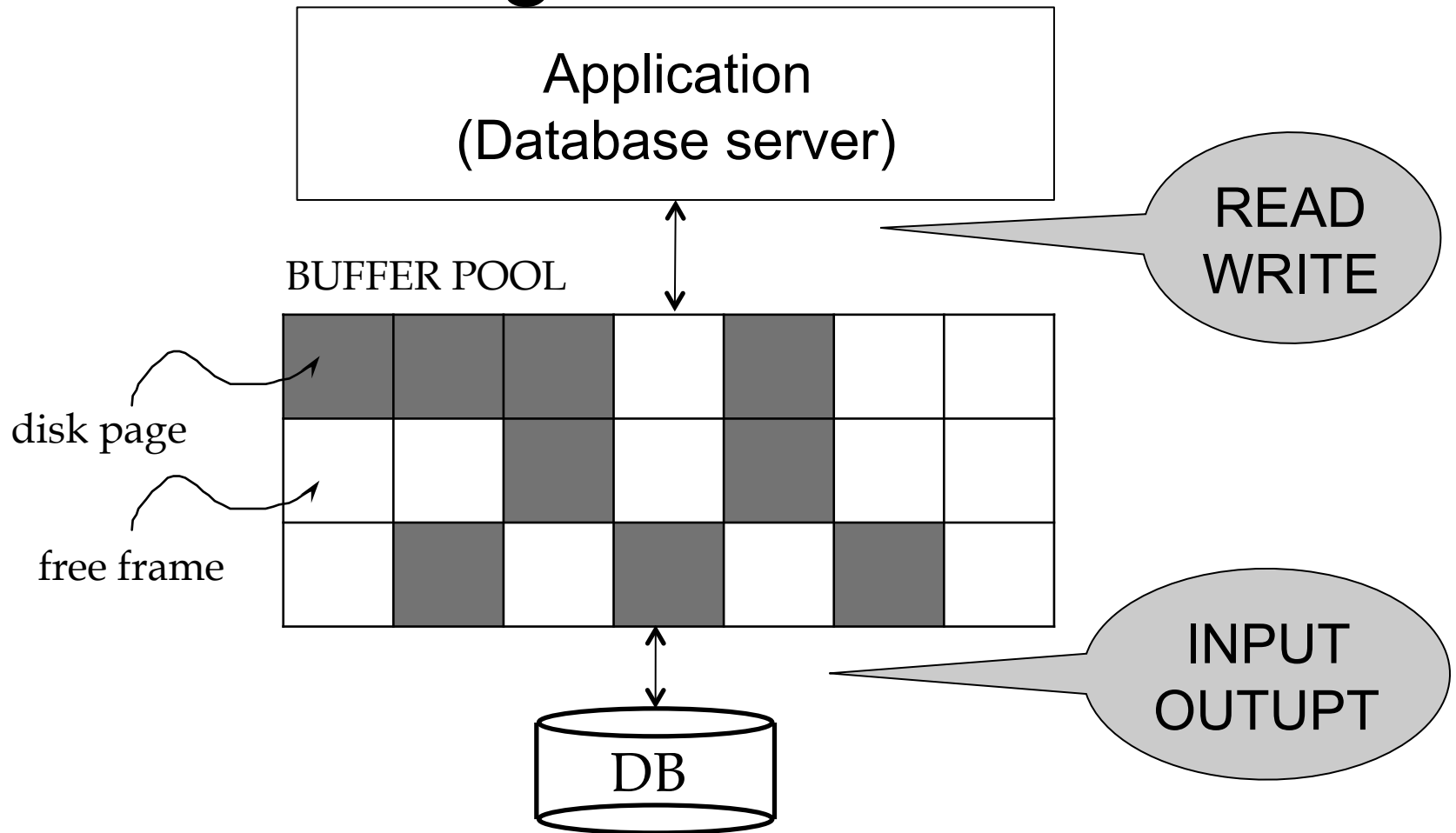


CSE544

Transactions: Recovery

Thursday, January 27, 2011

Buffer Management in a DBMS



Large gap between disk I/O and memory → Buffer pool

Page Replacement Policies

- LRU = expensive
 - Next slide
- Clock algorithm = cheaper alternative
 - Read in the book

Both work well in OS, but not always in DB

Least Recently Used (LRU)

