




CS 525: Advanced Database Organization

13: Failure and Recovery

Boris Glavic



Slides: adapted from a [course](#) taught by [Hector Garcia-Molina](#), Stanford InfoLab



CS 525  Notes 13 - Failure and Recovery 1  IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY



Now

- Crash recovery

CS 525  Notes 13 - Failure and Recovery 2  IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY



Correctness (informally)

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

CS 525  Notes 13 - Failure and Recovery 3  IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY



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 525  Notes 13 - Failure and Recovery 4  IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Recovery



- First order of business: Failure Model

CS 525  Notes 13 - Failure and Recovery 5  IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Events — Desired

 — Undesired — Expected

 — Unexpected

CS 525  Notes 13 - Failure and Recovery 6  IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Our failure model

The diagram shows a central horizontal line representing a system bus. Above the bus is a box labeled 'CPU' with a dashed arrow pointing to it from the text 'processor'. Below the bus are two components: a box labeled 'M' with a dashed arrow pointing to it from the text 'memory' on the left, and a cylinder labeled 'D' with a dashed arrow pointing to it from the text 'disk' on the right.

CS 525 Notes 13 - Failure and Recovery 7 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Desired events: see product manuals....

Undesired expected events:

System crash

- memory lost
- cpu halts, resets

CS 525 Notes 13 - Failure and Recovery 8 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Desired events: see product manuals....

Undesired expected events:

System crash

- memory lost
- cpu halts, resets

————— that's it!! —————

Undesired Unexpected: Everything else!

CS 525 Notes 13 - Failure and Recovery 9 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Undesired Unexpected: Everything else!

Examples:

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

CS 525 Notes 13 - Failure and Recovery 10 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

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 525 Notes 13 - Failure and Recovery 11 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Second order of business:

Storage hierarchy

The diagram shows a box labeled 'Memory' containing a smaller box with an 'X', and a cylinder labeled 'Disk' also containing a smaller box with an 'X'. A dashed line connects the two boxes. Below the Memory box is the text 'DB Buffer'.

CS 525 Notes 13 - Failure and Recovery 12 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Operations:

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

CS 525



Notes 13 - Failure and Recovery

13

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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 525



Notes 13 - Failure and Recovery

14

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Key problem Unfinished transaction

Example

Constraint: A=B

T1: A ← A × 2

B ← B × 2

CS 525

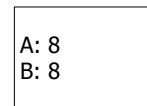


Notes 13 - Failure and Recovery

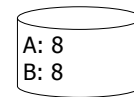
15

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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



memory



disk

CS 525

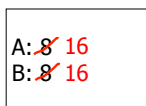


Notes 13 - Failure and Recovery

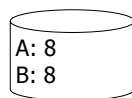
16

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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



memory



disk

CS 525

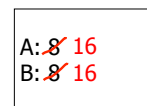


Notes 13 - Failure and Recovery

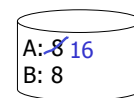
17

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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

CS 525





Notes 13 - Failure and Recovery

18



IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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

CS 525  Notes 13 - Failure and Recovery 19 



How to restore consistent state after crash?

- Desired state after recovery:
 - Changes of committed transactions are reflected on disk
 - Changes of unfinished transactions are not reflected on disk
- After crash we need to
 - **Undo** changes of unfinished transactions that have been written to disk
 - **Redo** changes of finished transactions that have not been written to disk

CS 525  Notes 13 - Failure and Recovery 20 

How to restore consistent state after crash?

- After crash we need to
 - **Undo** changes of unfinished transactions that have been written to disk
 - **Redo** changes of finished transactions that have not been written to disk
- We need to either
 - Store additional data to be able to Undo/Redo
 - Avoid ending up in situations where we need to Undo/Redo



CS 525  Notes 13 - Failure and Recovery 21 

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

T₁ is unfinished
 -> need to undo the write to A to recover to consistent state


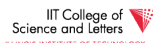
A: ~~8~~ 16
 B: 8

memory
disk

CS 525  Notes 13 - Failure and Recovery 22 



Logging

- After crash need to
 - **Undo**
 - **Redo**
- We need to know
 - Which operations have been executed
 - Which operations are reflected on disk
- -> **Log** upfront what is to be done

CS 525  Notes 13 - Failure and Recovery 23 

Buffer Replacement Revisited

- Now we are interested in knowing how buffer replacement influences recovery!

CS 525  Notes 13 - Failure and Recovery 24 

Buffer Replacement Revisited

- **Steal:** all pages with fix count = 0 are replacement candidates
 - Smaller buffer requirements
- **No steal:** pages that have been modified by active transaction -> not considered for replacement
 - No need to undo operations of unfinished transactions after failure

Buffer Replacement Revisited

- **Force:** Pages modified by transaction are flushed to disk at end of transaction
 - No redo required
- **No force:** modified (dirty) pages are allowed to remain in buffer after end of transaction
 - Less repeated writes of same page

Effects of Buffer Replacement

	force	No force
No steal	<ul style="list-style-type: none"> • No Undo • No Redo 	<ul style="list-style-type: none"> • No Undo • Redo
steal	<ul style="list-style-type: none"> • Undo • No Redo 	<ul style="list-style-type: none"> • Redo • Undo

Schedules and Recovery

- Are there certain schedules that are easy/hard/impossible to recover from?

