

CS 245: Database System Principles

Notes 08: Failure Recovery

Hector Garcia-Molina

CS 245

Notes 08

1

PART II

- Crash recovery (2 lectures) Ch.17[17]
- Concurrency control (3 lectures) Ch.18[18]
- Transaction processing (2 lects) Ch.19[19]
- Information integration (1 lect) Ch.20[21,22]

CS 245

Notes 08

2

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

CS 245

Notes 08

3

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

CS 245

Notes 08

4

Definition:

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

CS 245

Notes 08

5

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

CS 245

Notes 08

6

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

account

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

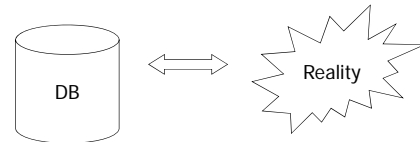
CS 245

Notes 08

7

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

Example 2 Database should reflect real world



CS 245

Notes 08

8

☞ 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 : $\begin{cases} a_2 \leftarrow a_2 + 100 \\ \text{TOT} \leftarrow \text{TOT} + 100 \end{cases}$

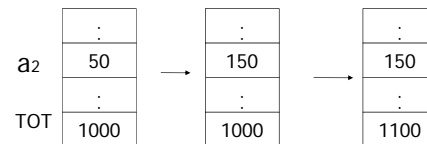
CS 245

Notes 08

9

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

Deposit \$100 in a_2 : $\begin{aligned} a_2 &\leftarrow a_2 + 100 \\ \text{TOT} &\leftarrow \text{TOT} + 100 \end{aligned}$

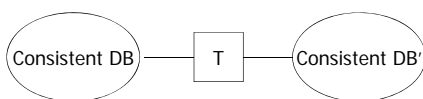


CS 245

Notes 08

10

Transaction: collection of actions that preserve consistency



CS 245

Notes 08

11

Big assumption:

If T starts with consistent state +

T executes in isolation

⇒ T leaves consistent state

CS 245

Notes 08

12

Correctness (informally)

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

CS 245

Notes 08

13

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

CS 245

Notes 08

14

How can we prevent/fix violations?

- Chapter 8[17]: due to failures only
- Chapter 9[18]: due to data sharing only
- Chapter 10[19]: due to failures and sharing

CS 245

Notes 08

15

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

CS 245

Notes 08

16

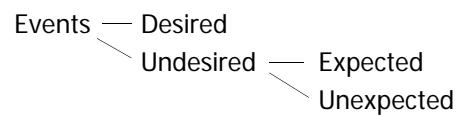
Chapter 8[17]: Recovery

- First order of business:
Failure Model

CS 245

Notes 08

17

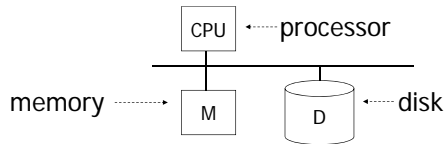


CS 245

Notes 08

18

Our failure model



CS 245

Notes 08

19

Desired events: see product manuals....

Undesired expected events:

System crash

- memory lost
- cpu halts, resets

CS 245

Notes 08

20

Desired events: see product manuals....

Undesired expected events:

System crash

- memory lost
- cpu halts, resets

— that's it!! —

Undesired Unexpected: Everything else!

CS 245

Notes 08

21

Undesired Unexpected: Everything else!

Examples:

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

CS 245

Notes 08

22

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

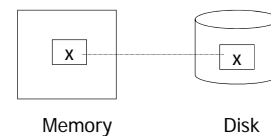
CS 245

Notes 08

23

Second order of business:

Storage hierarchy



CS 245

Notes 08

24

Operations:

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

CS 245

Notes 08

25

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

CS 245

Notes 08

26

Key problem Unfinished transaction

Example

Constraint: A=B

T1: A ← A × 2

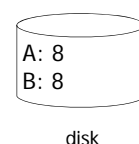
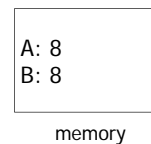
B ← B × 2

CS 245

Notes 08

27

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

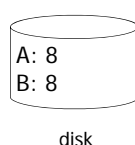
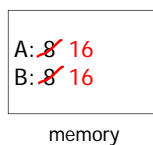


