

Lecture 15 – ARIES Recovery



Stable Diffusion: a ram in cyberpunk style in space

Nov 2, 2022
Lab 4 Due 11/9

Recovery Recap

- What happens during crash:
 - Memory is reset
 - State on disk persists
- After a crash, recovery ensures:
 - **Atomicity**: partially finished xactions are rolled back
 - **Durability**: committed xactions are on stable storage (disk)
- Brings database into a transaction consistent state, where committed transactions are fully reflected, and uncommitted transactions are completely undone

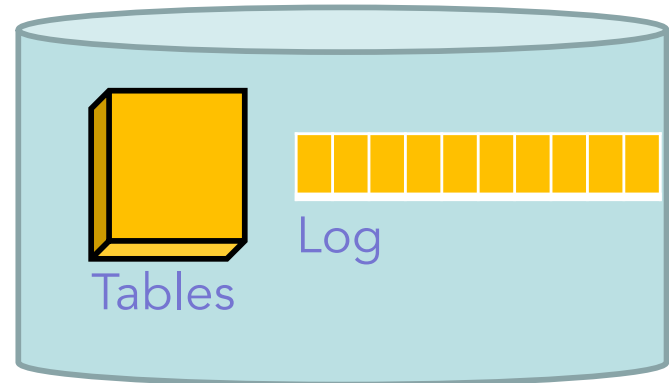
Database State

Log records start and end of transactions, and contents of writes done to tables so we can solve both problems



Memory

After crash, memory is gone!



Disk

Problem 1: Some transactions may have written their uncommitted state to tables – need to **UNDO**

Problem 2: Some transactions may not have flushed all of their state to tables prior to commit – need to **REDO**

Types of Log Records

- **Start (SOT)** Log Sequence Number (LSN), Transaction ID (TID)
 - LSN is a monotonically increasing log record number
- **End (EOT)** LSN, TID, outcome (commit or abort)
- **UNDO** LSN, TID, before image
- **REDO** LSN, TID, after image

For ARIES:

- **CHECKPOINT** LSN, TID, state to limit how much is logged
- **CLR** LSN, TID, allows us to restart recovery

Write Ahead Logging

Write what we plan to do, before we do it

Recovery with NO FORCE / STEAL

- After crash, we must:
 - REDO "winner" transactions that had committed
 - UNDO "loser" transactions that had not committed
- Winners are transactions with SOT and COMMIT in log
- Losers are those with SOT and no EOT, or ABORT
- Need to REDO winners from start to end
- Need to UNDO losers in reverse, from end to start
- Also need to UNDO aborted transactions

ARIES

- Gold standard in logging
 - Specifies **all** the details
- NO FORCE/STEAL
- Recoverable recovery
- “Physiological” Logging
- Low-overhead Checkpoints
- Support for escrow operations
 - E.g., increment / decrements

ARIES Approach: 3 Log Passes

FW = Forward Pass; BW = Backward Pass

- **Analysis**, to see what needs to be done (FW)
- **Redo**, to ensure DB reflects updates that are in the log but not in tables (FW)
 - Including those that belong to txns that will eventually be rolled back!
 - Why? Ensures “action consistent” state -- which will allow logical undo.
 - “Repeating History”
- **Undo**, to rollback losers (BW)

Log Record Format

Every log record has an LSN associated with it.

Log Record

LSN	TID	prevLSN
Undo Image (logical)		
Redo Image (physical)		

The previous LSN written by this transaction

Update records have both UNDO and REDO information.

Disk Page

pageLSN

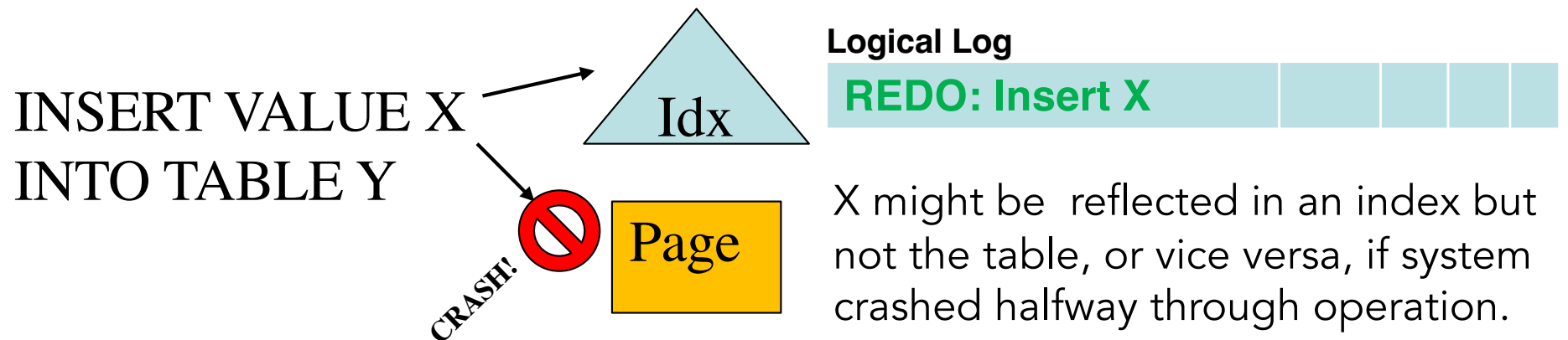
Every time a page is written, the latest LSN associated with that log record is included as the pageLSN.

“Physiological” Logging

- REDO is Physical
- UNDO is Logical

REDO must be physical

- At time of crash, database may not be in an "action consistent" state
- Some ops can encompass multiple non-atomic physical operations

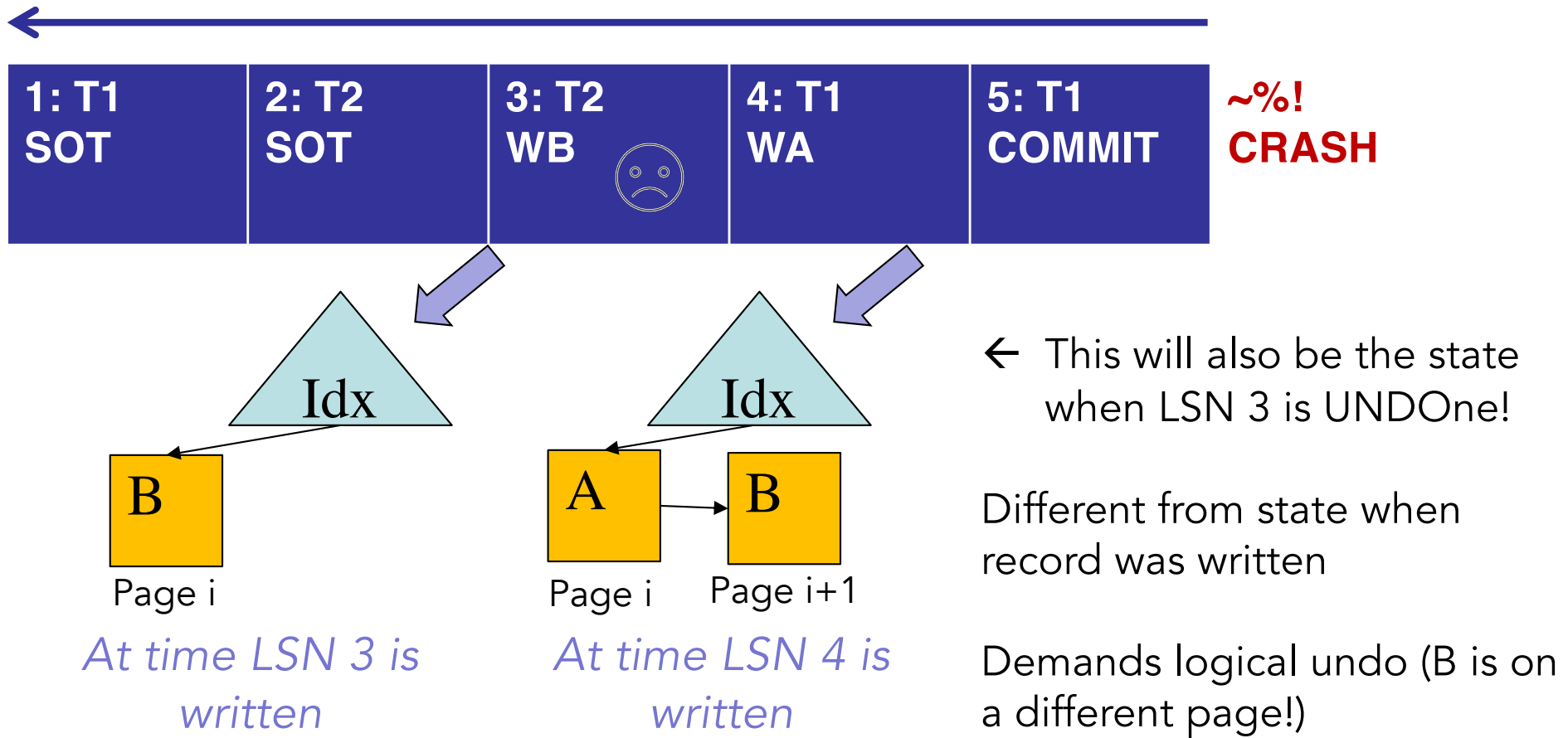


- Implies logging must be physical

UNDO must be logical

- We only UNDO some actions
- Implies state when UNDOing may not be same as when the log was written
- Physical logging (e.g., of the specific before and after images in the page) would fail

UNDO Example



Not needed in REDO, because we "repeat history" and replay everything

- Physical modifications made to the database since last time will still be correct.

ARIES Normal Operation

- Two key data structures:
 - Transaction table -- list of active transactions
 - Dirty page table -- List of pages that have been modified and not yet written to disk
- Data structures updates as system runs:
 - Pages asynchronously flushed to disk
 - Log forced before flush (but not before write)
 - Flushed are not logged
 - Log forced before COMMIT ack'd

Transaction Table

TransactionTable

lastLSN	TID
13	3

- All active transactions in table
- lastLSN: most recent log record written by that transaction

Dirty Page Table

dirtyPgTable

pgNo	recLSN
D	8
B	10
A	11
E	13

Recall, dirty pages are periodically flushed to disk by a background process.

On flush, remove from dirtyPageTable

- One entry for each page that has been modified but not flushed to disk
- recLSN: log record that **first** dirtied the page

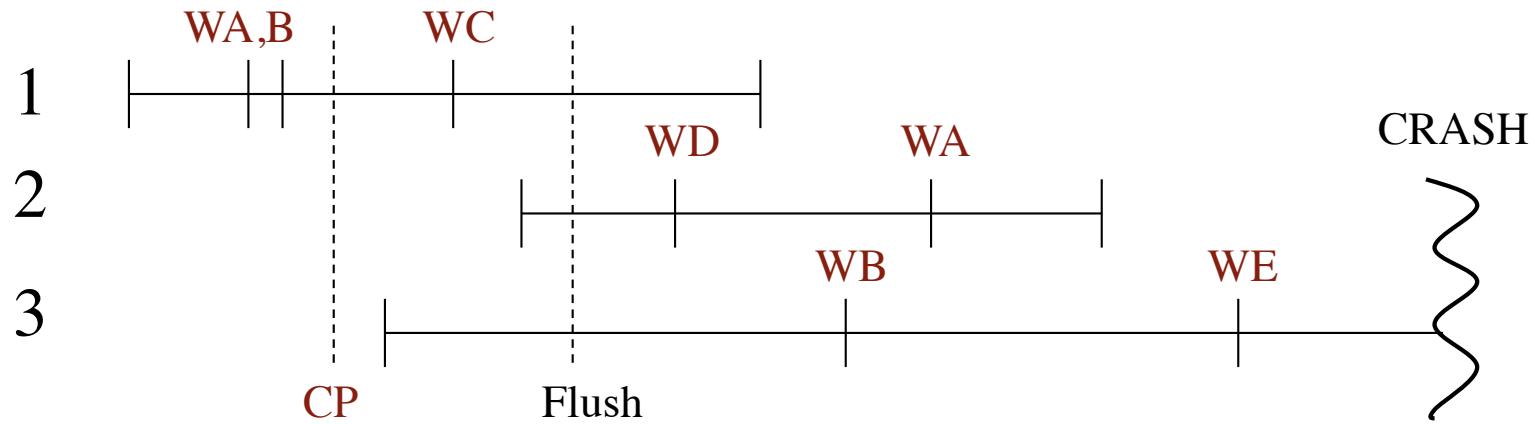
Checkpoints

- Taken periodically
- Log record that contains:
 1. the state of the dirty page table and
 2. the transaction table