Recoverable Schedules

- We should never have to rollback an already committed transaction (D in ACID)
- **Recoverable** schedules require that
 - A transaction does not commit before every transaction that is has read from has committed
 - A transaction **T** reads from another transaction **T'** if it reads an item X that has last been written by T' and T' has not aborted before the read

$$T_1 = w_1(X), c_1$$

$$T_2 = r_2(X), w_2(X), c_2$$

Recoverable Schedule

$$S_1 = w_1(X), r_2(X), w_2(X), c_1, c_2$$

Nonrecoverable Schedule

$$S_2 = w_1(X), r_2(X), w_2(X), c_2, c_1$$

Cascading Abort

- Transaction **T** has written an item that is later read by **T'** and **T** aborts after that
 - we have to also abort **T'** because the value it read is no longer valid anymore
 - This is called a **cascading abort**
 - Cascading aborts are complex and should be avoided

$$S = \dots w_1(X) \dots r_2(X) \dots a_1$$

CS 525



Notes 13 - Failure and Recovery

31

Cascadeless Schedules

- Cascadeless** schedules guarantee that there are no cascading aborts
 - Transactions only read values written by already committed transactions

CS 525



Notes 13 - Failure and Recovery

32

$$T_1 = w_1(X), c_1$$

$$T_2 = r_2(X), w_2(X), c_2$$

Cascadeless Schedule

$$S_1 = w_1(X), c_1, r_2(X), w_2(X), c_2$$

Recoverable Schedule

$$S_2 = w_1(X), r_2(X), w_2(X), c_1, c_2$$

Nonrecoverable Schedule

$$S_3 = w_1(X), r_2(X), w_2(X), c_2, c_1$$

CS 525



Notes 12 - Transaction Management

33

$$T_1 = w_1(X), a_1$$

$$T_2 = r_2(X), w_2(X), c_2$$

Cascadeless Schedule

$$S_1 = w_1(X), a_1, r_2(X), w_2(X), c_2$$

Recoverable Schedule

$$S_2 = w_1(X), r_2(X), w_2(X), a_1, a_2$$

Nonrecoverable Schedule

$$S_3 = w_1(X), r_2(X), w_2(X), c_2, a_1$$

Consider what happens if T1 aborts!

CS 525



Notes 12 - Transaction Management

34

Strict Schedules

- Strict** schedules guarantee that to Undo the effect of a transaction we simply have to undo each of its writes
 - Transactions do not read nor write items written by uncommitted transactions

CS 525



Notes 13 - Failure and Recovery

35

$$T_1 = w_1(X), c_1$$

$$T_2 = r_2(X), w_2(X), c_2$$

Cascadeless Schedule

$$S_1 = w_1(X), c_1, r_2(X), w_2(X), c_2$$

Recoverable Schedule

$$S_2 = w_1(X), r_2(X), w_2(X), c_1, c_2$$

Nonrecoverable Schedule

$$S_3 = w_1(X), r_2(X), w_2(X), c_2, c_1$$

CS 525



Notes 12 - Transaction Management

36

Compare Classes

ST \subset CL \subset RC \subset ALL

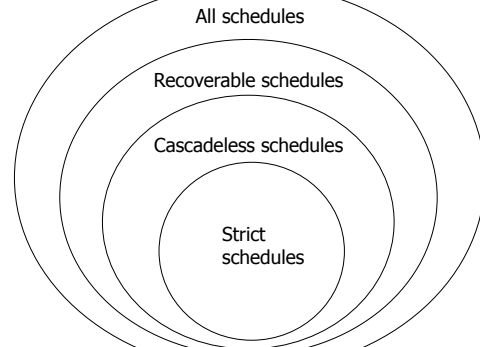
CS 525



Notes 13 - Failure and Recovery

37

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY



CS 525



Notes 13 - Failure and Recovery

38

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Logging and Recovery

- We now discuss approaches for logging and how to use them in recovery

CS 525



Notes 13 - Failure and Recovery

39

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

One solution: undo logging (immediate modification)

due to: Hansel and Gretel, 782 AD

CS 525



Notes 13 - Failure and Recovery

40

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

One solution: undo logging (immediate modification)

due to: Hansel and Gretel, 782 AD

- Improved in 784 AD to durable undo logging

CS 525



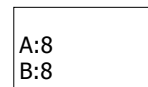
Notes 13 - Failure and Recovery

41

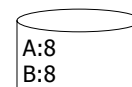
IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Undo logging (Immediate modification)