Most recent

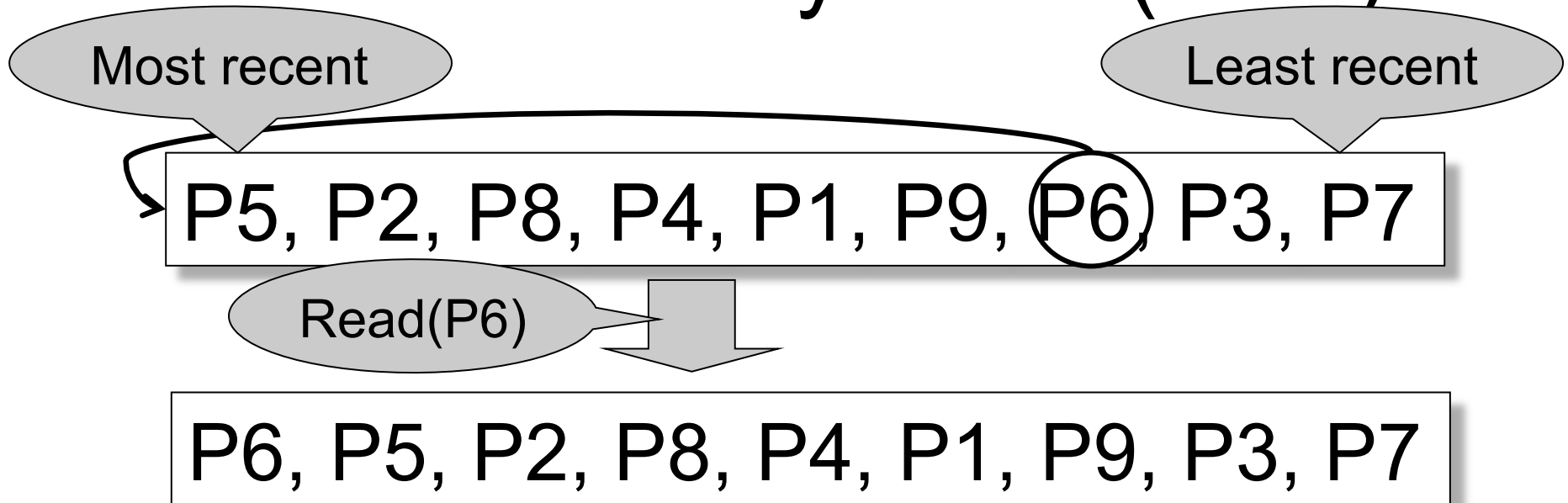
Least recent

P5, P2, P8, P4, P1, P9, P6, P3, P7

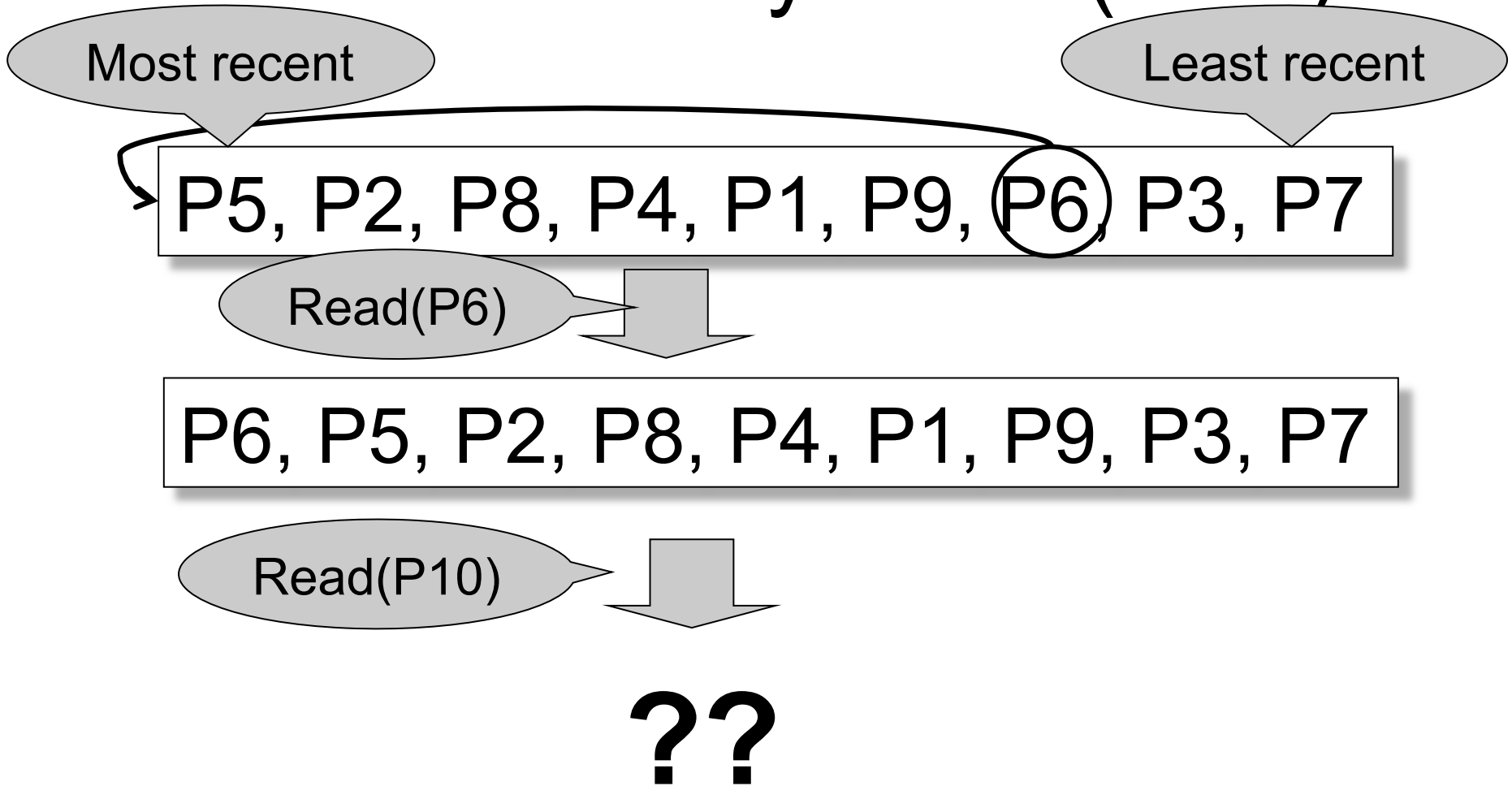
Read(P6)

??