CS 245

Notes 08

28

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

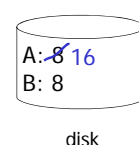
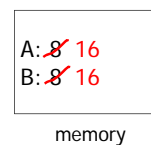


CS 245

Notes 08

29

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



CS 245

Notes 08

30

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

CS 245

Notes 08

31

One solution: undo logging (immediate modification)

due to: Hansel and Gretel, 1812 AD

CS 245

Notes 08

32

One solution: undo logging (immediate modification)

due to: Hansel and Gretel, 1812 AD

- Improved in 1813 AD to durable undo logging

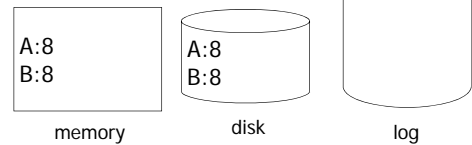
CS 245

Notes 08

33

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);



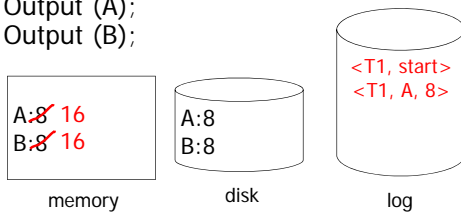
CS 245

Notes 08

34

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);



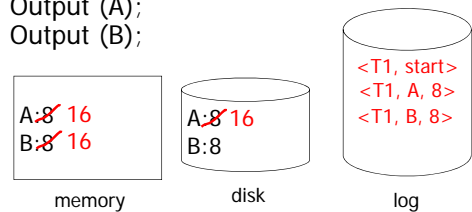
CS 245

Notes 08

35

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);



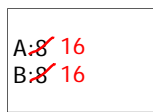
CS 245

Notes 08

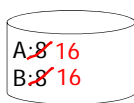
36

Undo logging (Immediate modification)

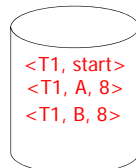
T1: Read (A,t); $t \leftarrow t \times 2$ A=B
 Write (A,t);
 Read (B,t); $t \leftarrow t \times 2$
 Write (B,t);
 Output (A);
 Output (B);



memory



disk



log

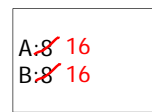
CS 245

Notes 08

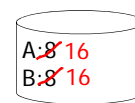
37

Undo logging (Immediate modification)

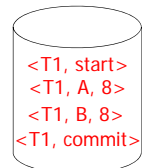
T1: Read (A,t); $t \leftarrow t \times 2$ A=B
 Write (A,t);
 Read (B,t); $t \leftarrow t \times 2$
 Write (B,t);
 Output (A);
 Output (B);



memory



disk



log

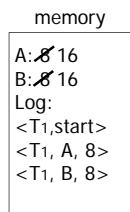
CS 245

Notes 08

38

One "complication"

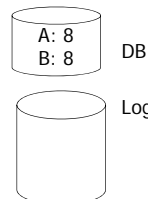
- Log is first written in memory
- Not written to disk on every action



CS 245

Notes 08

39

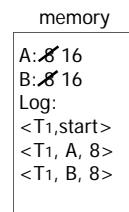


DB

Log

One "complication"

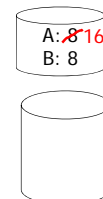
- Log is first written in memory
- Not written to disk on every action



CS 245

Notes 08

40



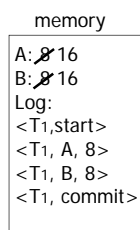
DB

Log

BAD STATE
1

One "complication"

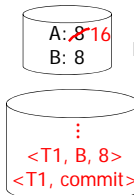
- Log is first written in memory
- Not written to disk on every action



CS 245

Notes 08

41



DB

Log

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

CS 245

Notes 08

42

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

CS 245

Notes 08

43

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??

CS 245

Notes 08

44

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

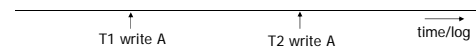
CS 245

Notes 08

45

Question

- Can writes of $\langle T_i, \text{abort} \rangle$ records be done in any order (in Step 3)?
 - Example: T_1 and T_2 both write A
 - T_1 executed before T_2
 - T_1 and T_2 both rolled-back
 - $\langle T_1, \text{abort} \rangle$ written but NOT $\langle T_2, \text{abort} \rangle$?
 - $\langle T_2, \text{abort} \rangle$ written but NOT $\langle T_1, \text{abort} \rangle$?