Doesn't require pages to be flushed to disk during checkpoint

- Allow us to limit amount of log we have to keep and replay during crash

ARIES Example



LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	1	A
3	UP	1	2	B
4	CP			
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

Flush

ARIES Data Structures

xactionTable

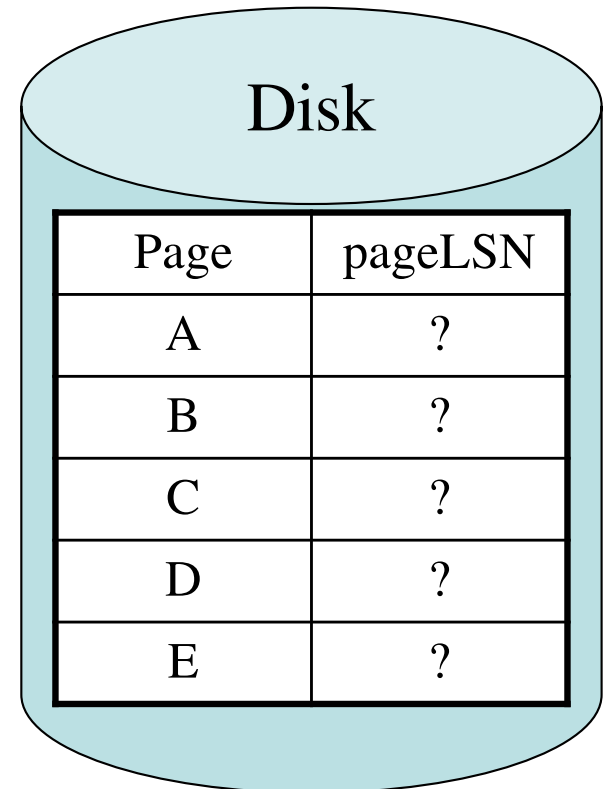
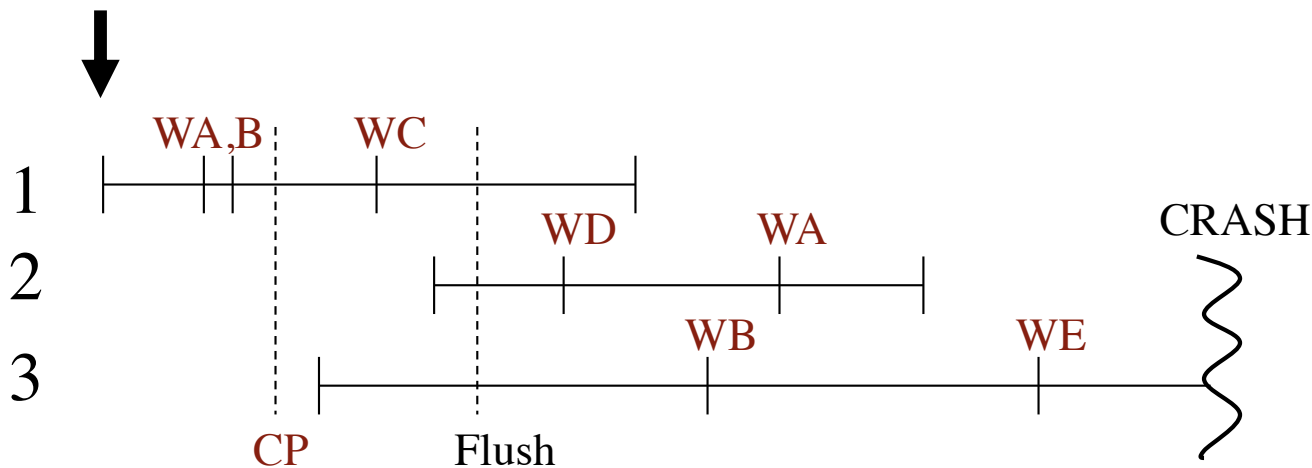
lastLSN	TID
1	1

dirtyPgTable

pgNo	recLSN

Checkpoint

xactionTable	
dirtyPgTable	



ARIES Data Structures

xactionTable

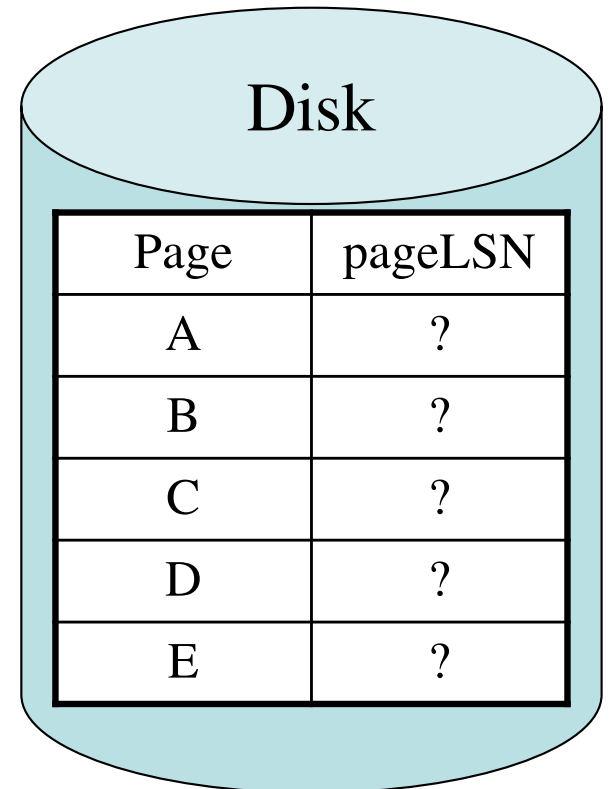
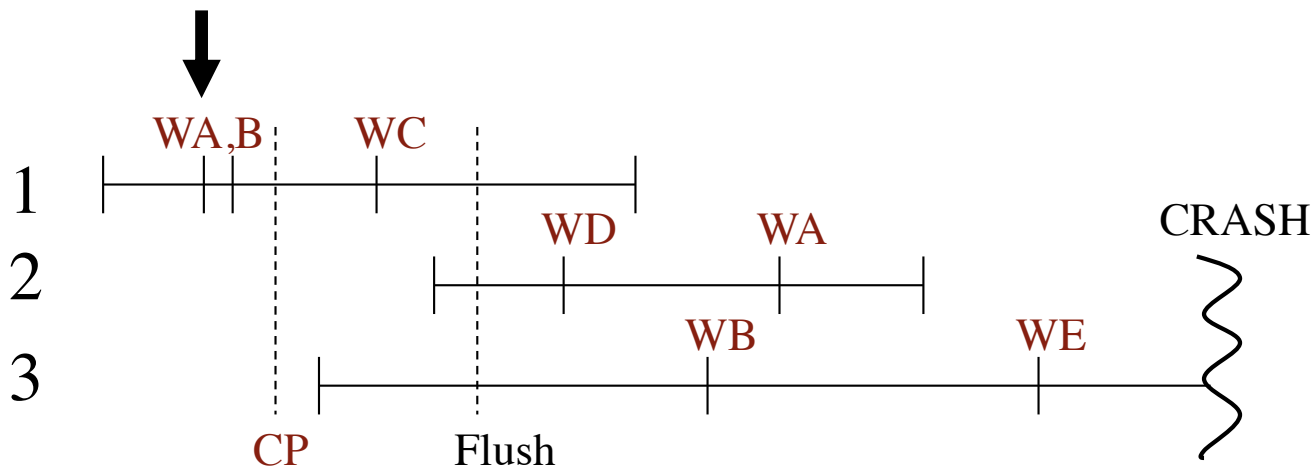
lastLSN	TID
2	1

dirtyPgTable

pgNo	recLSN
A	2

Checkpoint

xactionTable	
dirtyPgTable	



ARIES Data Structures

xactionTable

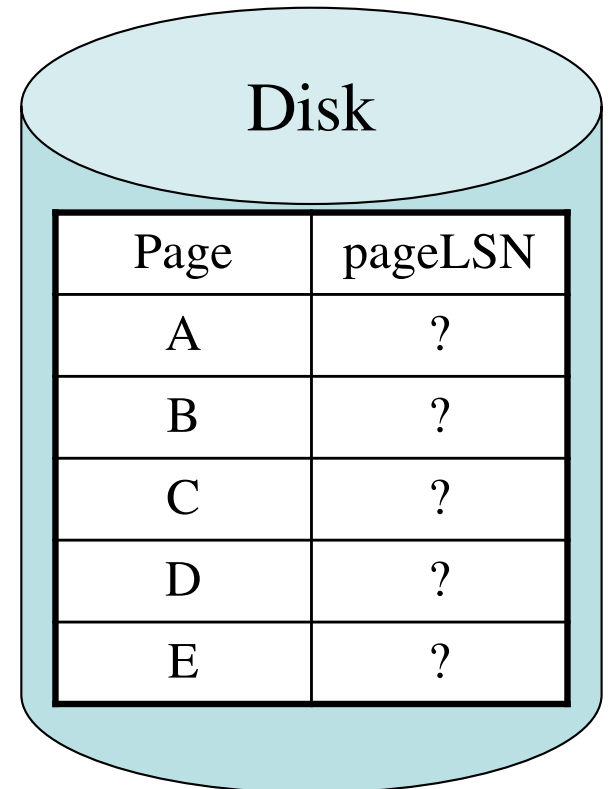
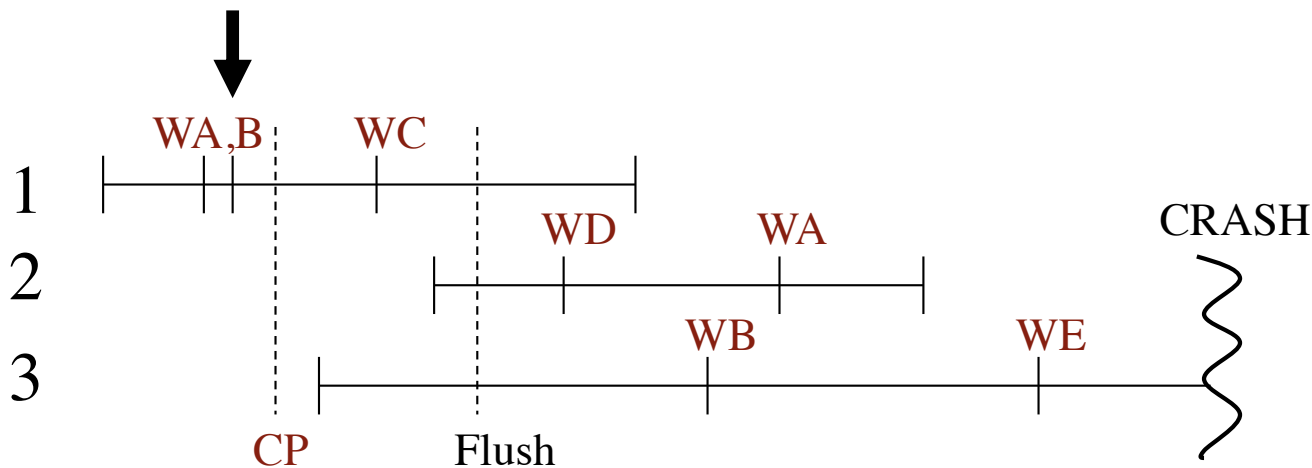
lastLSN	TID
3	1