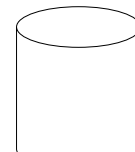
T1: Read (A,t); $t \leftarrow t \times 2$ A=B
Write (A,t);
Read (B,t); $t \leftarrow t \times 2$
Write (B,t);
Output (A);
Output (B);



memory



disk



log

CS 525



Notes 13 - Failure and Recovery

42

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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);

memory: A: 8, B: 8
 disk: A: 8, B: 8
 log: <T1, start>, <T1, A, 8>

CS 525 Notes 13 - Failure and Recovery 43 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

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);

memory: A: 8, B: 8
 disk: A: 8, B: 8
 log: <T1, start>, <T1, A, 8>, <T1, B, 8>

CS 525 Notes 13 - Failure and Recovery 44 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

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);

memory: A: 8, B: 8
 disk: A: 8, B: 8
 log: <T1, start>, <T1, A, 8>, <T1, B, 8>

CS 525 Notes 13 - Failure and Recovery 45 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

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);

memory: A: 8, B: 8
 disk: A: 8, B: 8
 log: <T1, start>, <T1, A, 8>, <T1, B, 8>, <T1, commit>

CS 525 Notes 13 - Failure and Recovery 46 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

One “complication”

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

memory: A: 8, B: 8, Log: <T1, start>, <T1, A, 8>, <T1, B, 8>
 DB: A: 8, B: 8
 Log: empty

CS 525 Notes 13 - Failure and Recovery 47 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

One “complication”

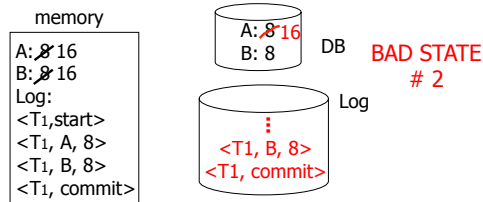
- Log is first written in memory
- Not written to disk on every action

memory: A: 8, B: 8, Log: <T1, start>, <T1, A, 8>, <T1, B, 8>
 DB: A: 8, B: 8
 Log: empty
BAD STATE # 1

CS 525 Notes 13 - Failure and Recovery 48 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

One “complication”

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



CS 525



Notes 13 - Failure and Recovery

49

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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 525



Notes 13 - Failure and Recovery

50

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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 525



Notes 13 - Failure and Recovery

51

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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 525



Notes 13 - Failure and Recovery

52

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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 525



Notes 13 - Failure and Recovery

53

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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 525



Notes 13 - Failure and Recovery

54

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

What if failure during recovery?

No problem! ⇒ Undo idempotent

- An operation is called **idempotent** if the number of times it is applied do not effect the result
- For Undo:
 - $\text{Undo}(\log) = \text{Undo}(\text{Undo}(\dots(\text{Undo}(\log)) \dots))$

Undo is idempotent

- We store the values of data items before the operation
- Undo can be executed repeatedly without changing effects
 - idempotent

Physical vs. Logical Logging

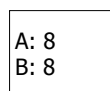
- How to represent values in log entries?
- Physical logging
 - Content of pages before and after
- Logical operations
 - Operation to execute for undo/redo
 - E.g., delete record x
- Hybrid (Physiological)
 - Delete record x from page y

To discuss:

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

Redo logging (deferred modification)

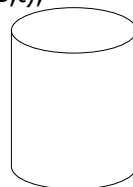
T1: Read(A,t); t- tx2; write (A,t);
Read(B,t); t- tx2; write (B,t);
Output(A); Output(B)



memory



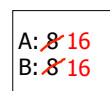
DB



LOG

Redo logging (deferred modification)

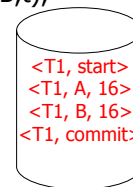
T1: Read(A,t); t- tx2; write (A,t);
Read(B,t); t- tx2; write (B,t);
Output(A); Output(B)



memory



DB



LOG

Redo logging (deferred modification)

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

memory DB LOG

CS 525 Notes 13 - Failure and Recovery 61 IIT College of Science and Letters

Redo logging (deferred modification)

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

memory DB LOG

CS 525 Notes 13 - Failure and Recovery 62 IIT College of Science and Letters

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 525 Notes 13 - Failure and Recovery 63 IIT College of Science and Letters

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 525 Notes 13 - Failure and Recovery 64 IIT College of Science and Letters

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 525 Notes 13 - Failure and Recovery 65 IIT College of Science and Letters

Recovery rules: Redo logging

- (1) Let S = set of transactions with $\langle T_i, \text{commit} \rangle$ (and no $\langle T_i, \text{end} \rangle$) in log
- (2) For each $\langle T_i, X, v \rangle$ in log, in forward order (earliest → latest) do:
 - if $T_i \in S$ then
 - Write(X, v)
 - Output(X)
- (3) For each $T_i \in S$, write $\langle T_i, \text{end} \rangle$

CS 525 Notes 13 - Failure and Recovery 66 IIT College of Science and Letters

Crash During Redo

- Since Redo log contains values after writes, repeated application of a log entry does not change result
--> idempotent