CS 245

Notes 08

46

What if failure during recovery?

No problem! ➡ Undo idempotent

CS 245

Notes 08

47

To discuss:

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

CS 245

Notes 08

48

Redo Logging



First send Gretel up with no rope,
then Hansel goes up safely with rope!



CS 245

Notes 08

49

Redo logging (deferred modification)

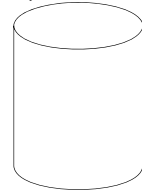
T1: Read(A,t); t ← t×2; write (A,t);
Read(B,t); t ← t×2; write (B,t);
Output(A); Output(B)

A: 8
B: 8

memory

A: 8
B: 8

DB



LOG

CS 245

Notes 08

50

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

memory

A: 8
B: 8

DB

<T1, start>
<T1, A, 16>
<T1, B, 16>
<T1, commit>

LOG

CS 245

Notes 08

51

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

memory

output

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

DB

<T1, start>
<T1, A, 16>
<T1, B, 16>
<T1, commit>

LOG

CS 245

Notes 08

52

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

memory

output

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

DB

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

LOG

CS 245

Notes 08

53

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

CS 245

Notes 08

54

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)

CS 245

Notes 08

55

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??

CS 245

Notes 08

56

Recovery rules: Redo logging

- Let S = set of transactions with $\langle T_i, \text{commit} \rangle$ (and no $\langle T_i, \text{end} \rangle$) in log
- 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)
- For each $T_i \in S$, write $\langle T_i, \text{end} \rangle$

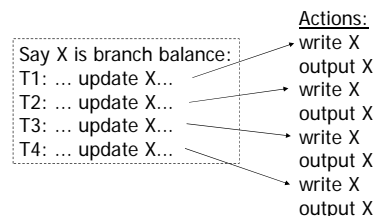
CS 245

Notes 08

57

Combining $\langle T_i, \text{end} \rangle$ Records

- Want to delay DB flushes for hot objects



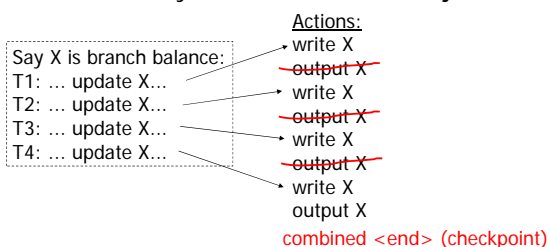
CS 245

Notes 08

58

Combining $\langle T_i, \text{end} \rangle$ Records

- Want to delay DB flushes for hot objects



CS 245

Notes 08

59

Solution: Checkpoint

- no $\langle t_i, \text{end} \rangle$ actions
- simple checkpoint

Periodically:

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

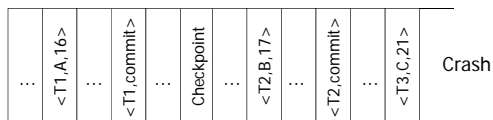
CS 245

Notes 08

60

Example: what to do at recovery?

Redo log (disk):



CS 245

Notes 08

61

Key drawbacks:

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

CS 245

Notes 08

62

Solution: undo/redo logging!

Update \Rightarrow <Ti, Xid, New X val, Old X val>
page X

CS 245

Notes 08

63

Rules

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

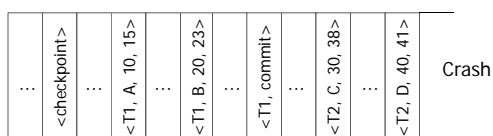
CS 245

Notes 08

64

Example: Undo/Redo logging what to do at recovery?

log (disk):

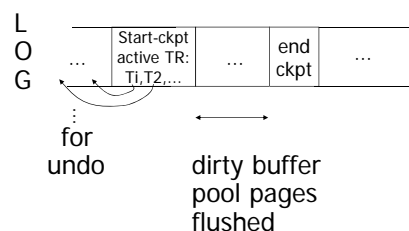


CS 245

Notes 08

65

Non-quiesce checkpoint



CS 245

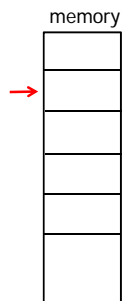
Notes 08

66

Non-quietse checkpoint

checkpoint process:
for $i := 1$ to M do
 output(buffer i)