dirtyPgTable

pgNo	recLSN
A	2
B	3

Checkpoint

xactionTable	
dirtyPgTable	



ARIES Data Structures

xactionTable

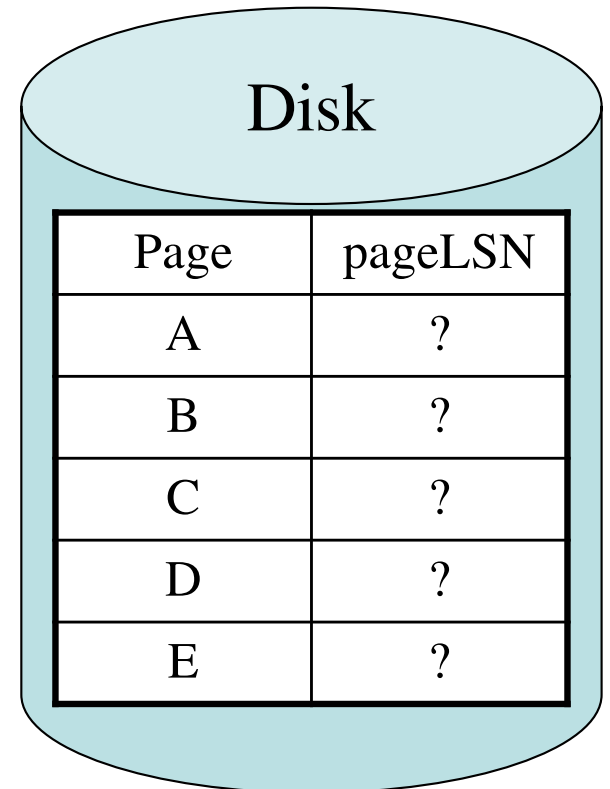
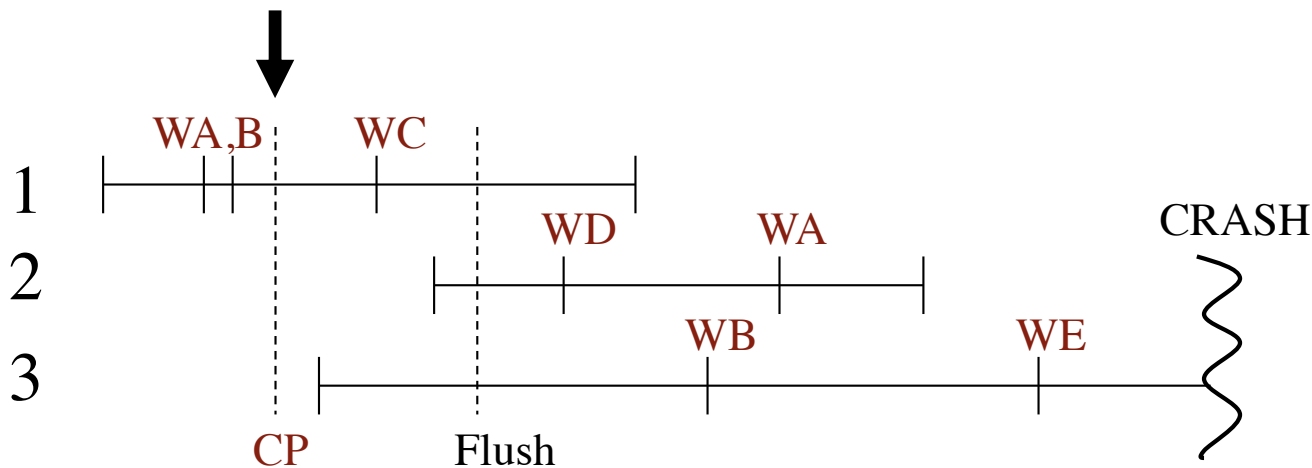
lastLSN	TID
3	1

dirtyPgTable

pgNo	recLSN
A	2
B	3

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

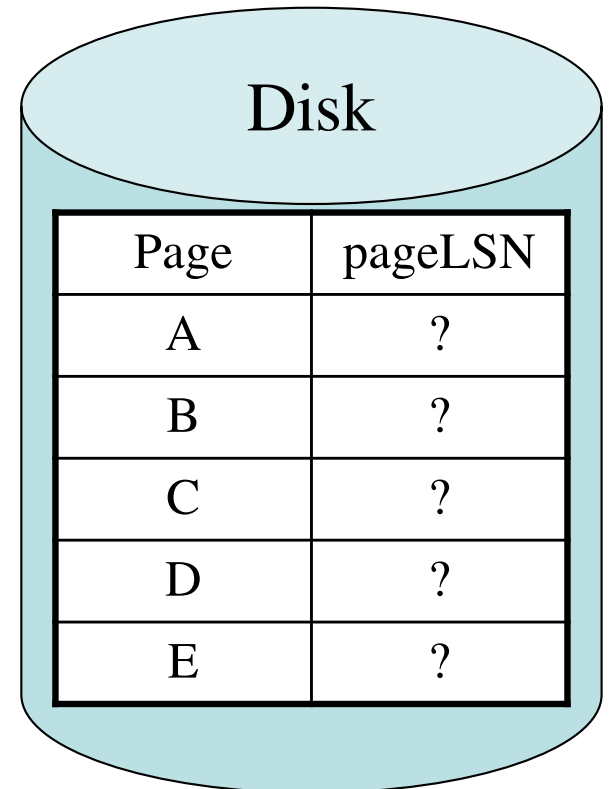
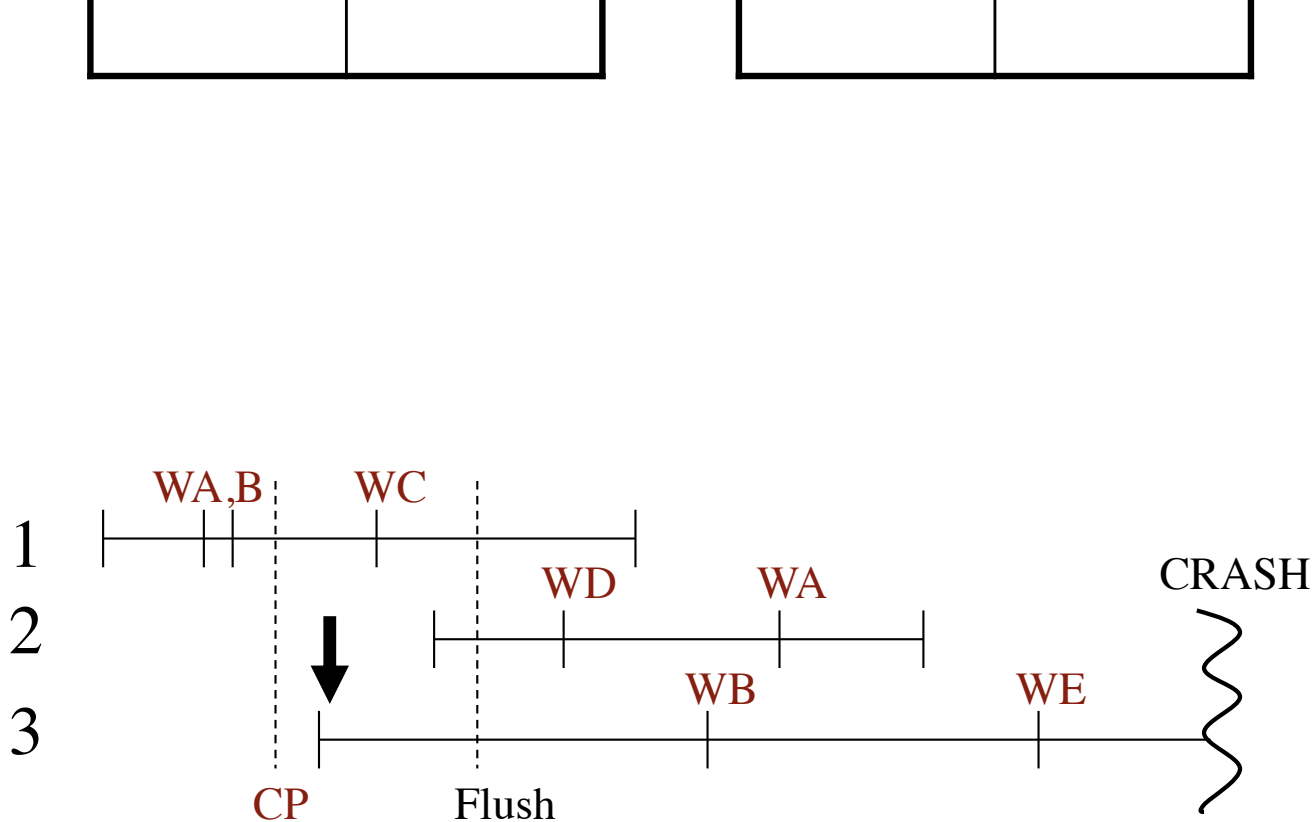
lastLSN	TID
3	1
5	3

dirtyPgTable

pgNo	recLSN
A	2
B	3

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

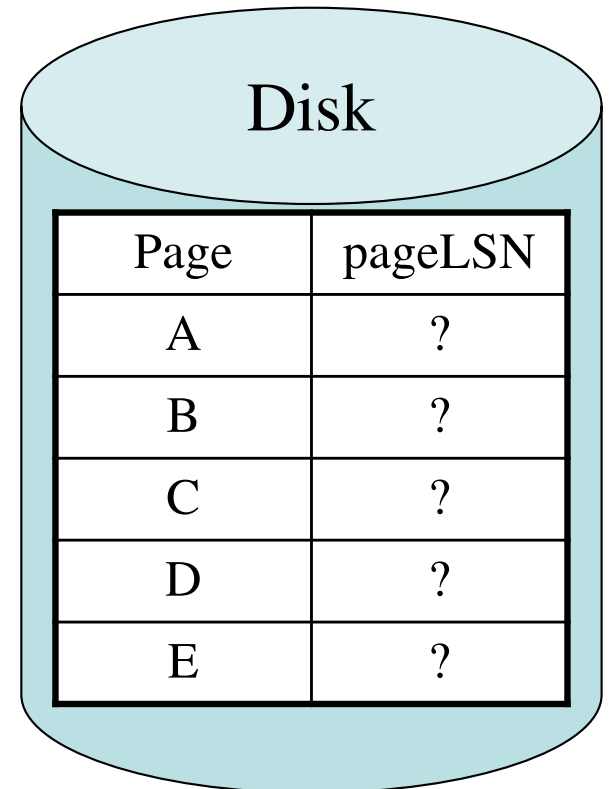
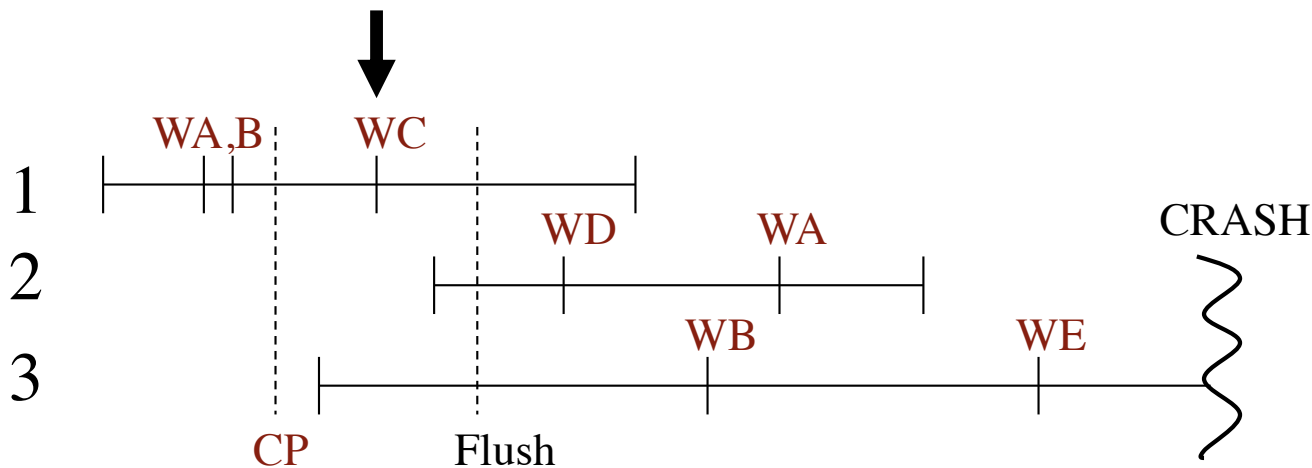
lastLSN	TID
6	1
5	3

