

# CS 245: Database System Principles

## Notes 08: Failure Recovery

Hector Garcia-Molina  
(Some modifications by Chris Olston)

CS 245

Notes 08

1

## PART II

- Crash recovery (2 lectures) Ch.8
- Concurrency control (3 lectures) Ch.9
- Transaction processing (2 lects) Ch.10
- Information integration (1 lect) Ch.11

CS 245

Notes 08

2

## ACID Properties

- Atomicity
  - Actions are never left partially executed
- Consistency
  - Actions leave the DB in a consistent state
- Isolation
  - Actions are not affected by other concurrent actions
- Durability
  - Effects of completed actions are resilient against system failures

CS 245

Notes 08

3

## Integrity or correctness of data

A  
Consistency  
I  
D

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

EMP

Name	Age
White	52
Green	3421
Gray	1

CS 245

Notes 08

4

## Integrity or consistency constraints

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

CS 245

Notes 08

5

## Definition:

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

CS 245

Notes 08

6

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

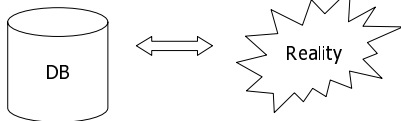
Note: could be "emulated" by simple constraints, e.g.,

account 

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

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

Example 2 Database should reflect real world



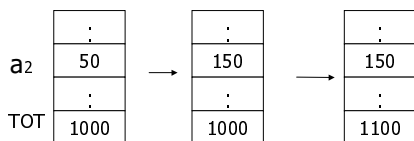
☞ in any case, continue with constraints...

Observation: DB cannot be consistent always!

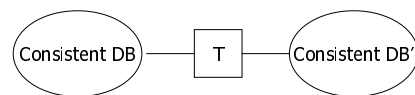
Example:  $a_1 + a_2 + \dots + a_n = \text{TOT}$  (constraint)  
 Deposit \$100 in  $a_2$ :  $\begin{cases} a_2 \leftarrow a_2 + 100 \\ \text{TOT} \leftarrow \text{TOT} + 100 \end{cases}$

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

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



Transaction: collection of actions that preserve consistency



### Big assumption:

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

CS 245

Notes 08

13

### Correctness (informally)

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

CS 245

Notes 08

14

### 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

CS 245

Notes 08

15

### 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

CS 245

Notes 08

16

### 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

17

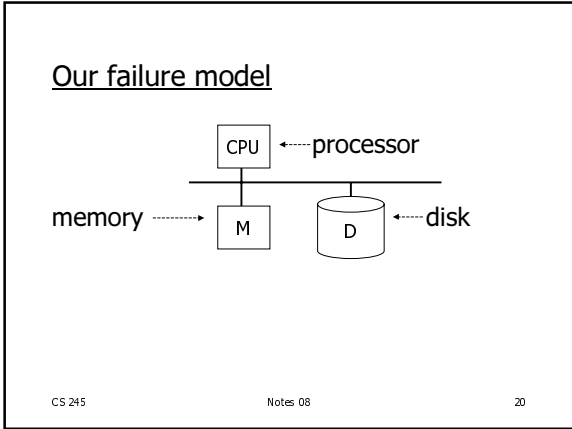
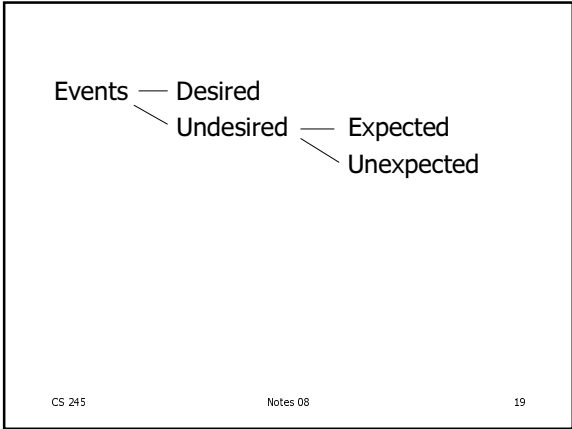
## Chapter 8 Recovery

- First order of business:  
Failure Model

CS 245

Notes 08

18



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 catches fire
- Car smashes into data center
- (You get the idea ...)

CS 245                      Notes 08                      22

Is this model reasonable?

In our model, data is either okay or "lost"

Assumes we can detect corruption and treat it as lost data

CS 245                      Notes 08                      23

Detecting corruption

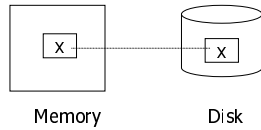
Approach: Add low level checks + redundancy to increase probability model holds

E.g., { Disk checksums  
 Memory parity  
 CPU checks: run code twice

CS 245                      Notes 08                      24

## Second order of business:

### Storage hierarchy



CS 245

Notes 08

25

## 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

CS 245

Notes 08

26

## Failure modes:

- Undesired/expected: system crash  
– Data on disk still there on restart  $\leftarrow$  next
- Undesired/unexpected:  
media failure  
catastrophic failure  
– Data on disk/disks lost!  $\leftarrow$  later

