\ \text{TOT} \leftarrow \text{TOT} + 100 \end{array} \right.$



9

Transaction: collection of actions that preserve consistency



10

Big assumption:

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

11

Correctness (informally)

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

12

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 ⇒ systems analysts

How can we prevent/fix violations?

- Chapter 8: due to failures only
- Chapter 9: due to data sharing only
- Chapter 10: 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 8 Recovery

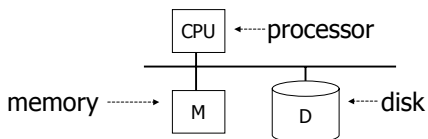
- First order of business:
Failure Model

16



17

Our failure model



18

Desired events: see product manuals...

Undesired expected events:

System crash

- memory lost
- cpu halts, resets

———— that's it!! ————

Undesired Unexpected: Everything else!

19

Undesired Unexpected: Everything else!

Examples:

- Disk data is lost
- Memory lost without CPU halt
- CPU implodes wiping out universe....

20

Is this model reasonable?

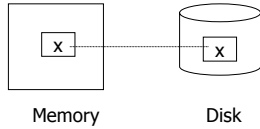
Approach: Add low level checks +
redundancy to increase
probability model holds

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

21

Second order of business:

Storage hierarchy



Operations:

- Input (x): block with x → memory
- Output (x): block with 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

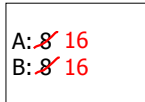
Constraint: $A=B$

$T_1: A \leftarrow A \times 2$

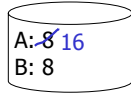
$B \leftarrow B \times 2$

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

failure!



memory



disk

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

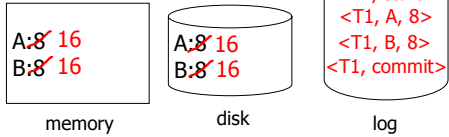
One solution: undo logging (immediate modification)

due to: Hansel and Gretel, 782 AD

- Improved in 784 AD to durable undo logging

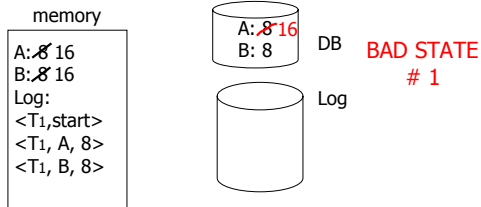
Undo logging (Immediate modification)

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);



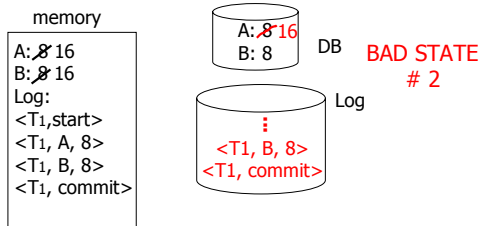
One "complication"

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



One "complication"

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



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

31

Recovery rules: Undo logging

- For every T_i with $\langle T_i, \text{start} \rangle$ in log:
 - If $\langle T_i, \text{commit} \rangle$ or $\langle T_i, \text{abort} \rangle$ in log, do nothing
 - Else { For all $\langle T_i, X, v \rangle$ in log:
 - { write (X, v)
 - { output (X)
 - Write $\langle T_i, \text{abort} \rangle$ to log

► IS THIS CORRECT??

32

Recovery rules: Undo logging

- (1) Let S = set of transactions with $\langle T_i, \text{start} \rangle$ in log, but no $\langle T_i, \text{commit} \rangle$ (or $\langle T_i, \text{abort} \rangle$) record in log
- (2) For each $\langle T_i, X, v \rangle$ in log, in reverse order (latest \rightarrow earliest) do:
 - if $T_i \in S$ then { - write (X, v)
 - output (X)
- (3) For each $T_i \in S$ do
 - write $\langle T_i, \text{abort} \rangle$ to log

33

What if failure during recovery?

No problem! = Undo idempotent

34

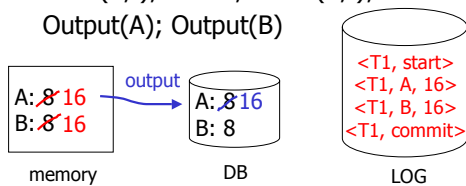
To discuss:

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

35

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)



36

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

37

Recovery rules: Redo logging

- For every T_i with $\langle T_i, \text{commit} \rangle$ in log:
 - For all $\langle T_i, X, v \rangle$ in log:
 - { Write(X, v)
 - { Output(X)

► IS THIS CORRECT??

38

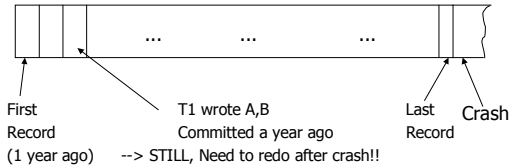
Recovery rules: Redo logging

- (1) Let S = set of transactions with $\langle T_i, \text{commit} \rangle$ in log
- (2) For each $\langle T_i, X, v \rangle$ in log, in forward order (earliest \rightarrow latest) do:
 - if $T_i \in S$ then { Write(X, v)
 - { Output(X) ← optional

39

Recovery is very, very **SLOW !**

Redo log:



40

Solution: Checkpoint (simple version)

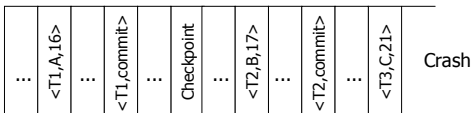
Periodically:

- (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

41

Example: what to do at recovery?