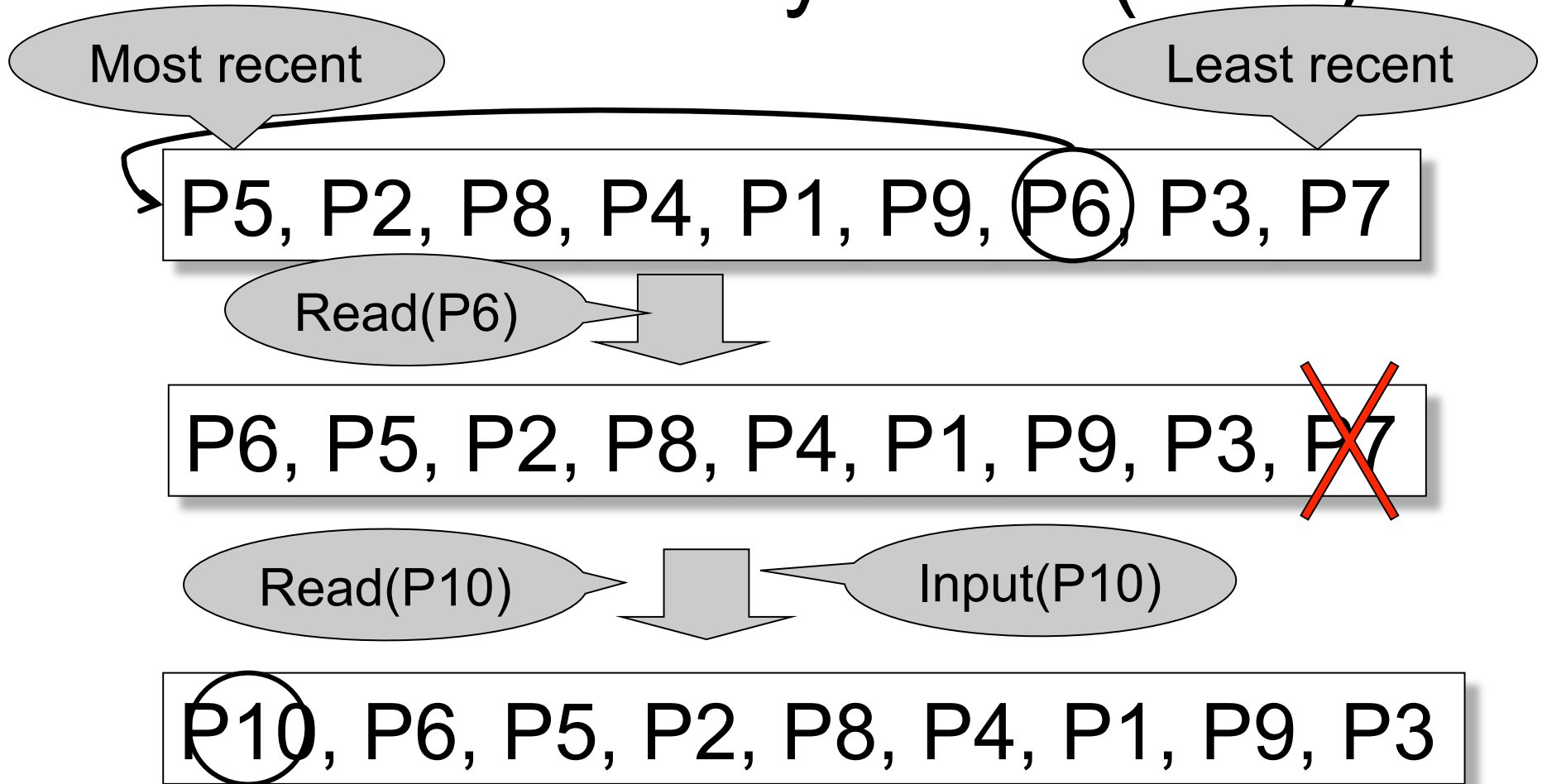
Least Recently Used (LRU)



Least Recently Used (LRU)



Least Recently Used (LRU)



Atomic Transactions

- **FORCE or NO-FORCE**
 - Should all updates of a transaction be forced to disk before the transaction commits?
- **STEAL or NO-STEAL**
 - Can an update made by an uncommitted transaction overwrite the most recent committed value of a data item on disk?

Atomicity: FORCE + NO-STEAL

Performance (e.g. LRU): NO-FORCE+STEAL

Atomic Transactions

- **NO-FORCE**: Pages of committed transactions not yet written to disk
- **STEAL**: Pages of uncommitted transactions already written to disk

In either case, *Atomicity* is violated

Notations

- **READ(X,t)**
 - copy element X to transaction local variable t
- **WRITE(X,t)**
 - copy transaction local variable t to element X
- **INPUT(X)**
 - read element X to memory buffer
- **OUTPUT(X)**
 - write element X to disk

Recovery

Write-ahead log = _A file that records every single action of all running transactions

- *Force* log entry to disk
- After a crash, transaction manager reads the log and finds out exactly what the transactions did or did not

Example

```
START TRANSACTION
READ(A,t);
t := t*2;
WRITE(A,t);
READ(B,t);
t := t*2;
WRITE(B,t)
COMMIT;
```

Atomicity:

Both A and B
are multiplied by 2,
or none is.

READ(A,t); t := t*2; WRITE(A,t);
 READ(B,t); t := t*2; WRITE(B,t)

Transaction

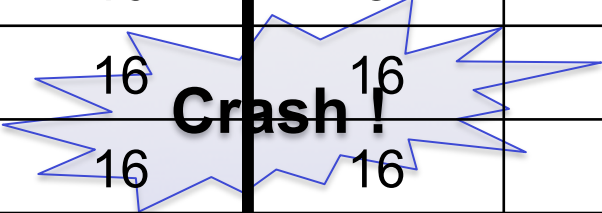
Buffer pool

Disk

| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |

Is this bad ?

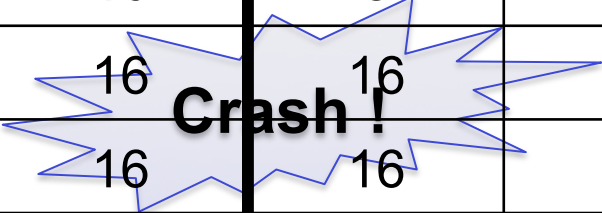
| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |



Is this bad ?

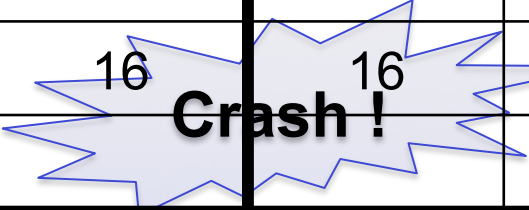
Yes it's bad: A=16, B=8....

| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |



Is this bad ?