CS 245

Notes 08

27

## System Crash

A  
C  
I  
Durability

Problem #1 Transaction completed,  
but results only in memory

If system crashes, effects of transaction lost

Solution: Don't tell user the transaction has  
"committed" until all effects are on disk

CS 245

Notes 08

28

## System Crash

### Problem #2 Unfinished transaction

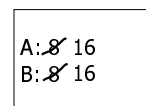
Example

Constraint:  $A=B$

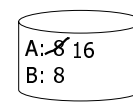
T1:  $A \leftarrow A \times 2$

$B \leftarrow B \times 2$

T1: 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); crash!



memory



disk

CS 245

Notes 08

29


CS 245

Notes 08

30

Atomicity  
C  
I  
D

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



CS 245 Notes 08 31

One variation: undo logging (immediate modification)

due to: Hansel and Gretel, 782 AD

- Improved in 784 AD to durable undo logging

CS 245 Notes 08 32

Undo logging (Immediate modification)

T1: Read (A,t); t ← tx2      A=B  
 Write (A,t);  
 Read (B,t); t ← tx2  
 Write (B,t);  
 Output (A);  
 Output (B);

A: 16  
B: 16

A: 16  
B: 16

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

memory                  disk                  log

CS 245 Notes 08 33

Assumption

- Logging “fakes” atomicity by undoing writes of partially completed xacts
- Assumption: writing a block is atomic
  - Atomicity of block writes is itself “faked” by low-level checks
  - And so it goes ...
  - At the bottom, there must be some intrinsically atomic action (writing a bit?)

CS 245 Notes 08 34

One “complication”

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

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

A: 16  
B: 8

Log

DB BAD STATE # 1

CS 245 Notes 08 35

One “complication”

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

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

A: 16  
B: 8

Log:  
:  
<T1, B, 8>  
<T1, commit>

DB BAD STATE # 2

CS 245 Notes 08 36

### 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

37

### 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

38

### 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

39

### What if failure during recovery?

No problem! = Undo idempotent

CS 245

Notes 08

40

### To discuss:

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

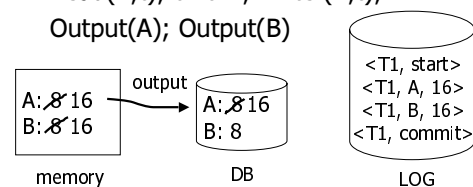
CS 245

Notes 08

41

### 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( $B$ )



CS 245

Notes 08

42

### 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

CS 245

Notes 08

43

### 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

44

### 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  $\left\{ \begin{array}{l} \text{Write}(X, v) \\ \text{Output}(X) \leftarrow \text{optional} \end{array} \right.$

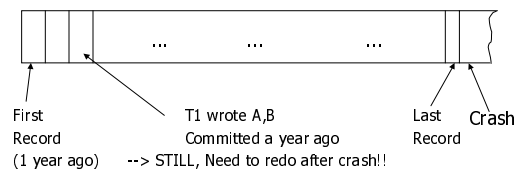
CS 245

Notes 08

45

### Recovery is very, very SLOW !

Redo log:



CS 245

Notes 08

46

### 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

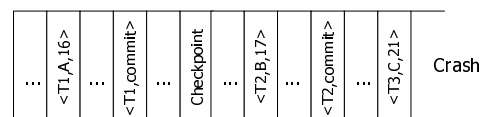
CS 245

Notes 08

47

### Example: what to do at recovery?

Redo log (disk):



CS 245

Notes 08

48



## Better: Fuzzy Checkpointing

- Want to be able to checkpoint without bringing the system to a halt
- Especially if transactions run for a long time
- Solution: Nonquiescent checkpointing
  - Will get to this soon...

CS 245

Notes 08

49

## Key drawbacks (with respect to buffer management)

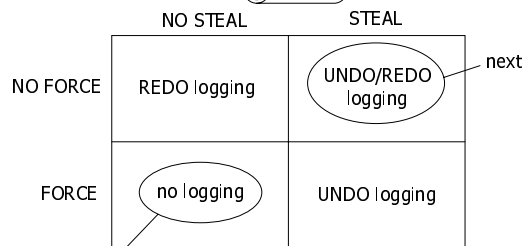
- *Undo logging:*
  - must FORCE data to disk to commit
- *Redo logging:*
  - need to keep all modified blocks in memory until commit
  - (other xacts can't STEAL buffer pages)

CS 245

Notes 08

50

## Logging Taxonomy



Only works if xacts only modify one block, and no two concurrent xacts modify same block

CS 245

Notes 08

51

## Solution: undo/redo logging!

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

CS 245

Notes 08

52

## Undo/Redo Logging Rules

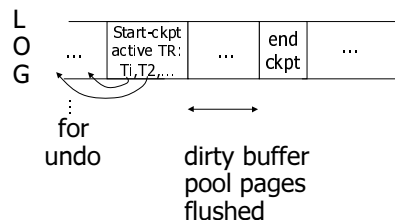
- Element X can be flushed before or after  $T_i$  commit
  - (Buffer manager has complete freedom)
- Before flushing a page, all corresponding log records flushed (WAL)
- Flush log before commit

