

Chapter 17: Coping with System Failures

(Slides by Hector Garcia-Molina,
<http://www-db.stanford.edu/~hector/cs245/notes.htm>)

Chapter 17

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

Chapter 17

2

Integrity or consistency constraints

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

Chapter 17

3

Definition:

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

Chapter 17

4

Observation: DB cannot be consistent always!

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

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

Chapter 17

5

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

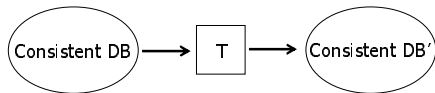
Deposit \$100 in a_2 : $a_2 \leftarrow a_2 + 100$
 $\text{TOT} \leftarrow \text{TOT} + 100$

	<table><tr><td>:</td></tr><tr><td>50</td></tr></table>	:	50	→	<table><tr><td>:</td></tr><tr><td>150</td></tr></table>	:	150	→	<table><tr><td>:</td></tr><tr><td>150</td></tr></table>	:	150
:											
50											
:											
150											
:											
150											
TOT	<table><tr><td>:</td></tr><tr><td>1000</td></tr></table>	:	1000		<table><tr><td>:</td></tr><tr><td>1000</td></tr></table>	:	1000		<table><tr><td>:</td></tr><tr><td>1100</td></tr></table>	:	1100
:											
1000											
:											
1000											
:											
1100											

Chapter 17

6

Transaction: collection of actions
that preserve consistency



Chapter 17

7

Big assumption:

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

Chapter 17

8

Correctness (informally)

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

Chapter 17

9

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

Chapter 17

10

We 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

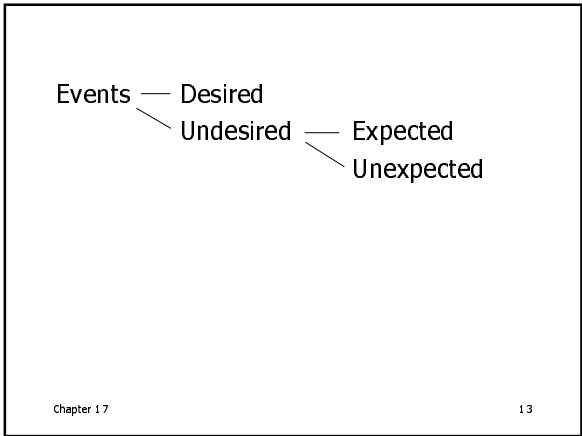
11

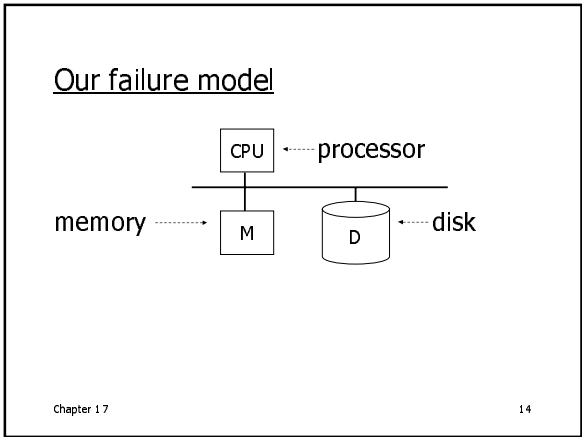
Chapter 17 Recovery Management

- First order of business:
Failure Model

Chapter 17

12





Desired events: see product manuals....

Undesired expected events:

- System crash
 - memory lost
 - cpu halts, resets

that's it!!

Undesired Unexpected: Everything else!

Chapter 17 15

Undesired Unexpected: Everything else!

Examples:

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

Chapter 17

16

Is this model reasonable?

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

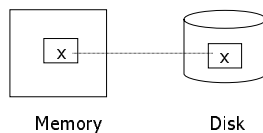
E.g., { Replicate disk storage (stable store)
Memory parity
CPU checks

Chapter 17

17

Second order of business:

Storage hierarchy



Chapter 17

18

Operations:

- Input (x): block with $x \rightarrow$ memory
- Output (x): block with $x \rightarrow$ disk
- Read (x,t): do input(x) if necessary
 $t \leftarrow$ value of x in block
- Write (x,t): do input(x) if necessary
value of x in block $\leftarrow t$

Chapter 17

19

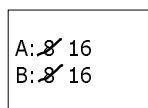
Key problem Unfinished transaction

Example Constraint: $A=B$
 $T_1: A \leftarrow A \times 2$
 $B \leftarrow B \times 2$

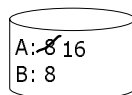
Chapter 17

20

T_1 : Read (A,t); $t \leftarrow t \times 2$
 Write (A,t);
 Read (B,t); $t \leftarrow t \times 2$
 Write (B,t);
 ~~Output (A);~~
 Output (B); failure!