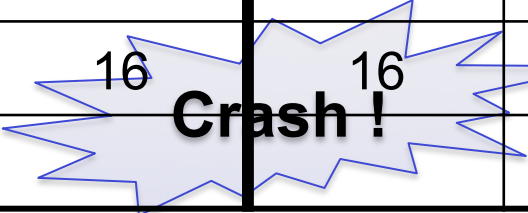
| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |



Is this bad ?

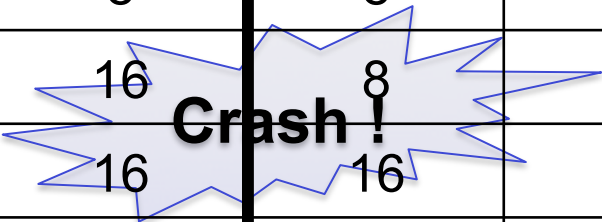
Yes it's bad: A=B=16, but not committed

| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |



Is this bad ?

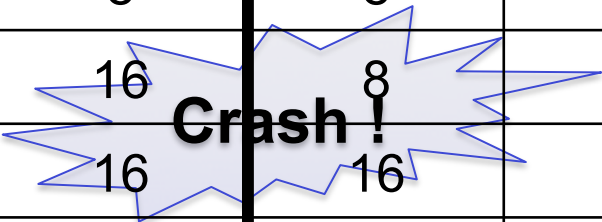
| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |



Is this bad ?

No: that's OK

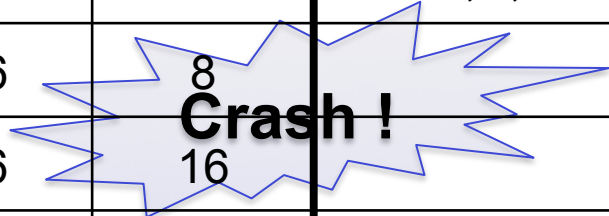
| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| INPUT(A) | | 8 | | 8 | 8 |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |
| COMMIT | | | | | |



UNDO Log

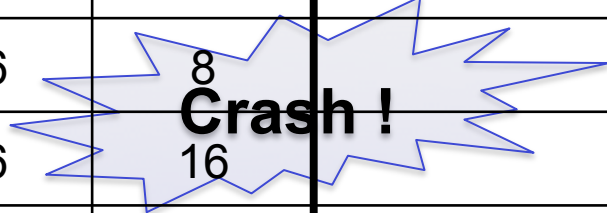
| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |

| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |



WHAT DO WE DO ?

| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |



WHAT DO WE DO ?

We UNDO by setting B=8 and A=8

| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |

Crash!

What do we do now ?

| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |

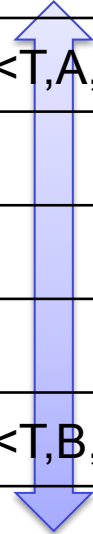


What do we do now ?

Nothing: log contains COMMIT

| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| INPUT(A) | | | | | 8 | |
| READ(A,t) | 8 | | | | 8 | |
| t:=t*2 | 16 | 8 | | | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |

When must we force pages to disk ?



| Action | t | Mem A | Mem B | Disk A | Disk B | UNDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| INPUT(A) | | 8 | | 8 | 8 | |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,8> |
| INPUT(B) | 16 | 16 | 8 | 8 | 8 | |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,8> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |
| COMMIT | | | | | | <COMMIT T> |

RULES: log entry before OUTPUT before COMMIT

REDO Log

Is this bad ?

| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| | | | | | |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| COMMIT | | | | | |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |

Crash !

Is this bad ?

Yes, it's bad: A=16, B=8

| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| | | | | | |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| COMMIT | | | | | |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |

Crash !

Is this bad ?

| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| | | | | | |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| COMMIT | | | | | |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |

Crash !

Is this bad ?

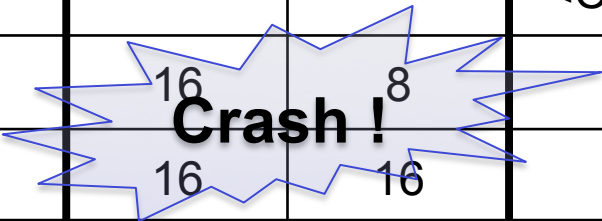
Yes, it's bad: T committed but A=B=8

| Action | t | Mem A | Mem B | Disk A | Disk B |
|------------|----|-------|-------|--------|--------|
| | | | | | |
| READ(A,t) | 8 | 8 | | 8 | 8 |
| t:=t*2 | 16 | 8 | | 8 | 8 |
| WRITE(A,t) | 16 | 16 | | 8 | 8 |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 |
| COMMIT | | | | | |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 |

Crash !

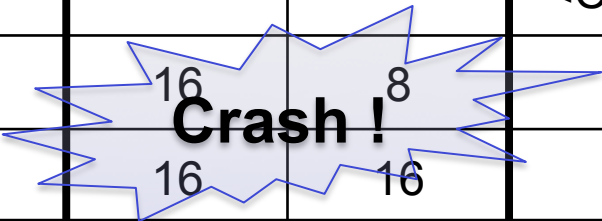
| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|------------|----|-------|-------|--------|--------|-----------------|
| | | | | | | <START T> |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT | | | | | | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |

| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT | | | | | | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |



How do we recover ?

| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT | | | | | | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |

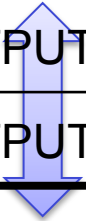


How do we recover ?