CS 245

Notes 08

53

## Non-quiesce checkpoint

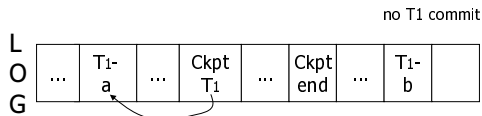


CS 245

Notes 08

54

## Examples what to do at recovery time?



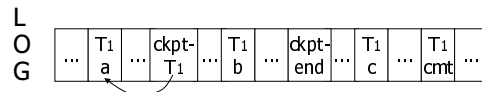
► Undo T1 (undo a,b)

CS 245

Notes 08

55

## Example



► Redo T1: (redo b,c)

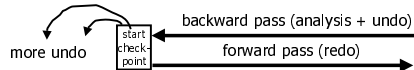
CS 245

Notes 08

56

## Recovery process:

- Backward pass (end of log → latest checkpoint start)
  - construct set C of transactions that committed since checkpoint start
  - Undo actions of uncommitted transactions
- Finish undo (follow undo chains)
  - Undo all actions of transactions not in C
- Forward pass (latest checkpoint start → end of log)
  - Redo actions of C transactions



CS 245

Notes 08

57

## Real world actions

E.g., fire missile

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



Can't undo!  
(Well, not yet anyway ... ☺)

CS 245

Notes 08

58

## Solution

Execute real-world actions after commit

Works for actions that are idempotent

If missile already launched, "fire missile" command has no effect (assuming one missile per launcher/silo)

CS 245

Notes 08

59

## Non-idempotent real world actions

E.g., dispense cash at ATM

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

↓  
\$

Can't undo!  
("Please reinsert cash" ☺)

Not idempotent!  
(Avoid dispensing cash twice)

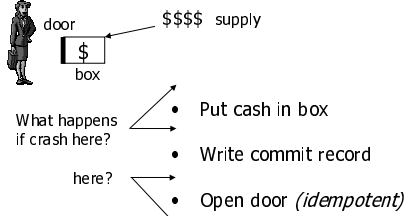
CS 245

Notes 08

60

## Solution

Design system & transaction so final action is idempotent



CS 245

Notes 08

61

## Failure modes:

- Undesired/expected: system crash
  - Data on disk still there on restart
  - Solution: atomicity via logging
- Undesired/unexpected:
  - media failure
  - catastrophic failure
  - Data on disk/disks lost! ← next

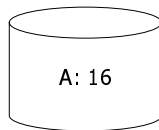
CS 245

Notes 08

62

## Media failure

(loss of non-volatile storage)



Solution: Make copies of data!

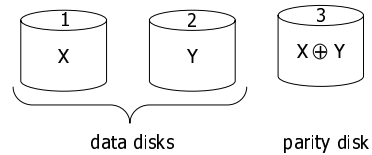
CS 245

Notes 08

63

## Example 1 RAID

- Array of disks with redundancy
- RAID-4: one extra disk for "parity"
- Output(X) --> write to disk 1, update 3
- Input(X) --> read disk 1



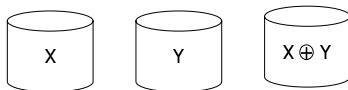
CS 245

Notes 08

64

## Recovery in RAID-4

- If one disk fails, reconstruct it using the other two disks
- ⇒ Assumes bad data can be detected



CS 245

Notes 08

65

## Example #2 Fully 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

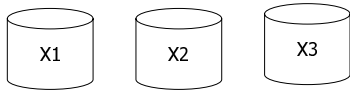
Notes 08

66

### Example 3 Triple modular redundancy

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

⇒ Protects against undetected bad data

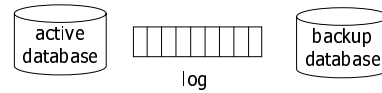


CS 245

Notes 08

67

### Example 4: DB Dump + Log



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

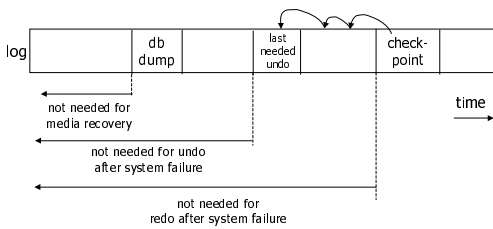
☆ Can't use undo-only logging!

CS 245

Notes 08

68

### When can log be discarded?



CS 245

Notes 08

69

### Catastrophic Failure

- Array of disks won't help in case of:

- Fire
- Explosion
- Earthquake
- Godzilla

- Also: vandalism, viruses

CS 245

Notes 08

70

### Solution

A  
C  
I  
Durability++

- Geographically distributed copies!



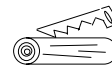
CS 245

Notes 08

71

### Summary

- Consistency of data
- One source of problems: failures
  - Logging



- Redundancy



- Another source of problems: Data Sharing..... next

CS 245

Notes 08

72