dirtyPgTable

pgNo	recLSN
A	2
B	3
C	6

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

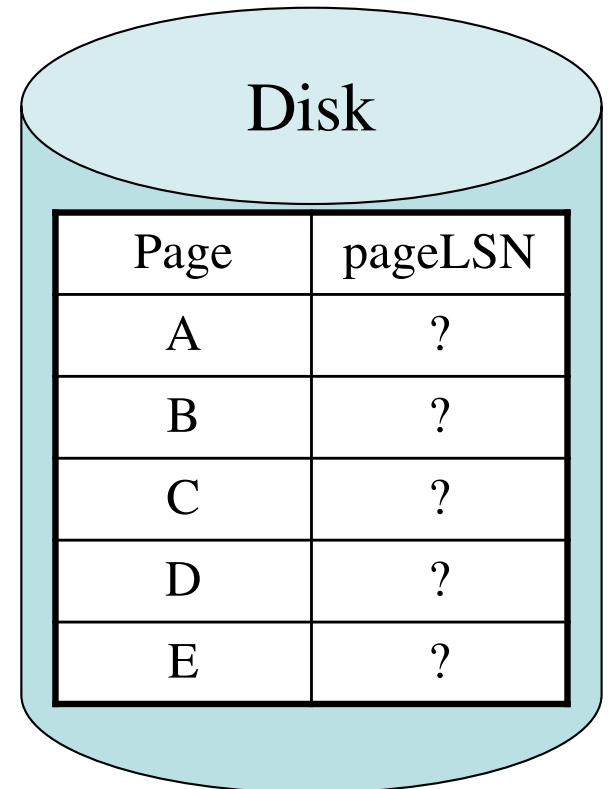
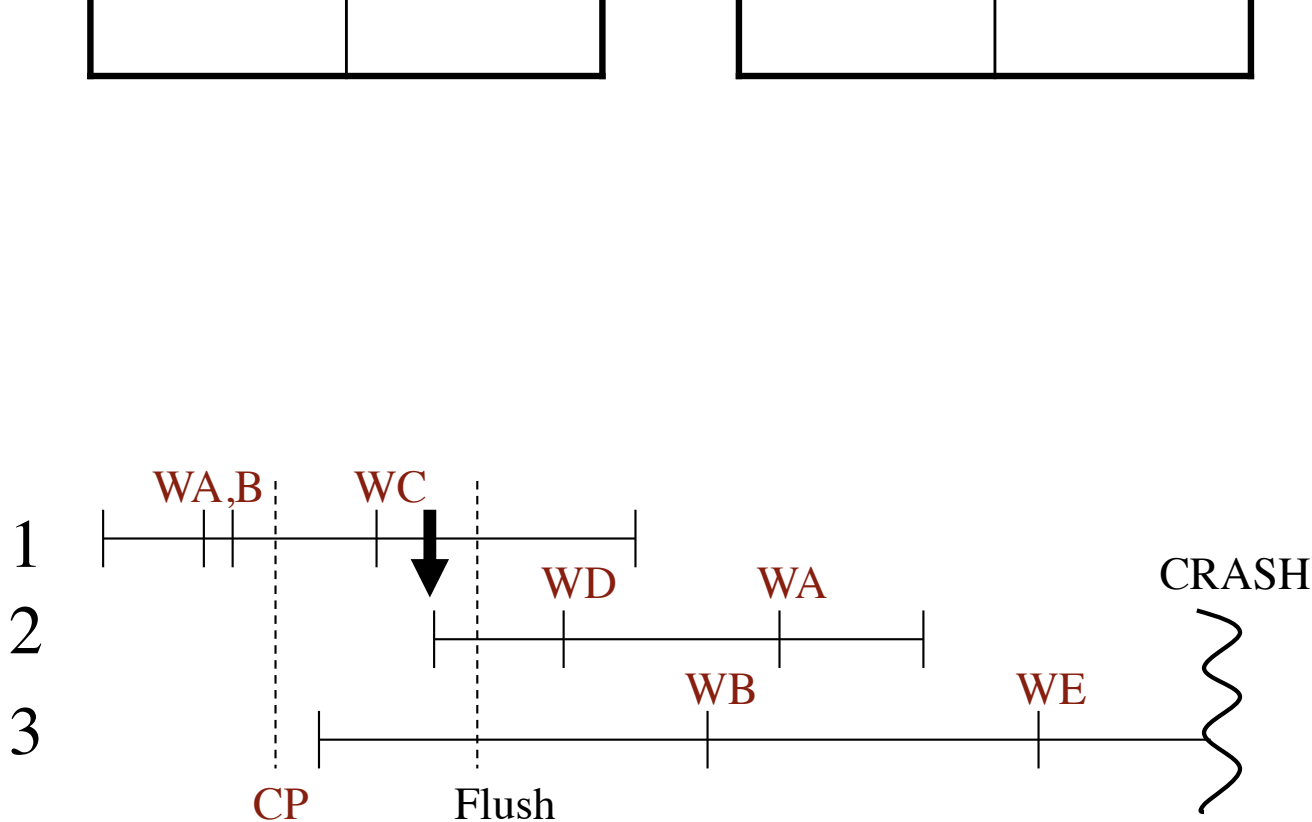
lastLSN	TID
6	1
5	3
7	2

dirtyPgTable

pgNo	recLSN
A	2
B	3
C	6

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

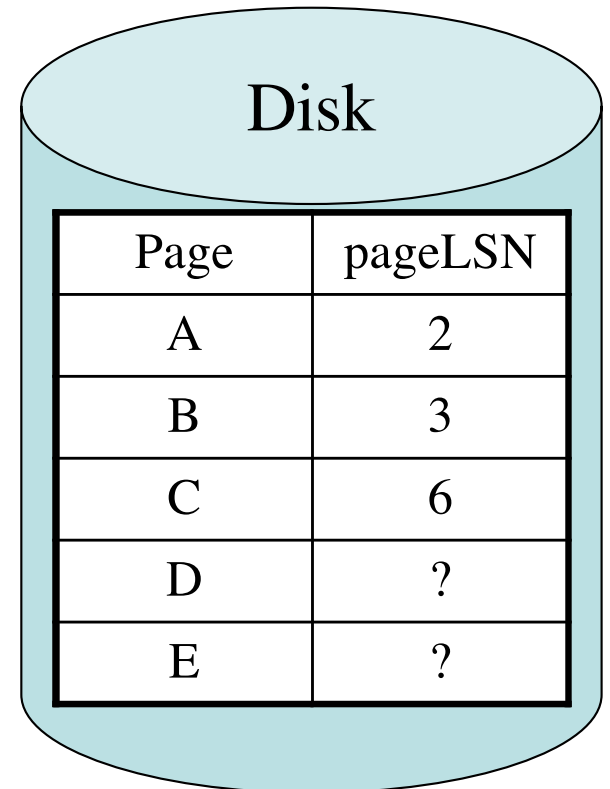
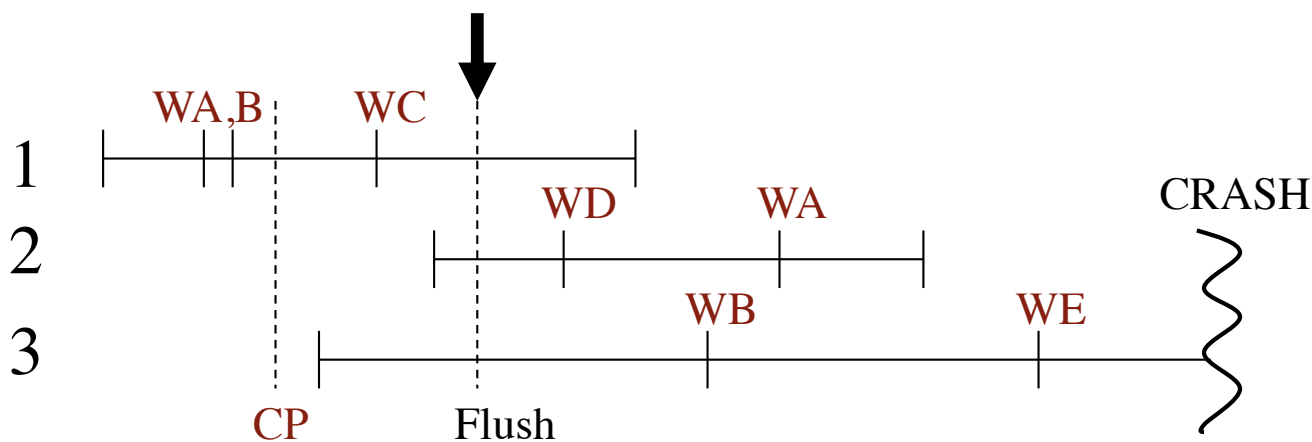
lastLSN	TID
6	1
5	3
7	2

dirtyPgTable

pgNo	recLSN

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

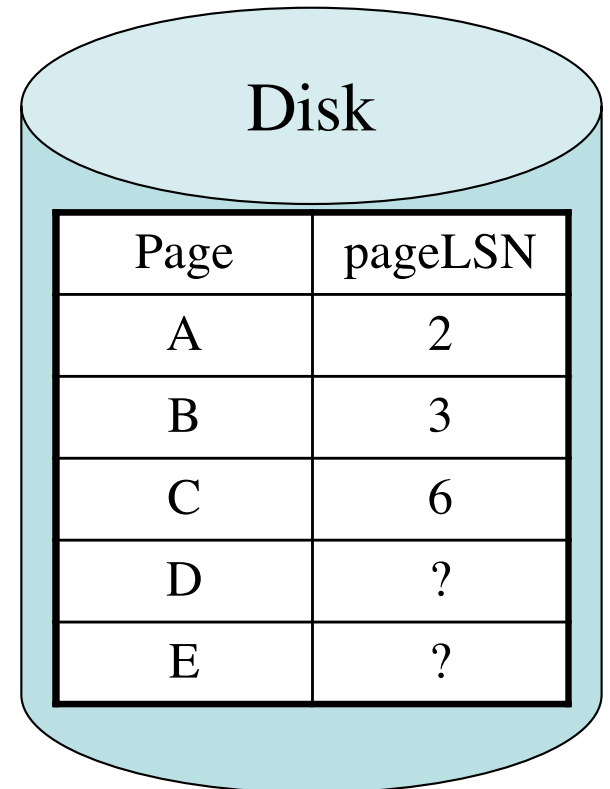
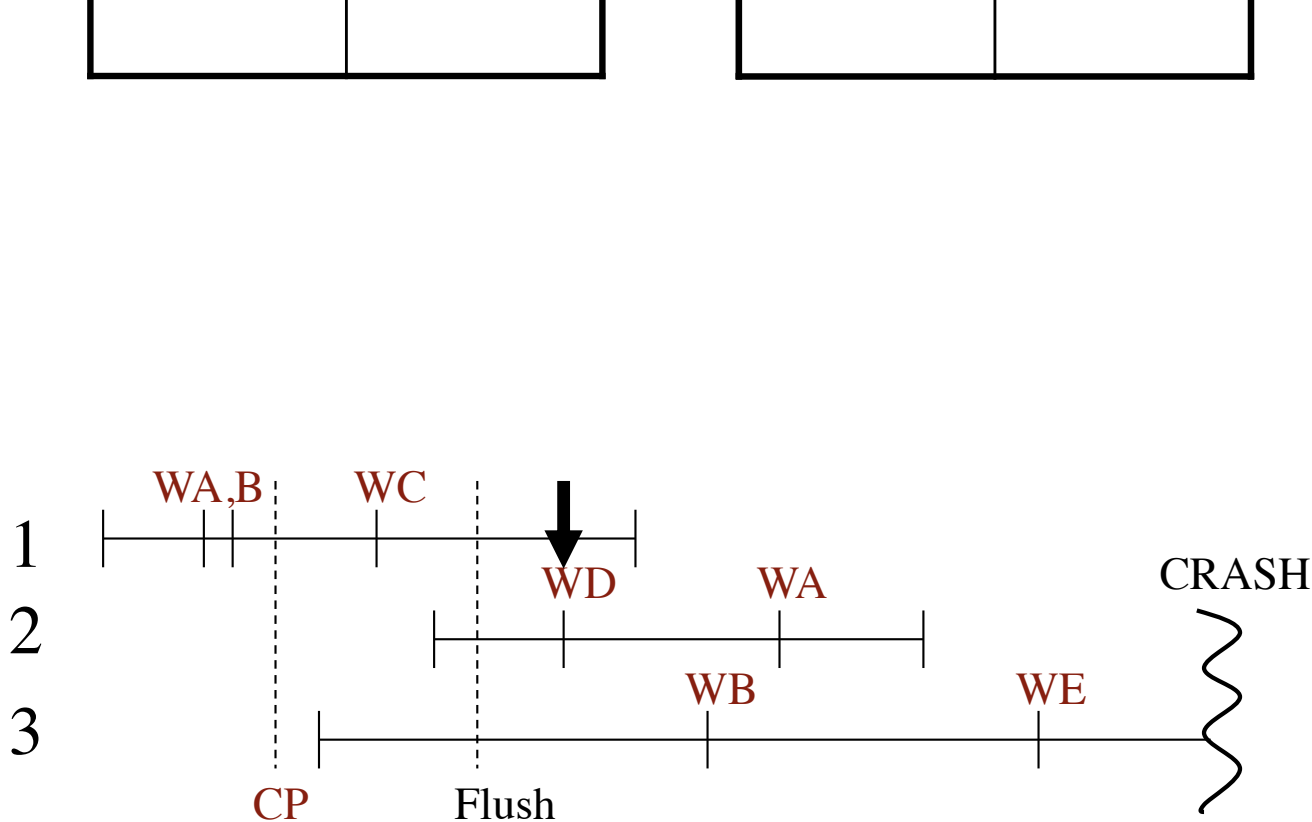
lastLSN	TID
6	1
5	3
8	2