memory



disk

Chapter 17

21

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

Chapter 17

22

One solution: undo logging (immediate modification)

due to: Hansel and Gretel, 782 AD

- Improved in 784 AD to durable undo logging

Chapter 17

23

Undo logging (Immediate modification)

T1: Read (A,t); $t \leftarrow t \times 2$ A=B
 Write (A,t);
 Read (B,t); $t \leftarrow t \times 2$
 Write (B,t);
 Output (A);
 Output (B);

A: ~~8~~ 16
 B: ~~8~~ 16

memory

A: ~~8~~ 16
 B: ~~8~~ 16

disk

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

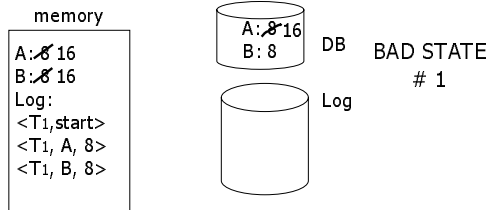
log

Chapter 17

24

One "complication"

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

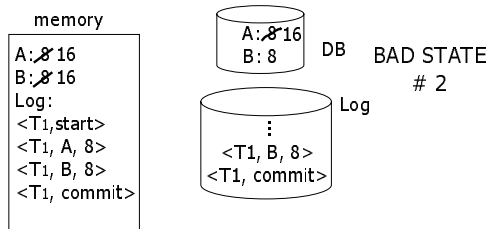


Chapter 17

25

One "complication"

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



Chapter 17

26

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

Chapter 17

27

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 $\left\{ \begin{array}{l} \text{- write } (X, v) \\ \text{- output } (X) \end{array} \right.$
- (3) For each $T_i \in S$ do
 - write $\langle T_i, \text{abort} \rangle$ to log

Chapter 17

28

What if failure during recovery?

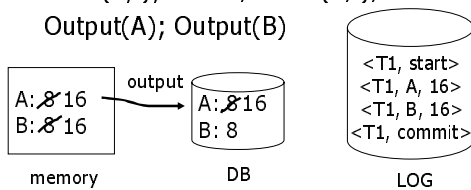
No problem! \Rightarrow Undo idempotent

Chapter 17

29

Redo logging (deferred modification)

T_1 : Read(A, t); $t \leftarrow t \times 2$; write (A, t);
 Read(B, t); $t \leftarrow t \times 2$; write (B, t);
 Output(A); Output(B)



Chapter 17

30

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

Chapter 17

31

Recovery rules: Redo logging

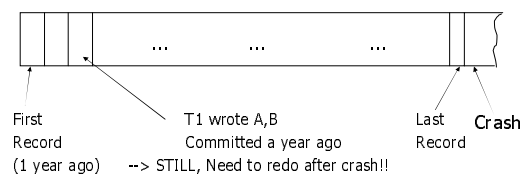
- (1) Let S = set of transactions with $\langle T_i, \text{commit} \rangle$ in log
- (2) For each $\langle T_i, X, v \rangle$ in log, in forward order (earliest \rightarrow latest) do:
 - if $T_i \in S$ then $\left\{ \begin{array}{l} \text{Write}(X, v) \\ \text{Output}(X) \leftarrow \text{optional} \end{array}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

Chapter 17

32

Recovery is very, very SLOW !

Redo log:



Chapter 17

33

Solution: Checkpoint (simple version)

Periodically:

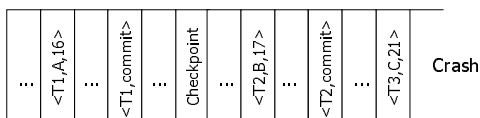
- (1) Do not accept new transactions
- (2) Wait until all transactions finish
- (3) Flush all log records to disk (log)
- (4) Flush all buffers to disk (DB) (do not discard buffers)
- (5) Write "checkpoint" record on disk (log)
- (6) Resume transaction processing

Chapter 17

34

Example: what to do at recovery?

Redo log (disk):



Chapter 17

35

Comparison of undo and redo logging

Key drawbacks:

- *Undo logging*: must flush data to disk immediately after a transaction finishes
 - Why is it bad?
- *Redo logging*: need to keep all modified blocks in memory until commit
 - Why is it bad?

Chapter 17

36

Solution: undo/redo logging!