Redo log (disk):



42

Key drawbacks:

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

43

Solution: undo/redo logging!

Update \Rightarrow $\langle T_i, X_{id}, \text{New } X \text{ val}, \text{Old } X \text{ val} \rangle$
page X

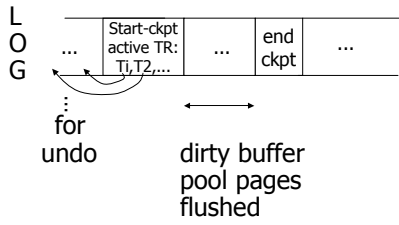
44

Rules

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

45

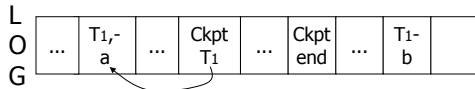
Non-quiet checkpoint



46

Examples what to do at recovery time?

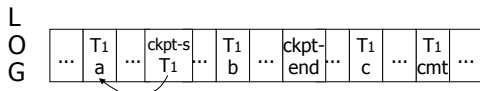
no T1 commit



► Undo T1 (undo a,b)

47

Example

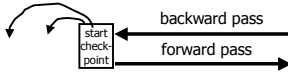


► Redo T1: (redo b,c)

48

Recovery process:

- **Backwards pass** (end of log → latest 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



49

Real world actions

E.g., dispense cash at ATM

$T_i = a_1 a_2 \dots a_j \dots a_n$

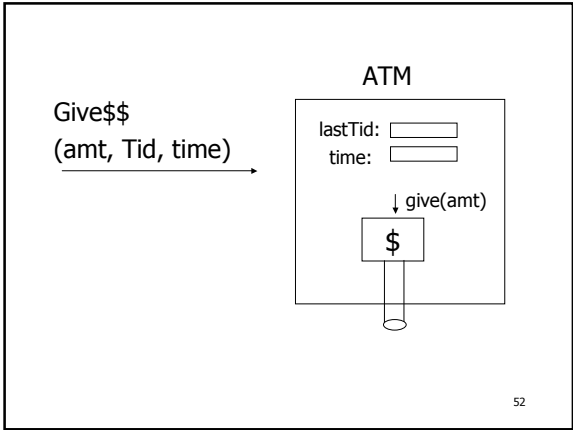
↓
\$

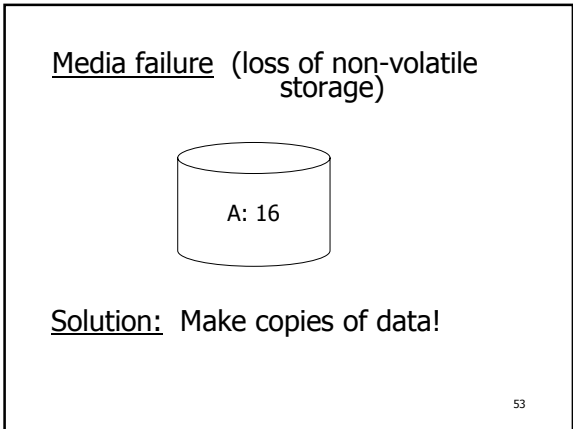
50

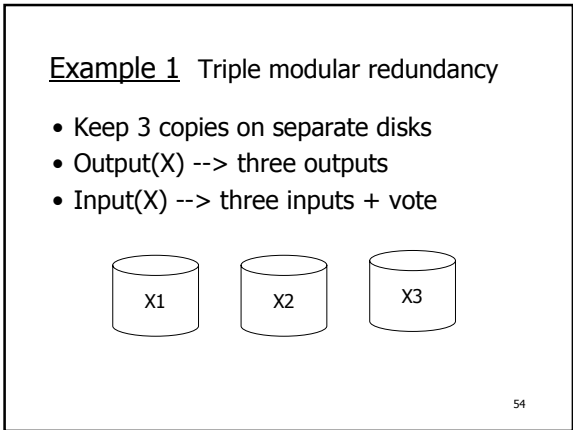
Solution

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

51





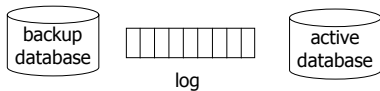


Example #2 Redundant writes, Single reads

- 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

55

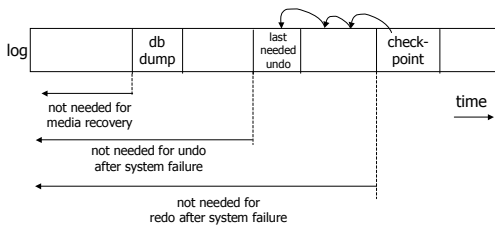
Example #3: DB Dump + Log



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

56

When can log be discarded?



57

Summary

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