dirtyPgTable

pgNo	recLSN
D	8

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

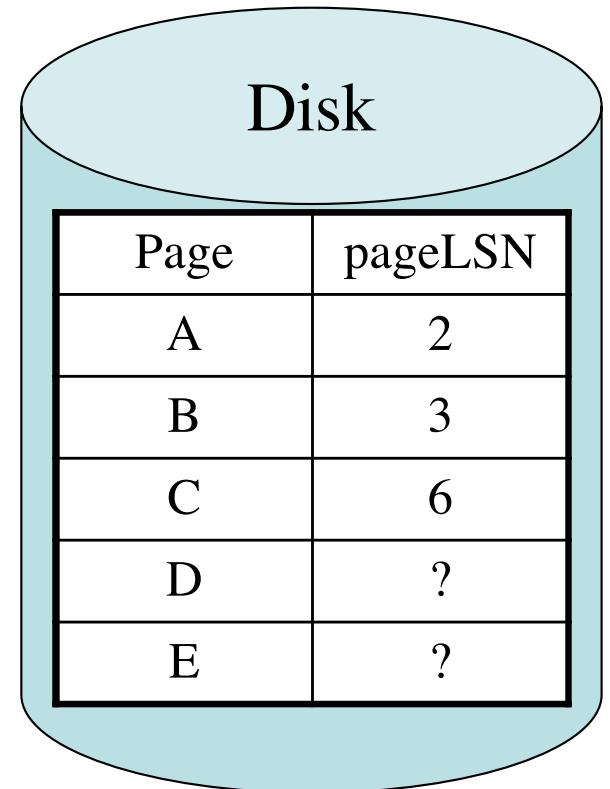
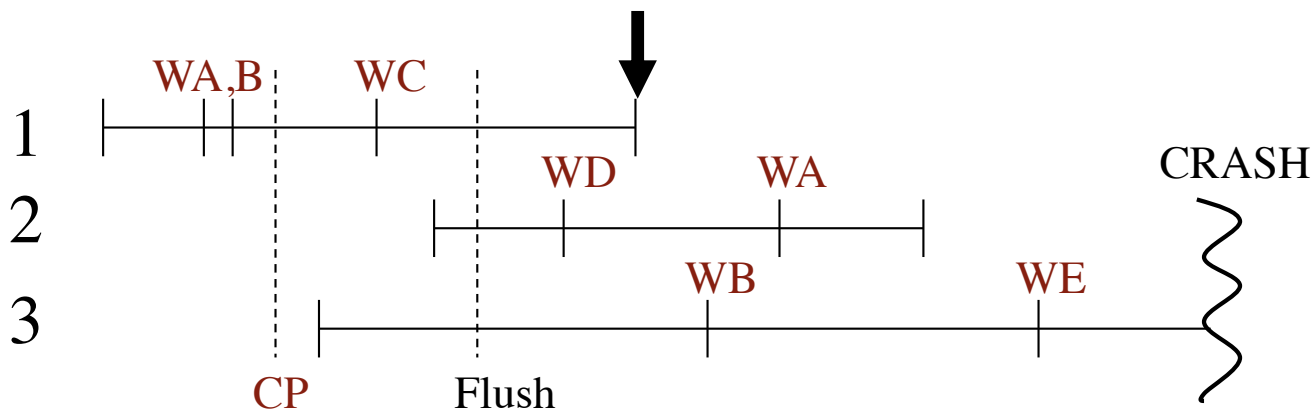
lastLSN	TID
6	1
5	3
8	2

dirtyPgTable

pgNo	recLSN
D	8

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

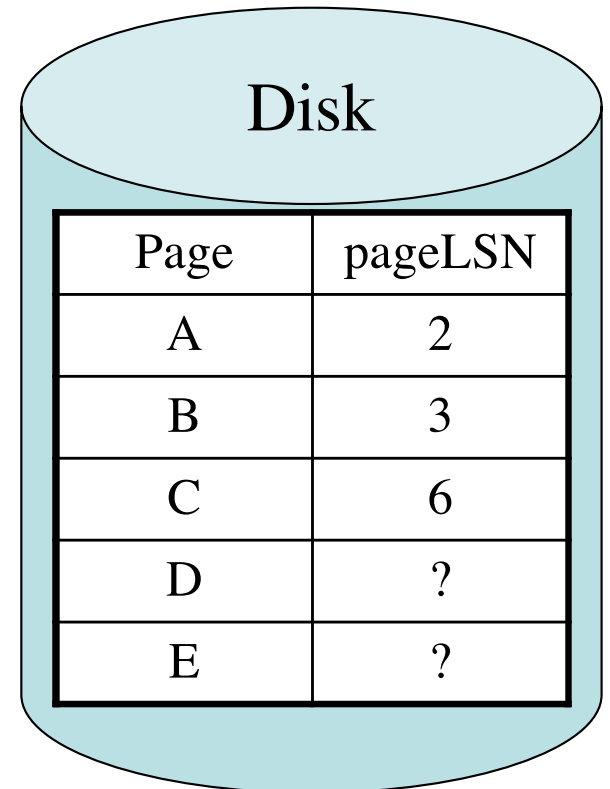
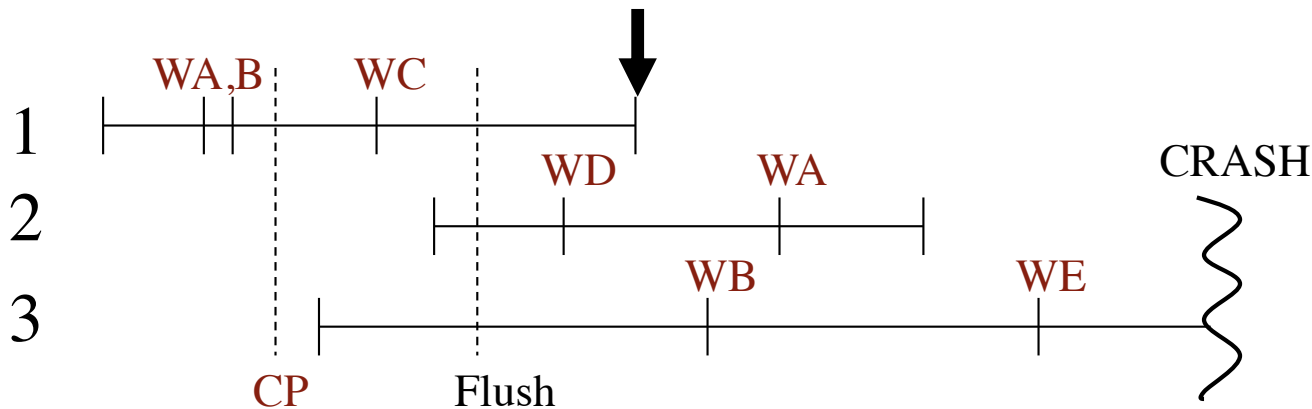
lastLSN	TID
5	3
8	2

dirtyPgTable

pgNo	recLSN
D	8

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

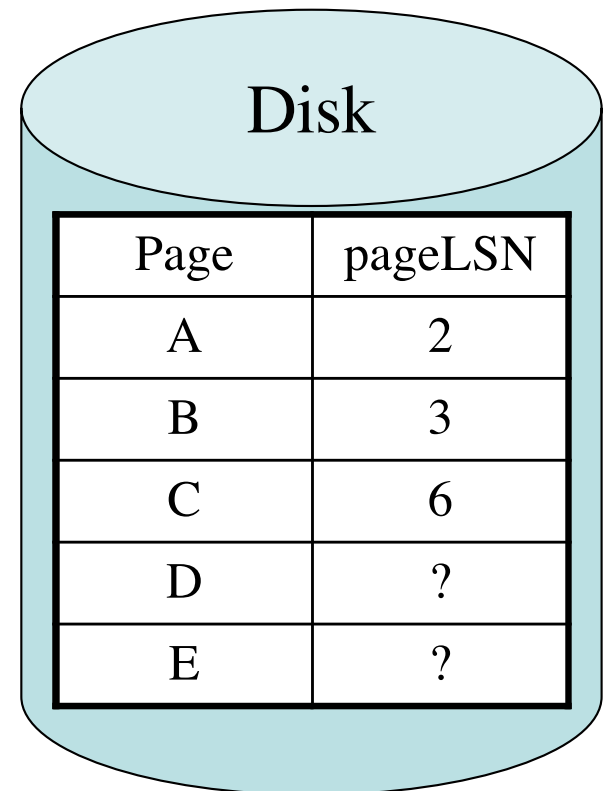
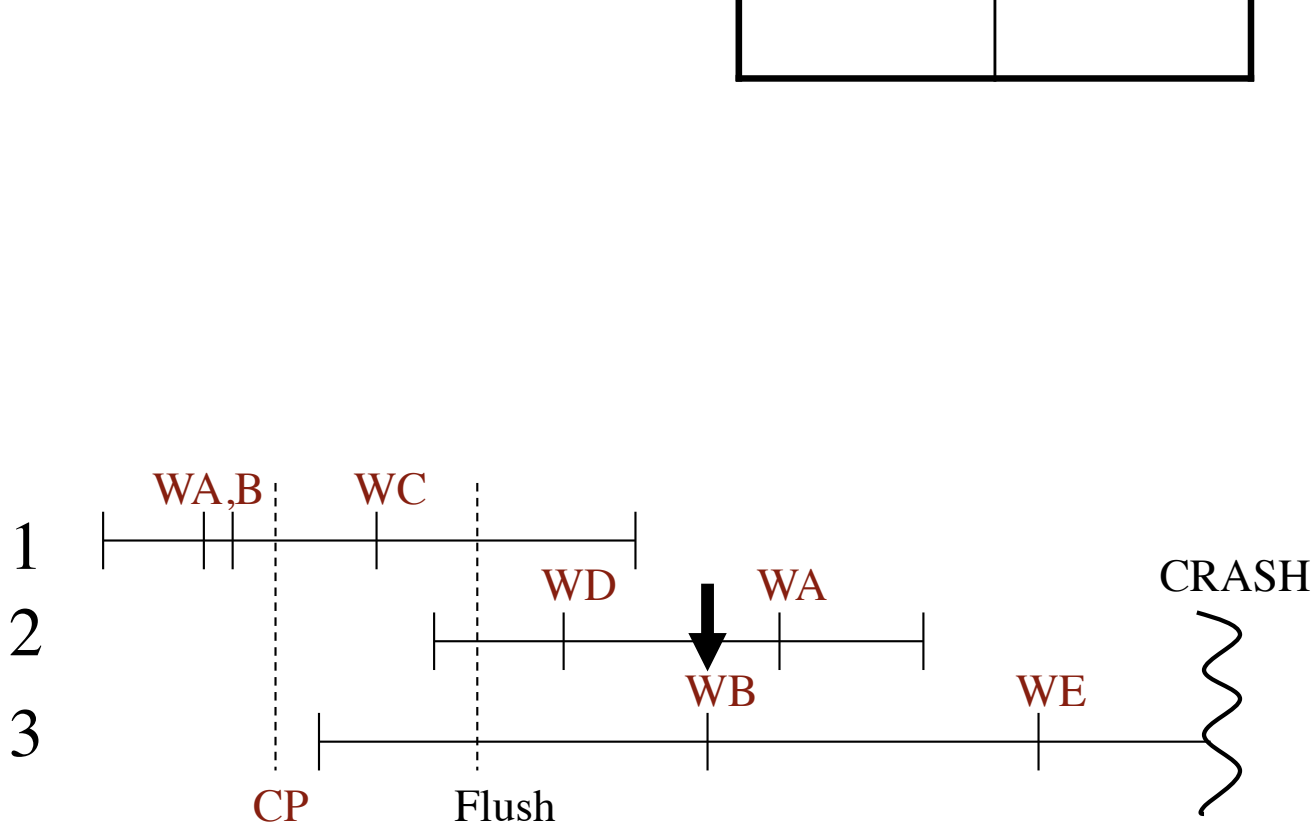
lastLSN	TID
10	3
8	2