We **REDO** by setting A=16 and B=16

| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| READ(A,t) | 8 | 8 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 8 | 8 | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | 8 | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT | | | | | | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |

When must we force pages to disk ?



| Action | t | Mem A | Mem B | Disk A | Disk B | REDO Log |
|------------|----|-------|-------|--------|--------|------------|
| | | | | | | <START T> |
| READ(A,t) | 8 | 8 | | 8 | 8 | |
| t:=t*2 | 16 | 8 | | 8 | 8 | |
| WRITE(A,t) | 16 | 16 | | 8 | 8 | <T,A,16> |
| READ(B,t) | 8 | 16 | 8 | 8 | 8 | |
| t:=t*2 | 16 | 16 | 8 | 8 | 8 | |
| WRITE(B,t) | 16 | 16 | 16 | 8 | 8 | <T,B,16> |
| COMMIT | | | | | | <COMMIT T> |
| OUTPUT(A) | 16 | 16 | 16 | 16 | 8 | |
| OUTPUT(B) | 16 | 16 | 16 | 16 | 16 | |

RULE: OUTPUT *after* COMMIT

Comparison Undo/Redo

- Undo logging: OUTPUT must be done early:
 - Inefficient
- Redo logging: OUTPUT must be done late:
 - Inflexible

Checkpointing

- To ensure recovery we must read from the beginning of the log: this is too inefficient
- Checkpointing: periodically write information to the log to allow us to process only the tail of the log

ARIES Recovery Manager

- A redo/undo log
- **Physiological logging**
 - Physical logging for REDO
 - Logical logging for UNDO
- Efficient checkpointing



Why ?

ARIES Recovery Manager

Log entries:

- $\langle \text{START } T \rangle$ -- when T begins
- Update: $\langle T, X, u, v \rangle$
 - T updates X , old value= u , new value= v
 - In practice: undo only and redo only entries
- $\langle \text{COMMIT } T \rangle$ or $\langle \text{ABORT } T \rangle$
- CLR's – we'll talk about them later.

ARIES Recovery Manager

Rule:

- If T modifies X, then $\langle T, X, u, v \rangle$ must be written to disk before OUTPUT(X)

We are free to OUTPUT early or late

LSN = Log Sequence Number

- LSN = identifier of a log entry
 - Log entries belonging to the same txn are linked
- Each page contains a **pageLSN**:
 - LSN of log record for latest update to that page
 - Will serve to determine if an update needs to be redone

ARIES Data Structures

- **Active Transactions Table**
 - Lists all running txns (active txns)
 - For each txn: **lastLSN** = most recent update by txn
- **Dirty Page Table**
 - Lists all dirty pages
 - For each dirty page: **recoveryLSN** (**recLSN**)= first LSN that caused page to become dirty
- **Write Ahead Log**
 - LSN, **prevLSN** = previous LSN for same txn

$W_{T100}(P7)$
 $W_{T200}(P5)$
 $W_{T200}(P6)$
 $W_{T100}(P5)$

ARIES Data Structures

Dirty pages

| pageID | recLSN |
|--------|--------|
| P5 | 102 |
| P6 | 103 |
| P7 | 101 |

Log (WAL)

| LSN | prevLSN | transID | pageID | Log entry |
|-----|---------|---------|--------|-----------|
| 101 | - | T100 | P7 | |
| 102 | - | T200 | P5 | |
| 103 | 102 | T200 | P6 | |
| 104 | 101 | T100 | P5 | |

Active transactions

| transID | lastLSN |
|---------|---------|
| T100 | 104 |
| T200 | 103 |

Buffer Pool

| | | |
|-------------------|-------------------|-------------------|
| P8 | P2 | ... |
| | ... | |
| P5 PageLSN=104 | P6 PageLSN=103 | P7 PageLSN=101 |
| | | |

ARIES Normal Operation

T writes page P

- What do we do ?

ARIES Normal Operation

T writes page P

- What do we do ?
- Write $\langle T, P, u, v \rangle$ in the Log
- **pageLSN=LSN**
- **lastLSN=LSN**
- **recLSN**=if isNull then **LSN**

ARIES Normal Operation

Buffer manager wants to OUTPUT(P)

- What do we do ?

Buffer manager wants INPUT(P)

- What do we do ?

ARIES Normal Operation

Buffer manager wants to OUTPUT(P)

- Flush log up to **pageLSN**
- Remove P from **Dirty Pages** table

Buffer manager wants INPUT(P)

- Create entry in **Dirty Pages** table
recLSN = NULL

ARIES Normal Operation

Transaction T starts

- What do we do ?

Transaction T commits/aborts

- What do we do ?

ARIES Normal Operation

Transaction T starts

- Write <START T> in the **log**
- New entry T in **Active TXN**;
lastLSN = null

Transaction T commits/aborts

- Write <COMMIT T>
- Flush **log** up to this entry

Checkpoints

Write into the log

- Entire **active transactions table**
- Entire **dirty pages table**

Recovery always starts by analyzing latest checkpoint

Background process periodically flushes dirty pages to disk

ARIES Recovery

1. Analysis pass

- Figure out what was going on at time of crash
- List of dirty pages and active transactions

2. Redo pass (repeating history principle)

- Redo all operations, even for transactions that will not commit
- Get back to state at the moment of the crash

3. Undo pass

- Remove effects of all uncommitted transactions
- Log changes during undo in case of another crash during undo