Update \Rightarrow $\langle T_i, Xid, Old\ X\ val, New\ X\ val \rangle$
page X

Chapter 17

37

Rules

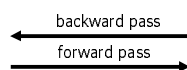
- Page X can be flushed before or after T_i commit
- Log record flushed before corresponding updated page (WAL)
- Flush at commit (log only)

Chapter 17

38

Recovery process:

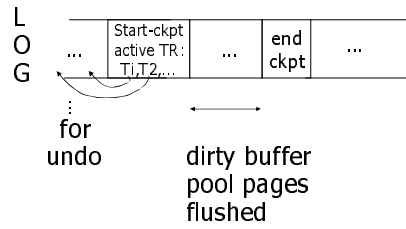
- **Backwards pass** (end of log \rightarrow latest checkpoint start)
 - construct set S of committed transactions
 - undo actions of transactions not in S
- **Forward pass** (latest checkpoint start \rightarrow end of log)
 - redo actions of S transactions



Chapter 17

39

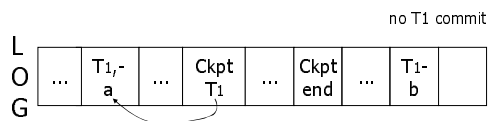
Non-quiet checkpoint



Chapter 17

40

Examples what to do at recovery time?

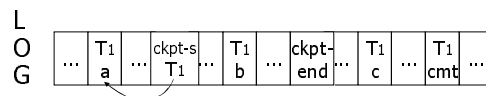


► Undo T1 (undo a,b)

Chapter 17

41

Example



► Redo T1: (redo b,c)

Chapter 17

42

Nonquiescent checkpoint

- The rules: section 17.4.3
- Example 17.12 p. 906

Chapter 17

43

-
-
-
-
-
-

43

Real world actions

E.g., dispense cash at ATM

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

↓

\$

Chapter 17

44

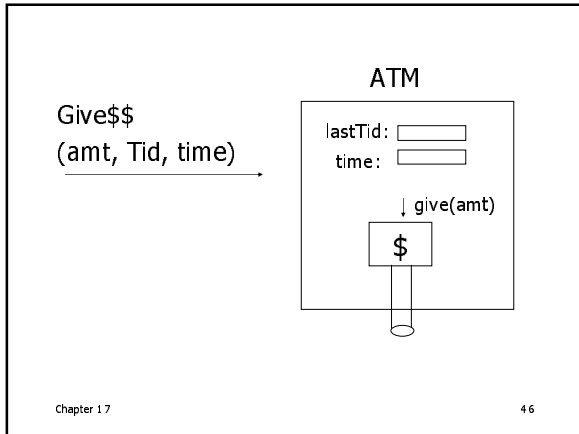
Solution

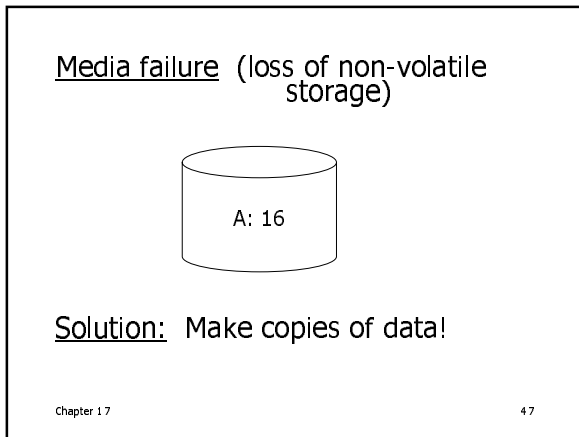
execute real-world actions after commit

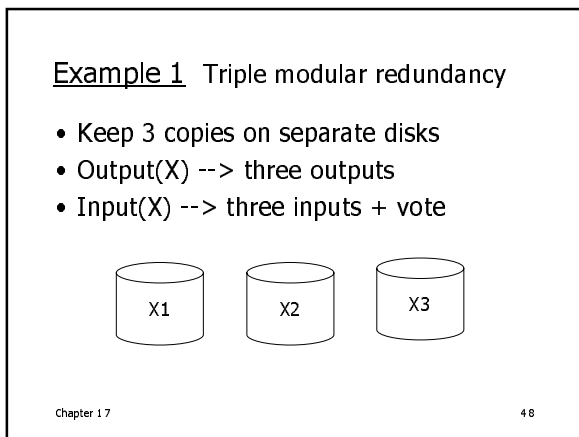
Chapter 17

45

45







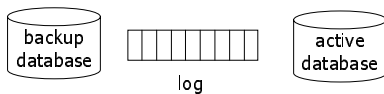
Example #2 Redundant writes, Single reads

- 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

Chapter 17

49

Example #3: DB Dump + Log

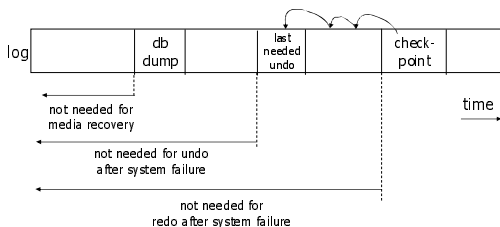


- If active database is lost,
 - restore active database from backup
 - bring up-to-date using redo entries in log

Chapter 17

50

When can log be discarded?



Chapter 17

51

Summary

- Consistency of data
- One source of problems: failures
 - Logging
 - Redundancy
- Another source of problems:
Data Sharing..... Next: Chapter 18

Chapter 17

52