dirtyPgTable

pgNo	recLSN
D	8
B	10

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

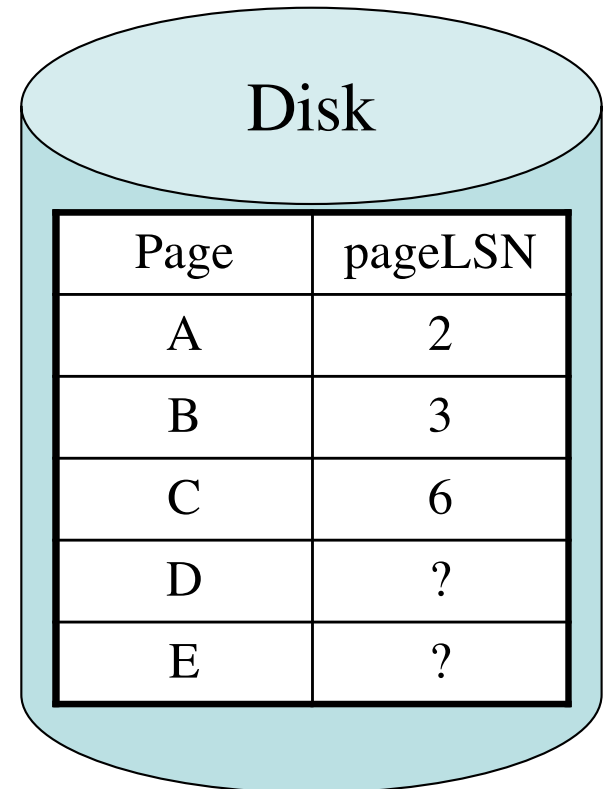
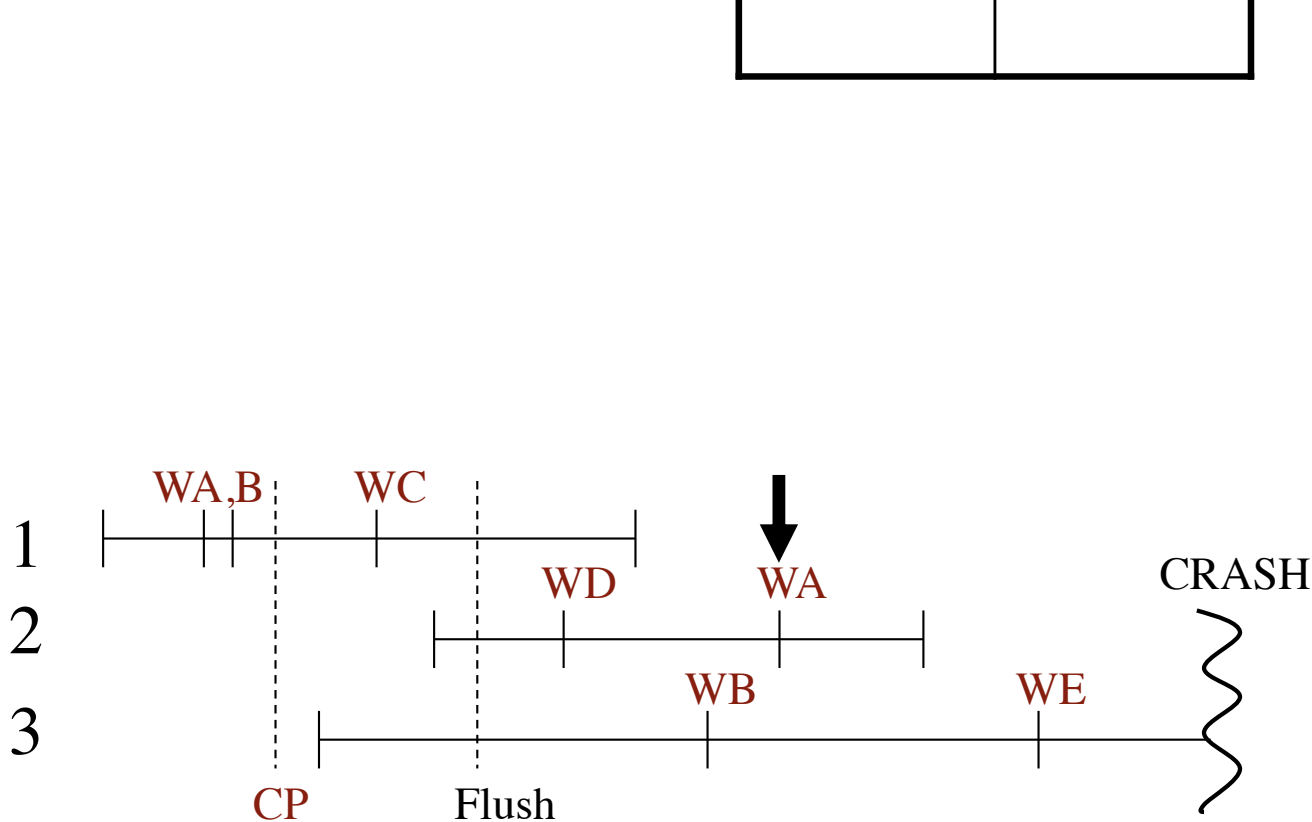
lastLSN	TID
10	3
11	2

dirtyPgTable

pgNo	recLSN
D	8
B	10
A	11

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

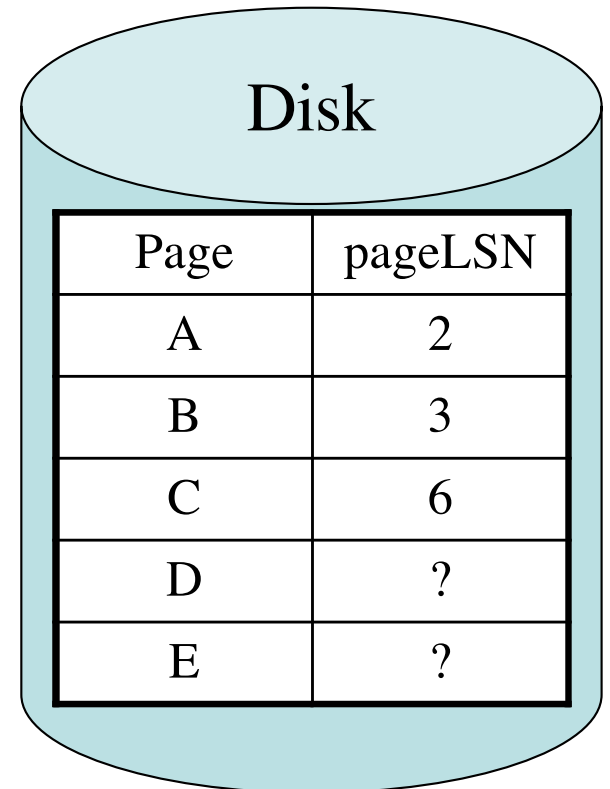
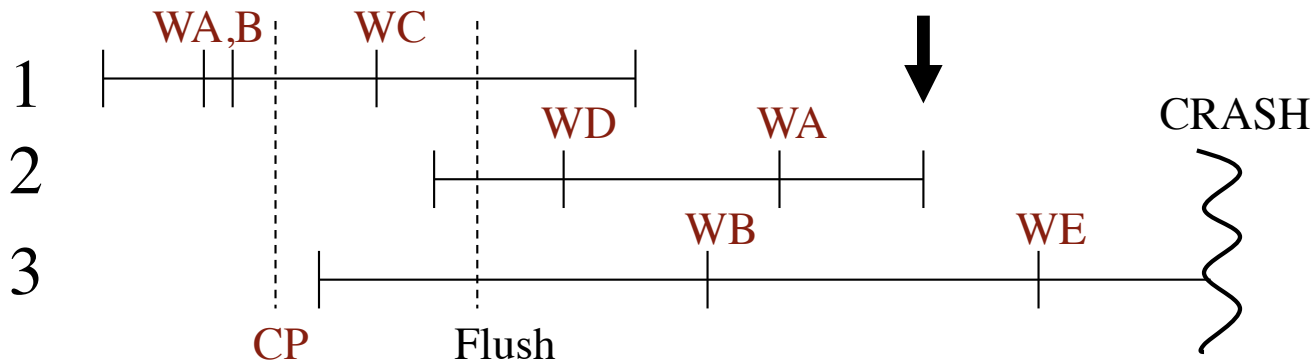
lastLSN	TID
10	3
11	2

dirtyPgTable

pgNo	recLSN
D	8
B	10
A	11

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

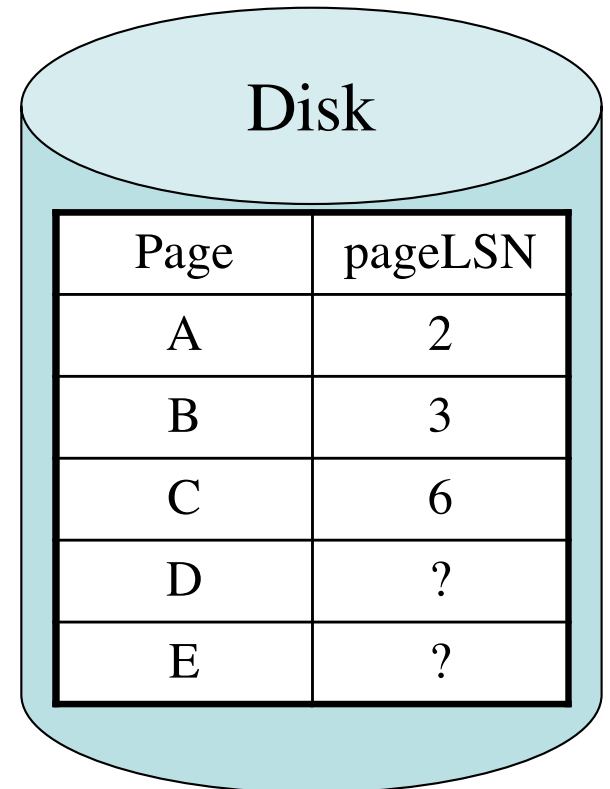
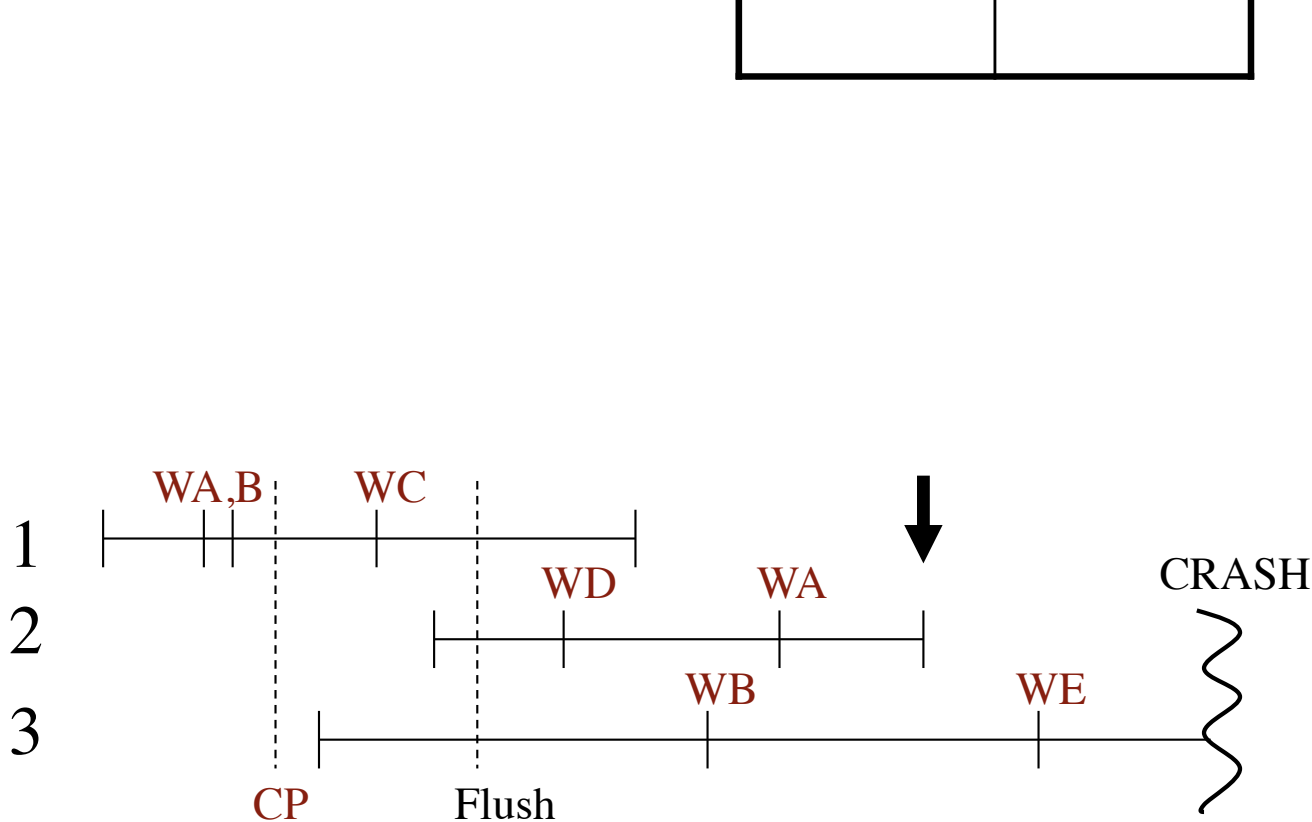
lastLSN	TID
10	3