ARIES Method Illustration

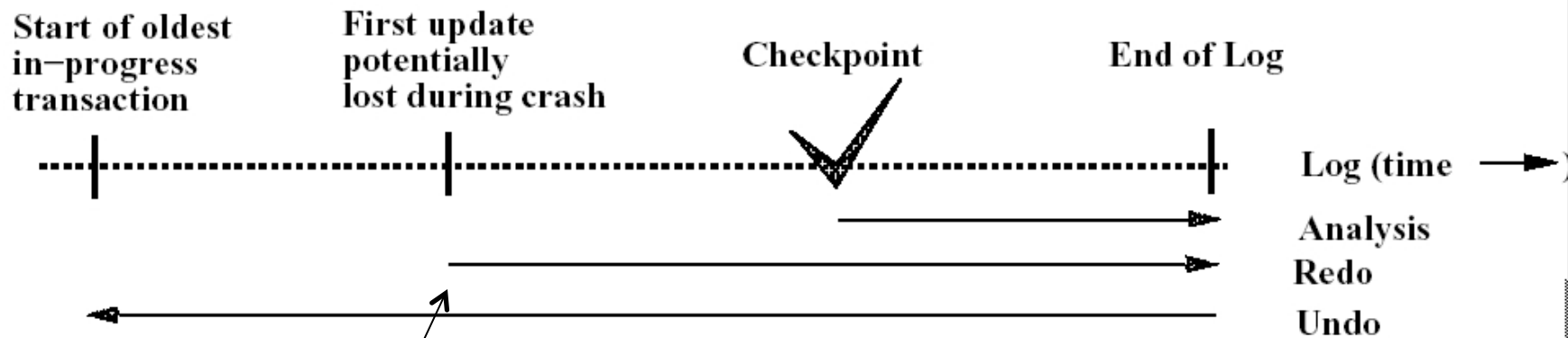


Figure 3: The Three Passes of ARIES Restart

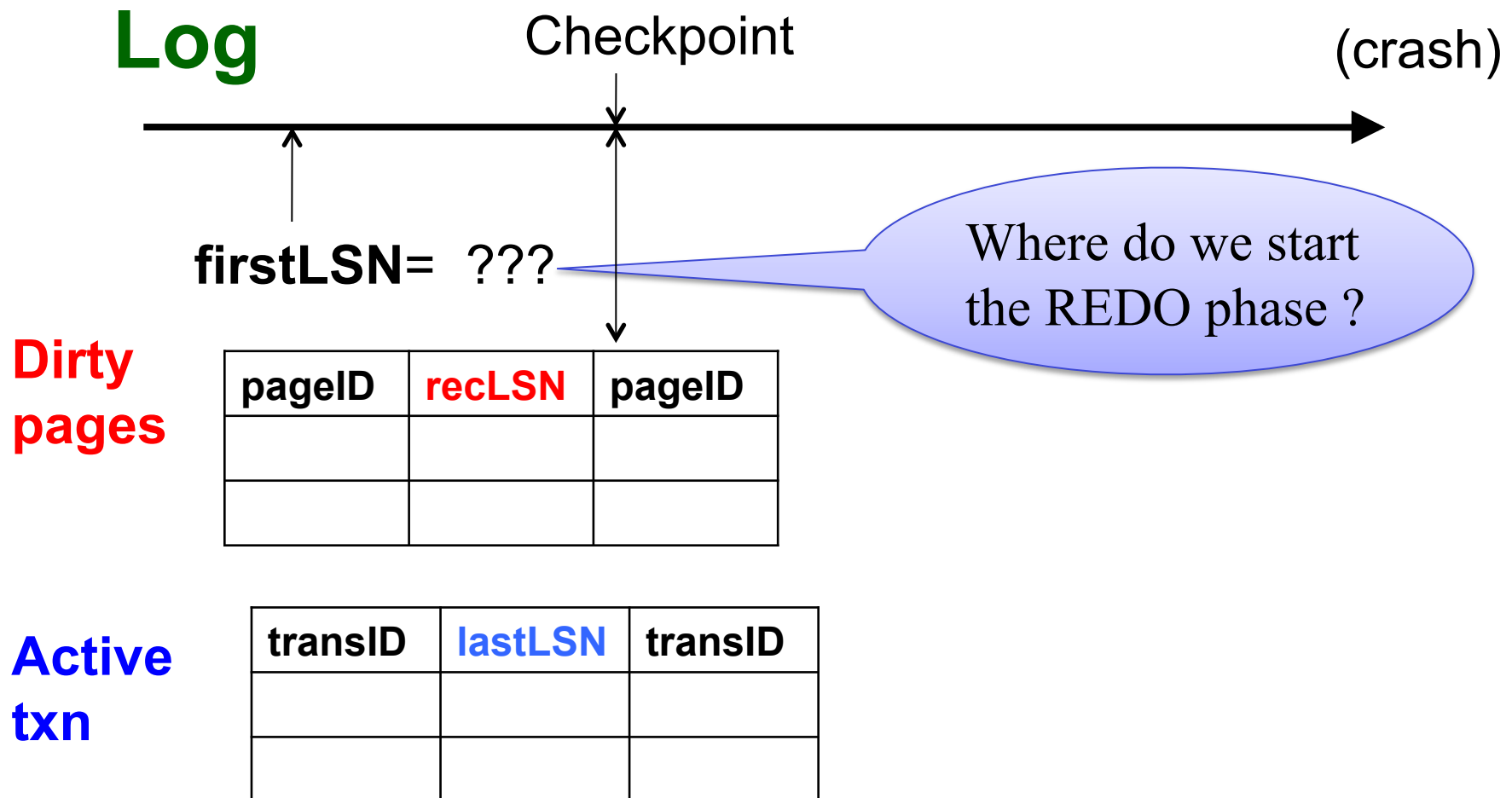
First undo and first redo log entry might be in reverse order

[Figure 3 from Franklin97]