CS 525



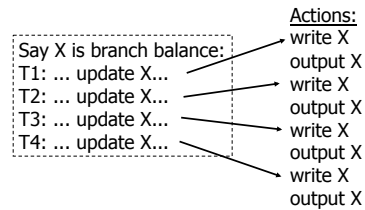
Notes 13 - Failure and Recovery

67

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Combining <Ti, end> Records

- Want to delay DB flushes for hot objects



CS 525



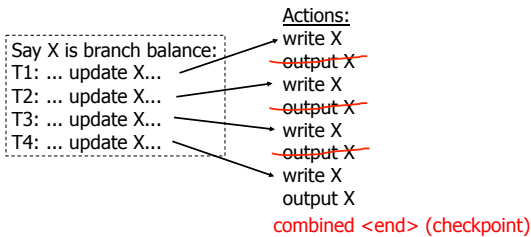
Notes 13 - Failure and Recovery

68

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Combining <Ti, end> Records

- Want to delay DB flushes for hot objects



CS 525



Notes 13 - Failure and Recovery

69

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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

CS 525



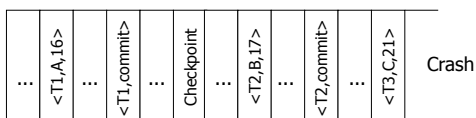
Notes 13 - Failure and Recovery

70

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Example: what to do at recovery?

Redo log (disk):



CS 525



Notes 13 - Failure and Recovery

71

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Advantage of Checkpoints

- Limits recovery to parts of the log after the checkpoint
 - Think about system that has been online for months
 - --> Analyzing the whole log is too expensive!
- Source of backups
 - If we backup checkpoints we can use them for media recovery!

CS 525



Notes 13 - Failure and Recovery

72

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Checkpoints Justification

- Checkpoint should be consistent DB state
 - No active transactions
 - Do not accept new transactions
 - Wait until all transactions finish
 - DB state reflected on disk
 - Flush log
 - Flush buffers

CS 525



Notes 13 - Failure and Recovery

73

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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 525



Notes 13 - Failure and Recovery

74

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Solution: undo/redo logging!

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

CS 525



Notes 13 - Failure and Recovery

75

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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)

CS 525



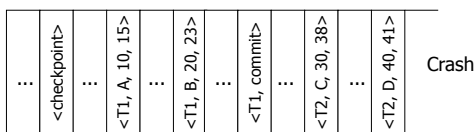
Notes 13 - Failure and Recovery

76

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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

log (disk):



CS 525



Notes 13 - Failure and Recovery

77

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Checkpoint Cost

- Checkpoints are expensive
 - No new transactions can start
 - A lot of I/O
 - Flushing the log
 - Flushing dirty buffer pages

CS 525

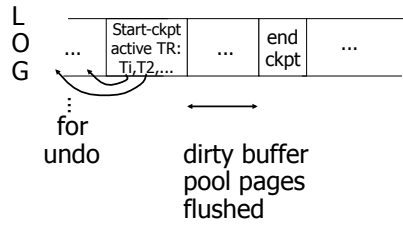


Notes 13 - Failure and Recovery

78

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Non-quietse checkpoint



CS 525

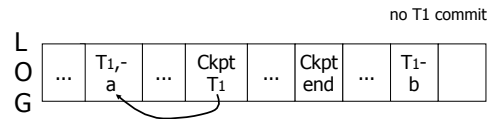


Notes 13 - Failure and Recovery

79

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Examples what to do at recovery time?



CS 525

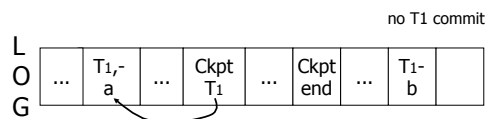


Notes 13 - Failure and Recovery

80

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Examples what to do at recovery time?



➡ Undo T1 (undo a,b)

CS 525

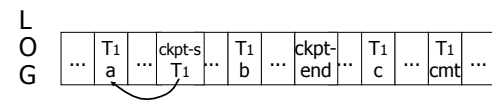


Notes 13 - Failure and Recovery

81

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Example



CS 525

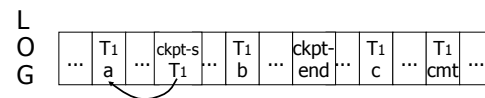


Notes 13 - Failure and Recovery

82

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Example



➡ Redo T1: (redo b,c)

CS 525

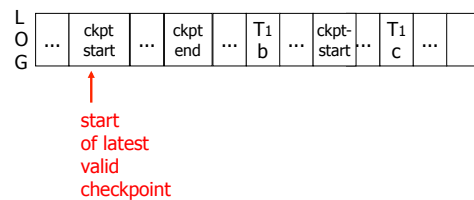


Notes 13 - Failure and Recovery

83

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Recover From Valid Checkpoint:



CS 525



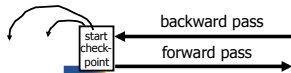
Notes 13 - Failure and Recovery

84

IIT College of Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

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



Real world actions

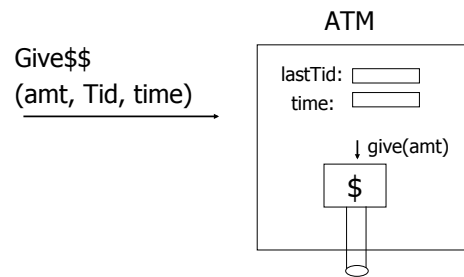
E.g., dispense cash at ATM

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

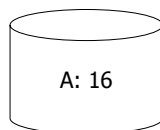
↓
\$

Solution

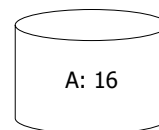
- (1) execute real-world actions after commit
- (2) try to make idempotent



Media failure (loss of non-volatile storage)



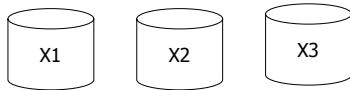
Media failure (loss of non-volatile storage)



Solution: Make copies of data!

Example 1 Triple modular redundancy

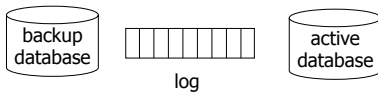
- Keep 3 copies on separate disks
- Output(X) --> three outputs
- Input(X) --> three inputs + vote



Example #2 Redundant writes, Single reads

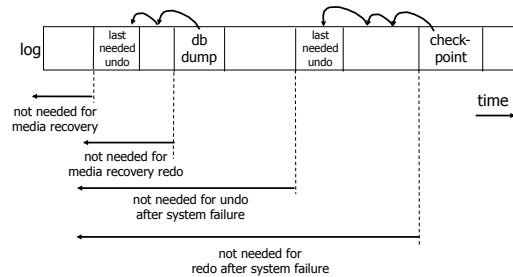
- Keep N copies on separate disks
 - Output(X) --> N outputs
 - Input(X) --> Input one copy
- done { - if ok,
- else try another one
- ⇒ Assumes bad data can be detected

Example #3: DB Dump + Log



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

When can log be discarded?



Practical Recovery with ARIES

- **ARIES**
 - Algorithms for Recovery and Isolation Exploiting Semantics
- Implemented in, e.g.,
 - DB2
 - MSSQL

Underlying Ideas

- Keep track of state of pages by relating them to entries in the log
- **WAL**
- Recovery in **three phases**
 - Analysis, Redo, Undo
- Log entries to track state of Undo for repeated failures
- **Redo**: page-oriented -> efficient
- **Undo**: logical -> permits higher level of concurrency

Log Entry Structure

- **LSN**
 - Log sequence number
 - Order of entries in the log
 - Usually **log file id** and **offset** for direct access

CS 525



Notes 13 - Failure and Recovery

97

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

- **LSN**
- **Entry type**
 - Update, compensation, commit, ...
- **TID**
 - Transaction identifier
- **PrevLSN**
 - LSN of previous log record for same transaction
- **UndoNxtLSN**
 - Next undo operation for CLR (later!)
- **Undo/Redo data**
 - Data needed to undo/redo the update

CS 525



Notes 13 - Failure and Recovery

98

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Page Header Additions

- **PageLSN**
 - **LSN** of the last update that modified the page
 - Used to know which changes have been applied to a page

CS 525



Notes 13 - Failure and Recovery

99

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Forward Processing

- Normal operations when no ROLLBACK is required
 - WAL: write redo/undo log record for each action of a transaction
- Buffer manager has to ensure that
 - changes to pages are not persisted before the corresponding log record has been persisted
 - Transactions are not considered committed before all their log records have been flushed

CS 525



Notes 13 - Failure and Recovery

100

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Dirty Page Table

- **PageLSN**
 - Entries **<PageID,ReclSN>**
 - Whenever a page is first fixed in the buffer pool with intention to modify
 - Insert **<PageId,ReclSN>** with **ReclSN** being the current end of the log
 - Flushing a page removes it from the Dirty page table