dirtyPgTable

pgNo	recLSN
D	8
B	10
A	11

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



ARIES Data Structures

xactionTable

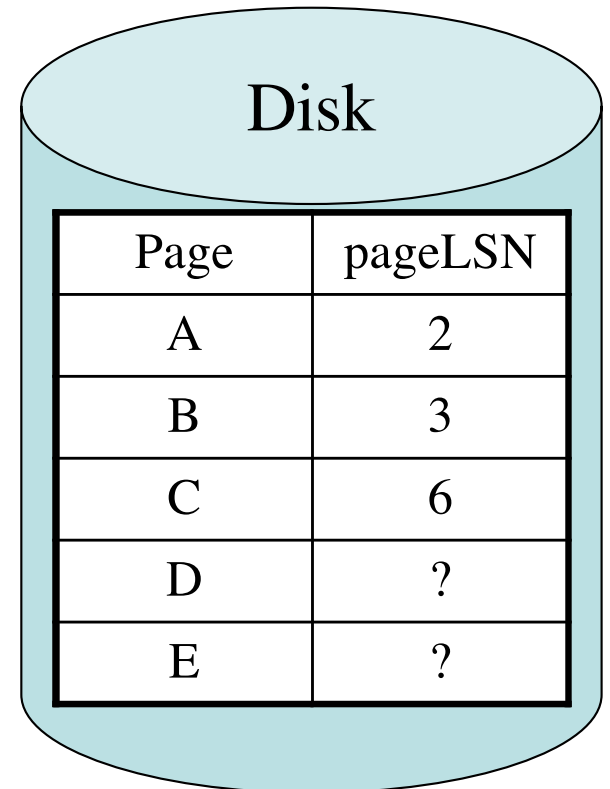
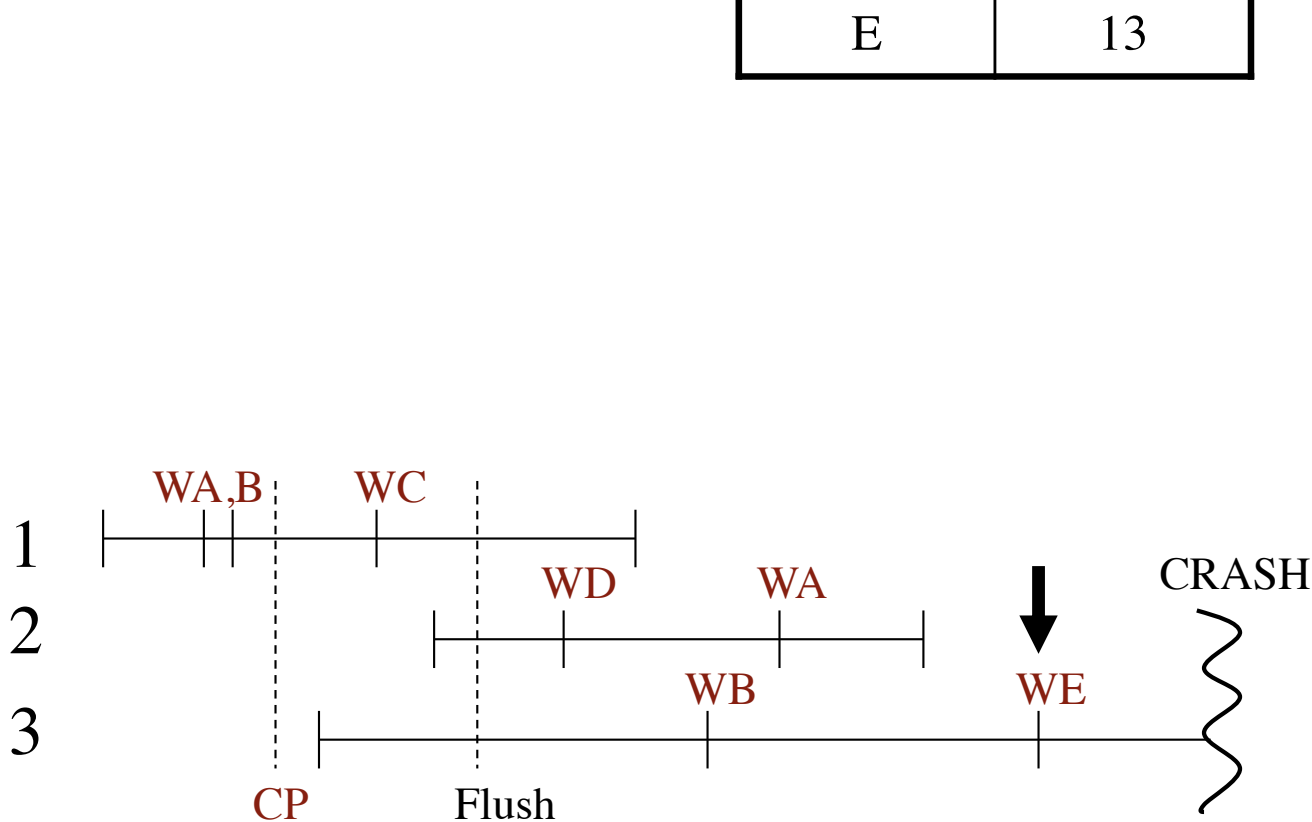
lastLSN	TID
13	3

dirtyPgTable

pgNo	recLSN
D	8
B	10
A	11
E	13

Checkpoint

xactionTable	3-1
dirtyPgTable	A-2, B-3



Crash Recovery

- 3 Phases
 - Analysis
 - Rebuild data structures
 - Determine winners & losers
 - Redo
 - “Repeat history”
 - Why?
 - Undo
 - Undo Losers

Analysis Pass

- Goal: reconstruct the state of the transaction table and the dirty page table at the time the crash occurred.
- Play log forward
 - Add and remove xactions to/from the transaction table on SOT and COMMIT/ABORT
 - Update the lastLSN on writes
 - Update the dirty page table as writes happen


State After Analysis

- After analysis, what can we say about dirty page table and transaction table?
- Txn table tells us what to UNDO
- Dirty pages is a conservative list of pages that need to be REDO
- Why is it conservative?
 - Because we don't actually know what is on disk; some pages may already have updates applied

Where to Begin Analysis


- Beginning of log?
 - Ok, but may require us to scan a lot of log
- Last checkpoint!
- How do we find it?
 - Keep a pointer to the checkpoint at a well-known place on disk

Analysis



LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	A
3	UP	1	3	B
4	CP			
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

Analysis



LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

xactionTable

lastLSN	TID
3	1


dirtyPgTable

pgNo	recLSN
A	2
B	3

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

Analysis



LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

xactionTable

lastLSN	TID
3	1
5	3

dirtyPgTable

pgNo	recLSN
A	2
B	3

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

Analysis



LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

xactionTable

lastLSN	TID
6	1
5	3


dirtyPgTable

pgNo	recLSN
A	2
B	3
C	6

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

Analysis



LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

xactionTable

lastLSN	TID
6	1
5	3
7	2

dirtyPgTable

pgNo	recLSN
A	2
B	3
C	6

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

Analysis

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E



xactionTable

lastLSN	TID
6	1
5	3
8	2

dirtyPgTable

pgNo	recLSN
A	2
B	3
C	6
D	8

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

Analysis

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3		
11	UP	2		
12	EOT	2		
13	UP	3		

Dirty page table doesn't reflect true state on disk.

Conservative: *at least* all previous LSNs are on disk



xactionTable

lastLSN	TID
13	3

Losers

dirtyPgTable

pgNo	recLSN
A	2
B	3
C	6
D	8
E	13

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

Redo

- Where to begin?
 - Checkpoint?
 - $\text{Min}(\text{recLSN})!$ – earliest unflushed update
- What to REDO
 - Everything?
 - Slow
 - Problematic if using logical (escrow) logging
 - Redo an update UNLESS:
 - Page is not in dirtyPgTable
 - Page flushed prior to checkpoint, didn't dirty
 - $\text{LSN} < \text{recLSN}$
 - Page flushed & redirtied prior to checkpoint
 - $\text{LSN} \leq \text{pageLSN}$
 - Page flushed after checkpoint

Only step that requires going to disk

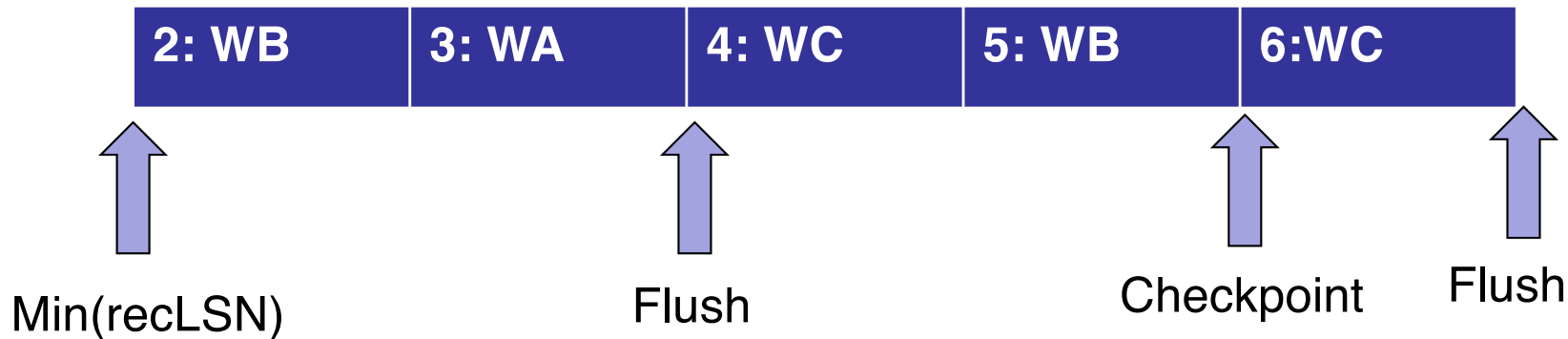
dirtyPgTable

pgNo	recLSN
A	2
B	3
C	6
D	8
E	13

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

REDO Conditions Example



dirtyPgTable
@ CP

pgNo	recLSN
B	5
C	4

Redo an update UNLESS:

Page is not in dirtyPgTable **A/LSN 3**

Page flushed prior to checkpoint, didn't dirty

LSN < recLSN **B/LSN 2**

Page flushed & redirtied prior to checkpoint

LSN <= pageLSN **C/LSN 6**

Page flushed after checkpoint

Disk

Page	pageLSN
A	3
B	5
C	6

Redo Example

Redo UNLESS

- Page is not in dirtyPgTable
- LSN < recLSN
- LSN <= pageLSN

DirtyPgTable

pgNo	recLSN
B	3
C	6
D	8
E	13

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	A
3	UP	1	3	B
4	CP			
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E



Flush



Redo Example

Redo UNLESS

- Page is not in dirtyPgTable
- LSN < recLSN
- LSN <= pageLSN

DirtyPgTable

pgNo	recLSN
C	6
D	8
E	13

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	A
3	UP	1	3	B
4	CP			
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E



Flush

Redo Example

Redo UNLESS

- Page is not in dirtyPgTable
- LSN < recLSN
- LSN <= pageLSN

DirtyPgTable

pgNo	recLSN
D	8
E	13

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	A
3	UP	1	3	B
4	CP			
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

Flush

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?