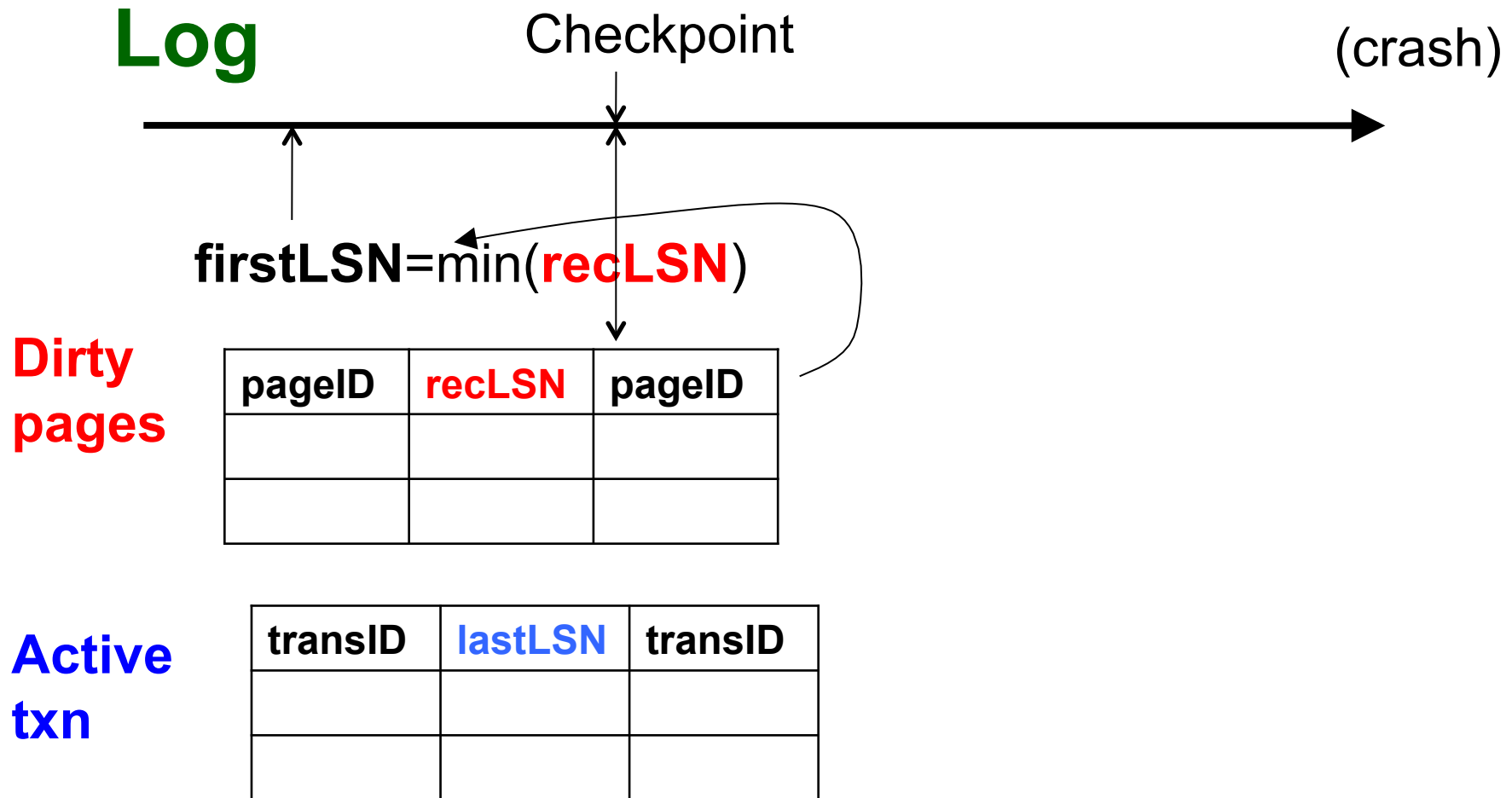
1. Analysis Phase

- Goal
 - Determine point in log where to start REDO
 - Determine set of dirty pages when crashed
 - Conservative estimate of dirty pages
 - Identify active transactions when crashed
- Approach
 - Rebuild **active transactions table** and **dirty pages table**
 - Reprocess the log from the checkpoint
 - Only update the two data structures
 - Compute: **firstLSN** = smallest of all **recoveryLSN**