CS 525



Notes 13 - Failure and Recovery

101

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Dirty Page Table

- Used for checkpointing
- Used for recovery to figure out what to redo

CS 525



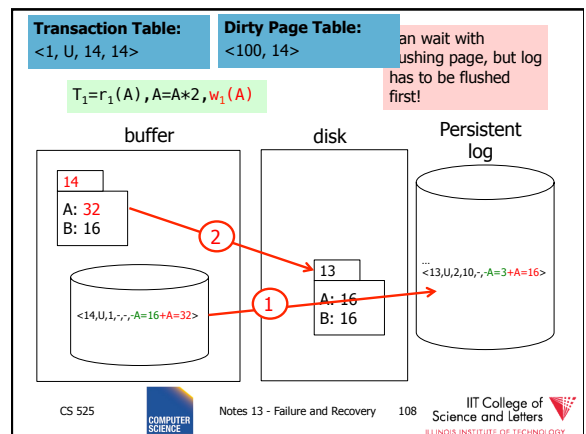
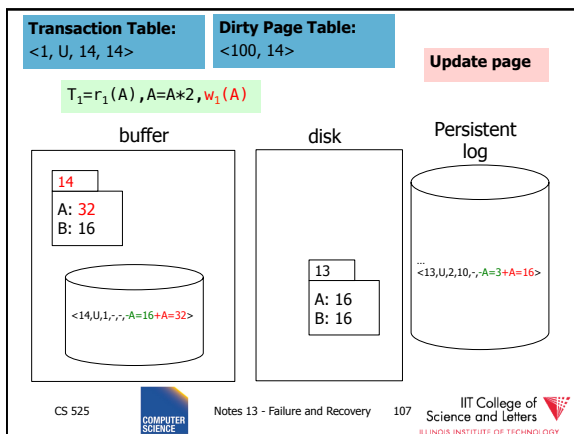
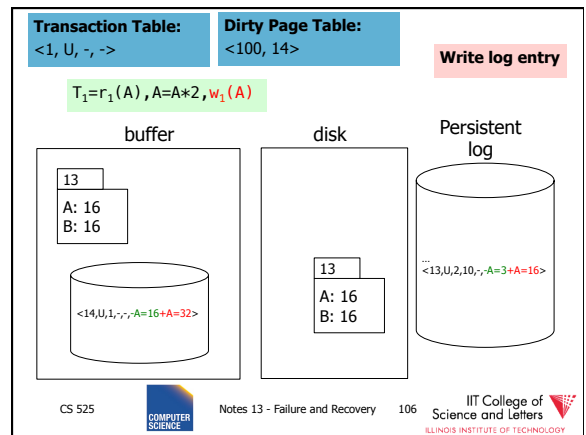
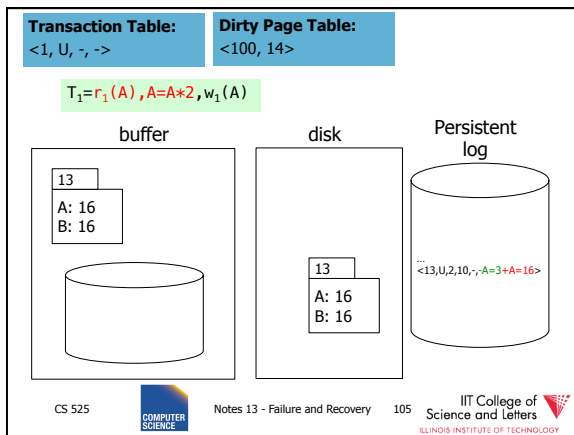
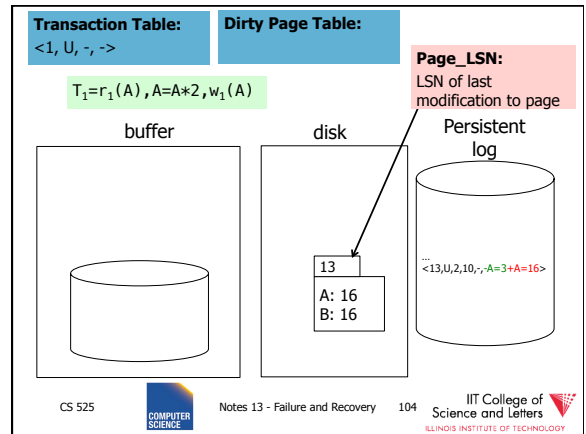
Notes 13 - Failure and Recovery

102

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Transaction Table

- TransID
 - Identifier of the transaction
- State
 - Commit state
- LastLSN
 - LSN of the last update of the transaction
- UndoNxtLSN
 - If last log entry is a CLR then UndoNxtLSN from that record
 - Otherwise = LastLSN



Undo during forward processing

- Transaction was rolled back
 - User aborted, aborted because of error, ...
- Need to undo operations of transaction
- During Undo
 - Write log entries for every undo
 - Compensation Log Records (CLR)**
 - Used to avoid repeated undo when failures occur

CS 525



Notes 13 - Failure and Recovery

109

Undo during forward processing

- Starting with the LastLSN of transaction from transaction table
 - Traverse log entries of transaction last to first using PrevLSN pointers
 - For each log entry use undo information to undo action
 - <LSN, Type, TID, PrevLSN, -, Undo/Redo data>
 - Before modifying data write an CLR that stores redo-information for the undo operation
 - UndoNextLSN** = PrevLSN of log entry we are undoing
 - Redo data** = How to redo the undo

CS 525



Notes 13 - Failure and Recovery

110

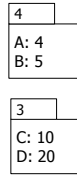
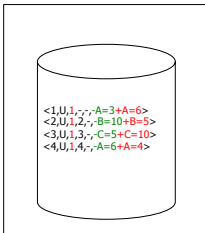
Transaction Table:

<1, U, 4, 4>

Undo T_1

$T_1 = w_1(A), w_1(B), w_1(C), w_1(A), a_1$

buffer



CS 525



Notes 13 - Failure and Recovery

111

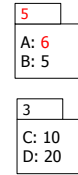
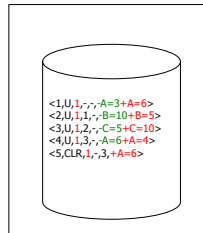
Transaction Table:

<1, U, 5, 3>

Undo T_1

$T_1 = w_1(A), w_1(B), w_1(C), w_1(A), a_1$

buffer



CS 525



Notes 13 - Failure and Recovery

112

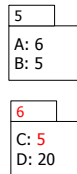
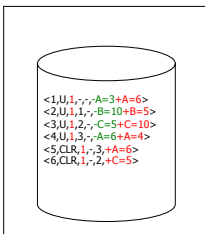
Transaction Table:

<1, U, 6, 2>

Undo T_1

$T_1 = w_1(A), w_1(B), w_1(C), w_1(A), a_1$

buffer



CS 525



Notes 13 - Failure and Recovery

113

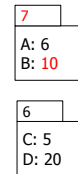
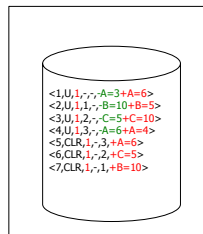
Transaction Table:

<1, U, 7, 1>

Undo T_1

$T_1 = w_1(A), w_1(B), w_1(C), w_1(A), a_1$

buffer



CS 525



Notes 13 - Failure and Recovery

114

Transaction Table:
 <1, U, 8, ->

Undo T₁

T₁ = w₁(A), w₁(B), w₁(C), w₁(A), a₁

buffer

The diagram shows a buffer containing the following log entries:
 <1,U,1,-,A=3+A=6>
 <2,U,1,1,-,B=10+B=5>
 <3,U,1,2,-,C=5+C=10>
 <4,U,1,3,-,A=6+A=4>
 <5,CLR,1,-,3,+A=6>
 <6,CLR,1,-,2,+C=5>
 <7,CLR,1,-,1,+B=10>
 <8,CLR,1,-,-,A=3>

Undo information for T₁ is shown as:
 8
 A: 3
 B: 10
 6
 C: 5
 D: 20

CS 525 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Fuzzy Checkpointing in ARIES

- **Begin of checkpoint**
 - Write **begin_cp** log entry
 - Write **end_cp** log entry with
 - Dirty page table
 - Transaction table
- **Master Record**
 - LSN of begin_cp log entry of last complete checkpoint

CS 525 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Restart Recovery

1. Analysis Phase
2. Redo Phase
3. Undo Phase

CS 525 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Analysis Phase

- 1) Determine LSN of last checkpoint using Master Record
- 2) Log Dirty Page Table and Transaction Table from checkpoint record
- 3) **RedoLSN** = min(ReclSN) from Dirty Page Table or checkpoint LSN if no dirty page

CS 525 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Analysis Phase

- 4) Scan log forward starting from RedoLSN
 - Update log entry from transaction
 - If necessary: Add to Page to Dirty Page Table
 - Add Transaction to Transaction Table or update LastLSN
 - Transaction end entry
 - Remove transaction from Transaction Table

CS 525 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Analysis Phase

- **Result**
 - Transaction Table
 - Transactions to be later undone
 - RedoLSN
 - Log entry to start Redo Phase
 - Dirty Page Table
 - Pages that may not have been written back to disk

CS 525 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Redo Phase

- Start at RedoLSN scan log forward
- Unconditional Redo
 - Even redo actions of transactions that will be undone later
- One redo once
 - Only redo operations that have not been reflected on disk (PageLSN)

CS 525



Notes 13 - Failure and Recovery

121

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Redo Phase

- For each update log entry
 - If affected page is not in Dirty Page Table or $ReLSN > LSN$
 - skip log entry
 - Fix page in buffer
 - If $PageLSN \geq LSN$ then operation already reflected on disk
 - Skip log entry
 - Otherwise apply update

CS 525



Notes 13 - Failure and Recovery

122

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Redo Phase

- Result
 - State of DB before Failure

CS 525



Notes 13 - Failure and Recovery

123

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Undo Phase

- Scan log backwards from end using Transaction Table
 - Repeatedly take log entry with max LSN from all the current action to be undone for each transaction
 - Write CLR
 - Update Transaction Table

CS 525



Notes 13 - Failure and Recovery

124

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Undo Phase

- All unfinished transactions have been rolled back

CS 525



Notes 13 - Failure and Recovery

125

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Idempotence?