[transactions run concurrently]



CS 245

Notes 08

67

Examples what to do at recovery time?

no T1 commit



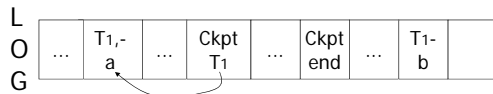
CS 245

Notes 08

68

Examples what to do at recovery time?

no T1 commit



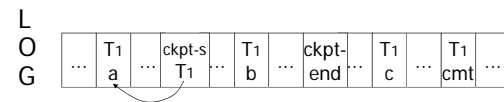
☒ Undo T1 (undo a,b)

CS 245

Notes 08

69

Example

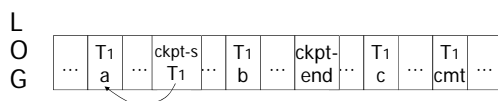


CS 245

Notes 08

70

Example



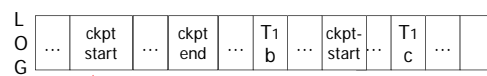
☒ Redo T1: (redo b,c)

CS 245

Notes 08

71

Recover From Valid Checkpoint:



start
of latest
valid
checkpoint

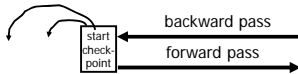
CS 245

Notes 08

72

Recovery process:

- **Backwards pass** (end of log \Rightarrow 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 \Rightarrow end of log)
 - redo actions of S transactions



CS 245

Notes 08

73

Real world actions

E.g., dispense cash at ATM

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

↓
\$

CS 245

Notes 08

74

Solution

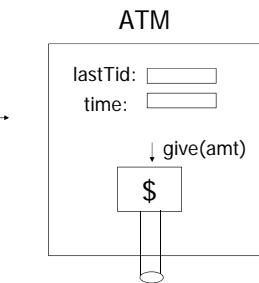
- (1) execute real-world actions after commit
- (2) try to make idempotent

CS 245

Notes 08

75

Give\$\$
(amt, Tid, time)

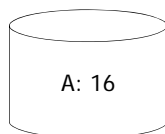


CS 245

Notes 08

76

Media failure (loss of non-volatile storage)

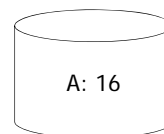


CS 245

Notes 08

77

Media failure (loss of non-volatile storage)



Solution: Make copies of data!

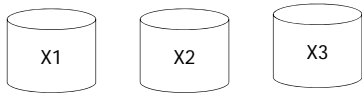
CS 245

Notes 08

78

Example 1 Triple modular redundancy

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



CS 245

Notes 08

79

Example #2 Redundant writes, Single reads

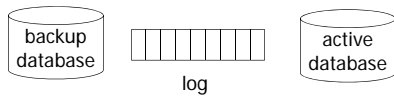
- Keep N copies on separate disks
 - Output(X) --> N outputs
 - Input(X) --> Input one copy
 - if ok, done
 - else try another one
- ⇔ Assumes bad data can be detected

CS 245

Notes 08

80

Example #3: DB Dump + Log



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

CS 245

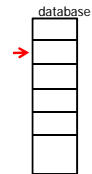
Notes 08

81

Backup Database

- Just like checkpoint, except that we write full database

create backup database:
 for i := 1 to DB_Size do
 [read DB block i; write to backup]
 [transactions run concurrently]



CS 245

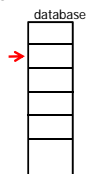
Notes 08

82

Backup Database

- Just like checkpoint, except that we write full database

create backup database:
 for i := 1 to DB_Size do
 [read DB block i; write to backup]
 [transactions run concurrently]



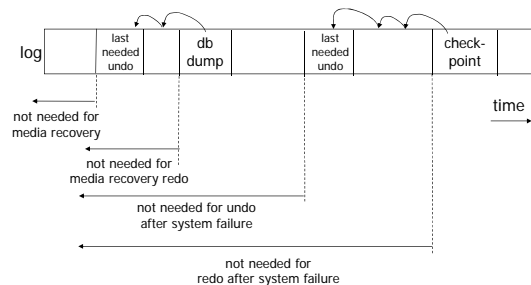
- Restore from backup DB and log:
 Similar to recovery from checkpoint and log

CS 245

Notes 08

83

When can log be discarded?



CS 245

Notes 08

84

Summary

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