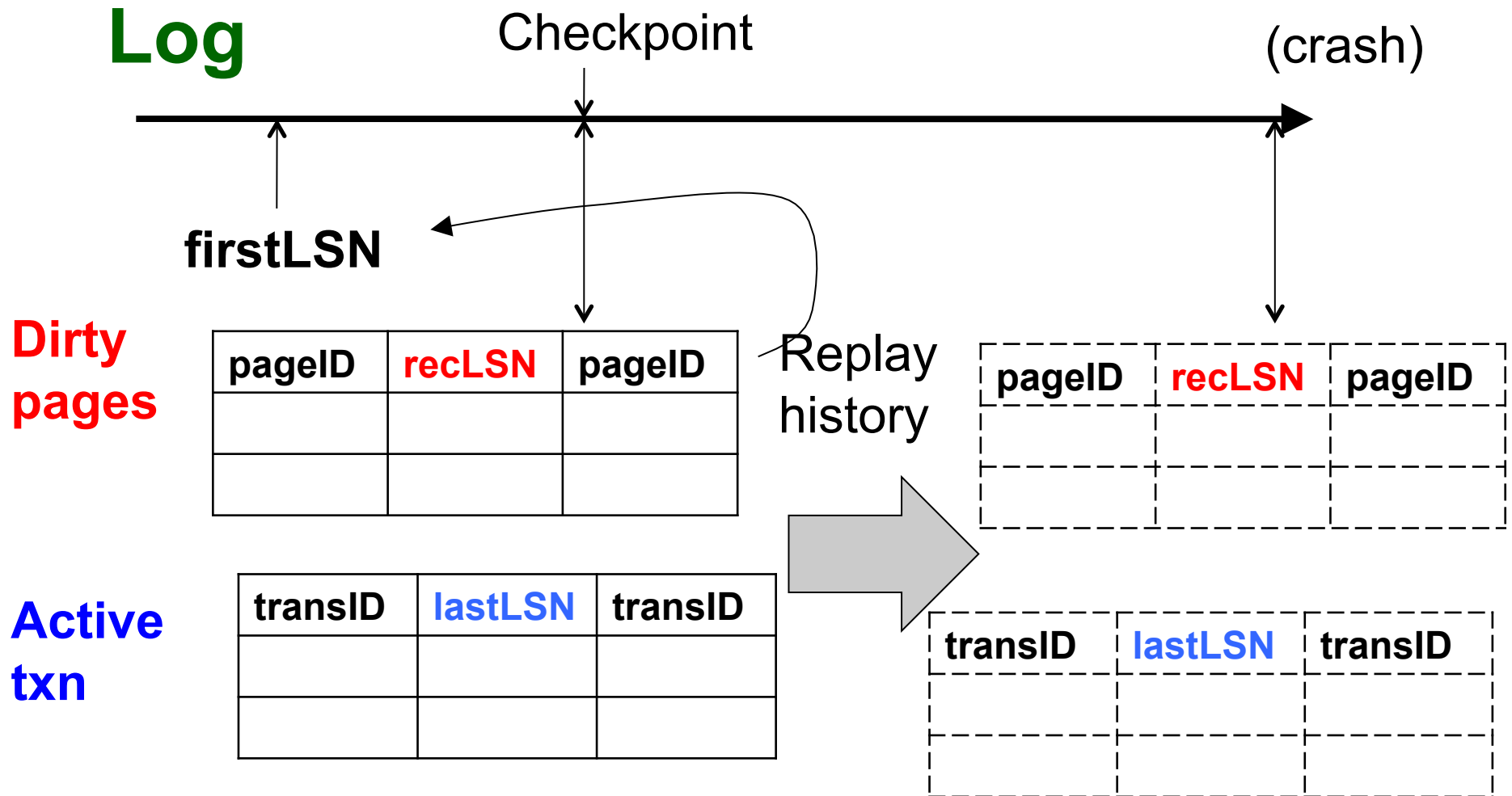
1. Analysis Phase



1. Analysis Phase



1. Analysis Phase



2. Redo Phase

Main principle: replay history

- Process Log forward, starting from **firstLSN**
- Read every log record, sequentially
- Redo actions are not recorded in the log
- Needs the **Dirty Page Table**

2. Redo Phase: Details

For each **Log** entry record **LSN**: $\langle T, P, u, v \rangle$

- Re-do the action $P=u$ and $WRITE(P)$
- But which actions can we skip, for efficiency ?

2. Redo Phase: Details

For each **Log** entry record **LSN**: $\langle T, P, u, v \rangle$

- If P is not in **Dirty Page** then **no update**
- If **recLSN** $>$ **LSN**, then **no update**
- **INPUT(P)** (read page from disk):
If **pageLSN** $>$ **LSN**, then **no update**
- Otherwise perform update

2. Redo Phase: Details

What happens if system crashes during REDO ?

2. Redo Phase: Details

What happens if system crashes during REDO ?

We REDO again ! Each REDO operation is *idempotent*: doing it twice is the as as doing it once.

3. Undo Phase

Main principle: “logical” undo

- Start from end of **Log**, move backwards
- Read only affected log entries
- Undo actions *are* written in the Log as special entries: **CLR** (Compensating Log Records)
- **CLRs** are redone, but never undone

3. Undo Phase: Details

- “Loser transactions” = uncommitted transactions in [Active Transactions Table](#)
- **ToUndo** = set of [lastLSN](#) of loser transactions

3. Undo Phase: Details

While **ToUndo** not empty:

- Choose most recent (largest) **LSN** in **ToUndo**
- If **LSN** = regular record $\langle T, P, u, v \rangle$:
 - Undo v
 - Write a **CLR** where **CLR.undoNextLSN** = **LSN.prevLSN**
- If **LSN** = **CLR record**:
 - Don't undo !
- if **CLR.undoNextLSN** not null, insert in **ToUndo**
otherwise, write $\langle \text{END TRANSACTION} \rangle$ in log

3. Undo Phase: Details

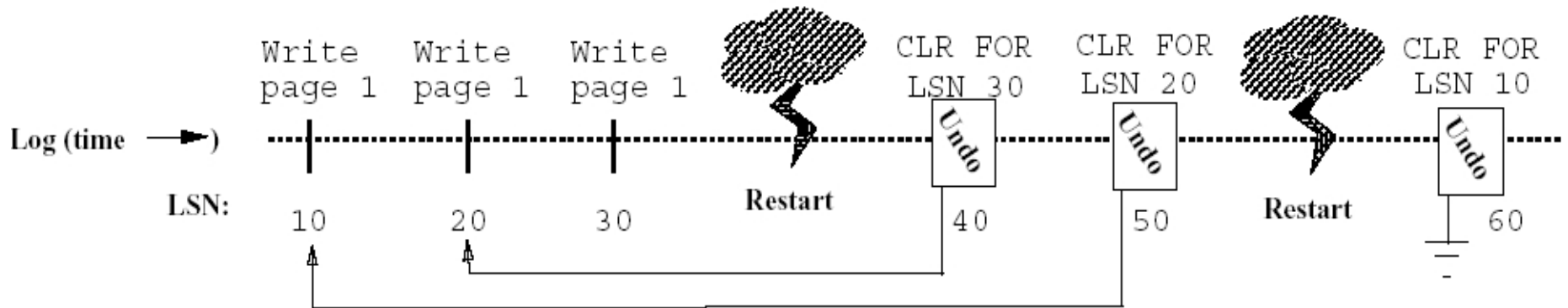


Figure 4: The Use of CLR for UNDO

[Figure 4 from Franklin97]

3. Undo Phase: Details

What happens if system crashes during
UNDO ?

3. Undo Phase: Details

What happens if system crashes during UNDO ?

We do not UNDO again ! Instead, each CLR is a REDO record: we simply redo the undo

Physical v.s. Logical Logging

Why are redo records physical ?

Why are undo records logical ?

Physical v.s. Logical Logging

Why are redo records physical ?

- Simplicity: replaying history is easy, and idempotent

Why are undo records logical ?

- Required for transaction rollback: this not “undoing history”, but selective undo