- Redo
 - We are not logging during Redo so repeated Redo will result in the same state
- Undo
 - If we see CLR we do not undo this action again

CS 525



Notes 13 - Failure and Recovery

126

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY


$T_1 = w_1(A), w_1(B), w_1(C), w_1(A), c_1$
 $T_2 = w_1(X), r(A), w(A)$

Log

```

<1,begin(T1),->
<2,begin(T2),->
<3,write(A,T1),1>
<4,write(X,T2),2>
<5,write(B,T1),3>
<6,write(C,T1),5>
<7,write(A,T1),6>
<8,commit(T1),7>
<9,write(A,T2),4>

```

CS 525  Notes 13 - Failure and Recovery 127 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

$T_1 = w_1(A), w_1(B), w_1(C), w_1(A), c_1$
 $T_2 = w_1(X), r(A), w(A)$

Analysis Phase:


- start at log entry 1
- add T_1 to transaction table (rec. 1)
- add T_2 to transaction table (rec. 2)
- add A to page table (RecLSN 3)
- add X to page table (RecLSN 4)
- add B to page table (RecLSN 5)
- add C to page table (RecLSN 6)
- remove T_1 from Transaction Table (rec. 8)

Log

```

<1,begin(T1),->
<2,begin(T2),->
<3,write(A,T1),1>
<4,write(X,T2),2>
<5,write(B,T1),3>
<6,write(C,T1),5>
<7,write(A,T1),6>
<8,commit(T1),7>
<9,write(A,T2),4>

```

CS 525  Notes 13 - Failure and Recovery 128 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

$T_1 = w_1(A), w_1(B), w_1(C), w_1(A), c_1$
 $T_2 = w_1(X), r(A), w(A)$

Analysis Phase Result:


- Transaction Table: $\langle T_2, 9 \rangle$
- Dirty Page Table: $\langle A, 3 \rangle, \langle B, 5 \rangle, \langle C, 6 \rangle, \langle X, 4 \rangle$
- RedoLSN = $\min(3,5,6,4) = 3$

Log

```

<1,begin(T1),->
<2,begin(T2),->
<3,write(A,T1),1>
<4,write(X,T2),2>
<5,write(B,T1),3>
<6,write(C,T1),5>
<7,write(A,T1),6>
<8,commit(T1),7>
<9,write(A,T2),4>

```

CS 525  Notes 13 - Failure and Recovery 129 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

$T_1 = w_1(A), w_1(B), w_1(C), w_1(A), c_1$
 $T_2 = w_1(X), r(A), w(A)$

Redo Phase (RedoLSN 3):


- Read A if PageLSN < 3 apply write
- Read X if PageLSN < 4 apply write
- Read B if PageLSN < 5 apply write
- Read C if PageLSN < 6 apply write
- Read A if PageLSN < 7 apply write
- Read A if PageLSN < 9 apply write

Log

```

<1,begin(T1),->
<2,begin(T2),->
<3,write(A,T1),1>
<4,write(X,T2),2>
<5,write(B,T1),3>
<6,write(C,T1),5>
<7,write(A,T1),6>
<8,commit(T1),7>
<9,write(A,T2),4>

```

CS 525  Notes 13 - Failure and Recovery 130 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

$T_1 = w_1(A), w_1(B), w_1(C), w_1(A), c_1$
 $T_2 = w_1(X), r(A), w(A)$

Undo Phase (T_2):


- Undo entry 9
 - write CLR with UndoNextLSN = 4
 - modify page A
- Undo entry 4
 - write CLR with UndoNextLSN = 2
 - modify page X
- Done

Log

```


<1,begin(T1),->
<2,begin(T2),->
<3,write(A,T1),1>
<4,write(X,T2),2>
<5,write(B,T1),3>
<6,write(C,T1),5>
<7,write(A,T1),6>
<8,commit(T1),7>
<9,write(A,T2),4>
<10,CLR(A,T2),4>
<11,CLR(X,T2),->

```

CS 525  Notes 13 - Failure and Recovery 131 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

ARIES take away messages

- Provide good performance by
 - Not requiring complete checkpoints
 - Linking of log records
 - Not restricting buffer operations (no-force/steal is ok)
- Logical Undo and Physical (Physiological) Redo
- Idempotent Redo and Undo
 - Avoid undoing the same

CS 525  Notes 13 - Failure and Recovery 132 IIT College of Science and Letters ILLINOIS INSTITUTE OF TECHNOLOGY

Media Recovery

- What if disks where log or DB is stored fails
 - ->keep backups of log + DB state

CS 525



Notes 13 - Failure and Recovery 133

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Log Backup

- Split log into several files
- Is append only, backup of old files cannot interfere with current log operations

CS 525



Notes 13 - Failure and Recovery 134

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Backup DB state

- Copy current DB state directly from disk
- May be inconsistent
- ->Use log to know which pages are up-to-date and redo operations not yet reflected

CS 525



Notes 13 - Failure and Recovery 135

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY

Summary

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

CS 525



Notes 13 - Failure and Recovery 136

IIT College of
Science and Letters
ILLINOIS INSTITUTE OF TECHNOLOGY