Redo Example

Redo UNLESS

- Page is not in dirtyPgTable
- LSN < recLSN
- LSN <= pageLSN

DirtyPgTable

pgNo	recLSN
D	8
E	13

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	
3	UP	1	3	
4	CP			
5	SOT	3		
6	UP	1	3	
7	SOT	2		
8	UP	2	7	D 
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

Flush 

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

Redo Example

Redo UNLESS

- Page is not in dirtyPgTable
- LSN < recLSN
- LSN <= pageLSN

DirtyPgTable

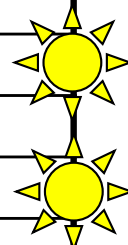
pgNo	recLSN
B	10
D	8
E	13

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	
3	UP	1	3	
4	CP			
5	SOT	3		
6	UP	1	3	
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

Flush



Redo Example

Redo UNLESS







- Page is not in dirtyPgTable
- LSN < recLSN
- LSN <= pageLSN

DirtyPgTable

pgNo	recLSN
A	11
B	10
D	8
E	13

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	
3	UP	1	3	
4	CP			
5	SOT	3		
6	UP	1	3	
7	SOT	2		
8	UP	2	7	D 
9	EOT	1	6	
10	UP	3	5	B 
11	UP	2	8	A 
12	EOT	2	11	
13	UP	3	10	E

Flush




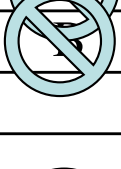
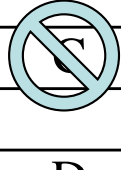


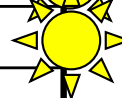

Redo Example

Redo UNLESS

- Page is not in dirtyPgTable
- LSN < recLSN
- LSN <= pageLSN

DirtyPgTable

pgNo	recLSN
A	11
B	10
D	8
E	13

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	
3	UP	1	3	
4	CP			
5	SOT	3		
6	UP	1	3	
7	SOT	2		
8	UP	2	7	D 
9	EOT	1	6	
10	UP	3	5	B 
11	UP	2	8	A 
12	EOT	2	11	
13	UP	3	10	E 

Flush

Disk

Page	pageLSN
A	2
B	3
C	6
D	?
E	?

State identical to pre-crash state



Undo

- Walk backwards, following prevLSNs to UNDO losers

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	A
3	UP	1	3	B
4	CP			
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E



xactionTable

lastLSN	TID
13	3

Undo

- Walk backwards, following prevLSNs to UNDO losers

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	A
3	UP	1	3	B
4	CP			
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E



xactionTable

lastLSN	TID
13	3

Undo

- Walk backwards, following prevLSNs to UNDO losers

LSN	Type	Tid	PrevLSN	Data
1	SOT	1		
2	UP	1	2	A
3	UP	1	3	B
4	CP			
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E



xactionTable

lastLSN	TID
13	3

- Why can we just blindly apply UNDOs? Repeated history!

Study Break

LSN	Txn ID	Type	Page ID / Object
10	Checkpoint		
11	T1	SOT	
12	T1	UP	P1/A
13	T2	SOT	
14	T3	SOT	
15	T2	UP	P5/B
16	T2	Commit	
17	T3	UP	P3/C
18	Checkpoint		
19	T3	UP	P3/C
20	T3	Commit	

No flushes occur during the execution of these transactions. At the time of checkpoint1, the dirty page table and the transaction table are both empty.

recLSN = first LSN that dirtied the page

lastLSN = most recent log record written by the transaction

Study Break

LSN	Txn ID	Type	Page ID / Object
10	Checkpoint		
11	T1	SOT	
12	T1	UP	P1/A
13	T2	SOT	
14	T3	SOT	
15	T2	UP	P5/B
16	T2	Commit	
17	T3	UP	P3/C
18	Checkpoint		
19	T3	UP	P3/C
20	T3	Commit	

PageID	recLSN
P1	12
P3	17
P5	15

Dirty Page Table

Transaction ID	lastLSN
T1	12

Transaction Table

- What must the status of the tables have been at the time of the crash?

Study Break # 2

LSN	Txn ID	Type	Page ID / Object
10	Checkpoint		
11	T1	SOT	
12	T1	UP	P1/A
13	T2	SOT	
14	T3	SOT	
15	T2	UP	P5/B
16	T2	Commit	
17	T3	UP	P3/C
18	Checkpoint		
19	T3	UP	P3/C
20	T3	Commit	

No flushes occur during the execution of these transactions. At the time of checkpoint1, the dirty page table and the transaction table are both empty.

1. At what LSN does the analysis phase begin? **18**
2. At what LSN does the REDO phase begin? **Min(recLSN) = 12**
3. What is the first LSN that is undone? **12**

Truncating Log

- Do we have to keep log forever?
- What is the earliest point in the log we will ever look at?

$\min(\text{last checkpoint}, \min(\text{recLSN}))$

→ we can safely truncate anything earlier

Compensation Log Records (CLR)s

- CLR record written after each UNDO
- Avoid repeating UNDO work
- Why?
 - Because UNDO is logical, and we don't check if records have already been UNDONE. Could get into trouble if re-undid some logical operation.

UNDO with CLR

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

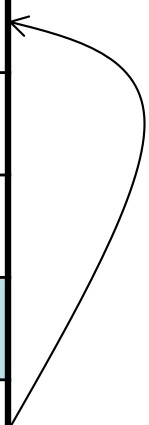
Losers: 3

UNDO with CLR

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E

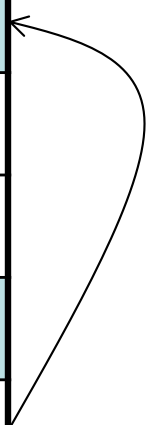
UNDO with CLR

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E
14	CLR	3	13	10



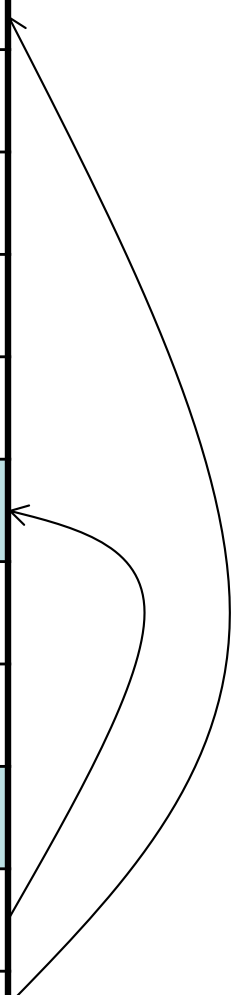
UNDO with CLR

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E
14	CLR	3	13	E, 10



UNDO with CLR

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E
14	CLR	3	13	E, 10
15	CLR	3	14	B, 5



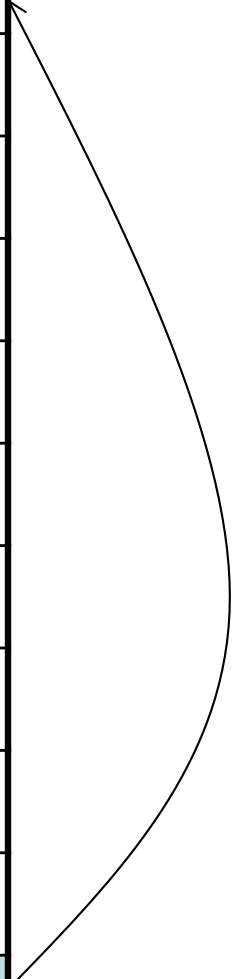
The diagram illustrates the undo process. Row 14 (CLR, Tid 3, PrevLSN 13, Data E, 10) and Row 15 (CLR, Tid 3, PrevLSN 14, Data B, 5) are shown with arrows pointing to Row 10 (UP, Tid 3, PrevLSN 5, Data B). This indicates that the undo operation for row 14 involves reverting to the state of row 10, and the undo operation for row 15 involves reverting to the state of row 14.

UNDO with CLR

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E
14	CLR	3	13	E, 10
15	CLR	3	14	B, 5

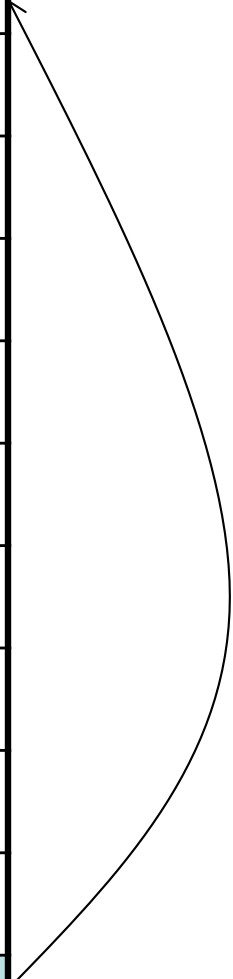
UNDO with CLR

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E
14	CLR	3	13	E, 10
15	CLR	3	14	B, 5



UNDO with CLR

LSN	Type	Tid	PrevLSN	Data
5	SOT	3		
6	UP	1	3	C
7	SOT	2		
8	UP	2	7	D
9	EOT	1	6	
10	UP	3	5	B
11	UP	2	8	A
12	EOT	2	11	
13	UP	3	10	E
14	CLR	3	13	E, 10
15	CLR	3	14	B, 5
16	EOT	3	15	



REDO with CLR

- REDO CLR on crash recovery
 - Use REDO rules to check if updates in CLR have already been done
 - Avoids repeating operational (escrow) operations
 - After processing CLR, update lastLSN field in dirtyPgTable
 - Allows UNDO to start from the right place, should we checkpoint while UNDOing

ARIES/Logging Recap

- NO FORCE, STEAL logging
- Use write ahead logging protocol
- Must FORCE log on COMMIT
- Periodically take (lightweight) checkpoints
- Asynchronously flush disk pages (without logging)

Disaster Recovery

- What if:
 - Disk on machine fails
 - Computer won't restart
 - Data center loses power
 - ...
- Solution: replication



Replication

- Typical approach: dedicated “hot standby”
 - Kept up to date via “log shipping” – it executes operations in the log in identical order to the primary
- May have several replicas, one nearby in local data center, one further away
 - “Half-width of a hurricane”
- Replicas often used for read-only queries
 - Have excess capacity because they are not processing xactions, just replaying log

Replica Fail Over

- On failure, start directing queries to replica
- Bring up new replica
 - Using, e.g., nightly backup + log
- Complex in practice:
 - Have to be really sure the database failed
 - Many organizations rely on manual failover
 - Failover needs to be tested frequently
 - Replication load can be significant

Transactions Summary

- Transactions provide a powerful way to isolate concurrent operations on the DB
- Studied two concurrency control methods two-phase locking and OCC
- Saw how write-ahead logging can be used to provide recoverability and roll-back
- Next time: distributed DBs, and distributed txns