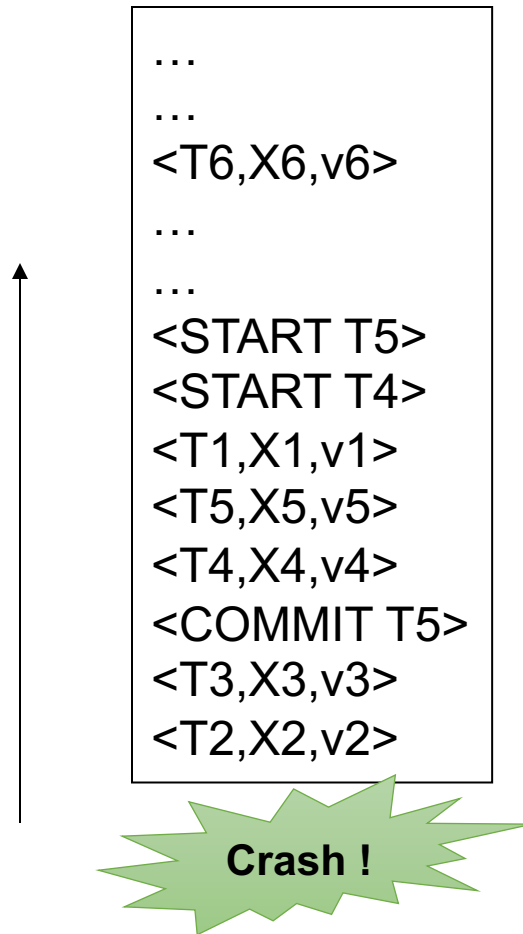


Database System Internals

Transactions: Recovery (part 2)

Paul G. Allen School of Computer Science and Engineering
University of Washington, Seattle

Recovery with Undo Log

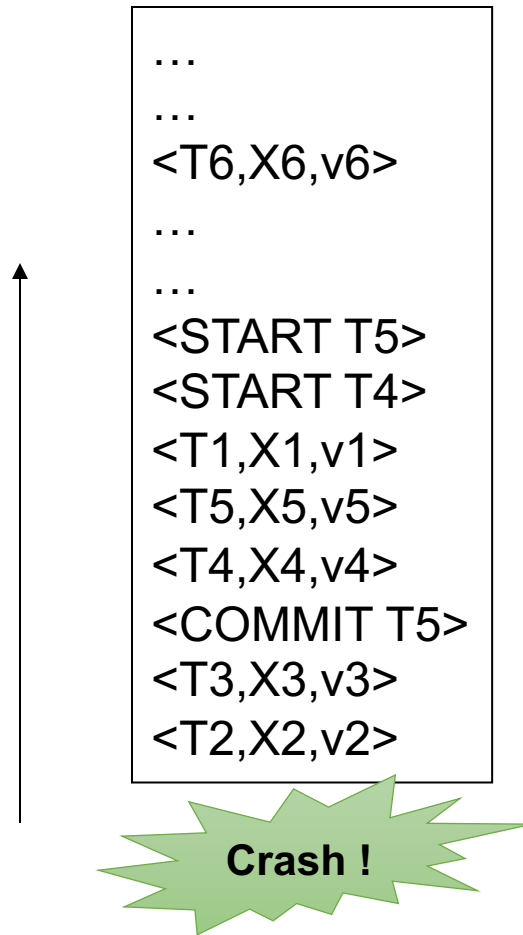


Question 1: Which updates are undone ?

Question 2:
How far back do we need to read in the log ?

Question 3:
What happens if second crash during recovery?

Recovery with Undo Log

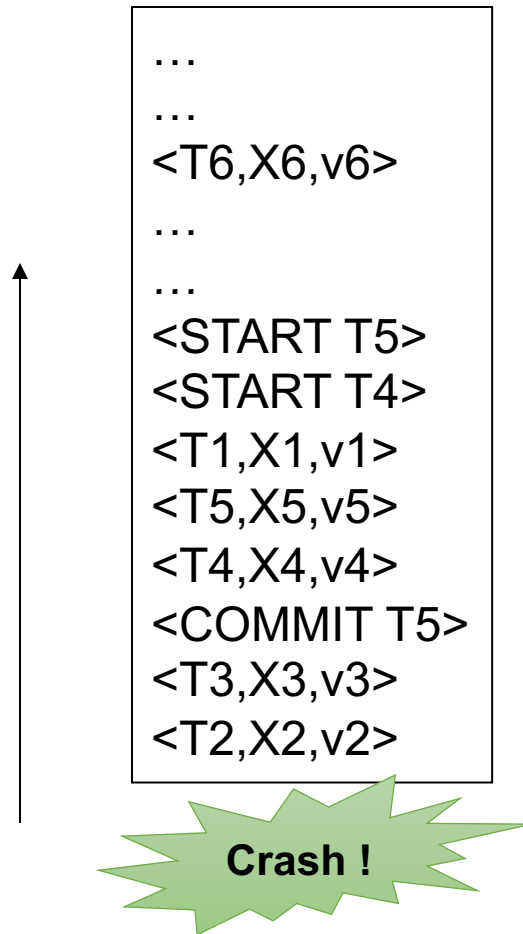


Question 1: Which updates are undone ?

Question 2:
How far back do we need to read in the log ?
To the beginning.

Question 3:
What happens if second crash during recovery?

Recovery with Undo Log



Question 1: Which updates are undone ?

Question 2:

How far back do we need to read in the log ?

To the beginning.

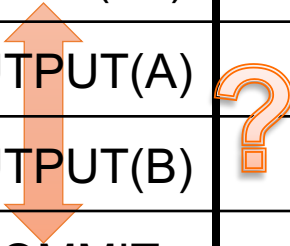
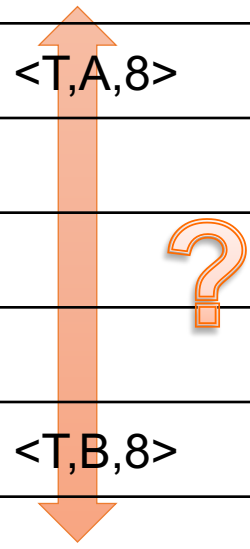
Question 3:

What happens if second crash during recovery?

No problem! Log records are idempotent. Can reapply.

When must we force log pages to disk ?

Act				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>

FORCE

RULES: log entry *before* OUTPUT *before* COMMIT

Undo-Logging Rules

U1: If T modifies X, then $\langle T, X, v \rangle$ must be written to disk before $\text{OUTPUT}(X)$

U2: If T commits, then $\text{OUTPUT}(X)$ must be written to disk before $\langle \text{COMMIT } T \rangle$

- Hence: OUTPUTs are done early, before the transaction commits



FORCE

Checkpointing

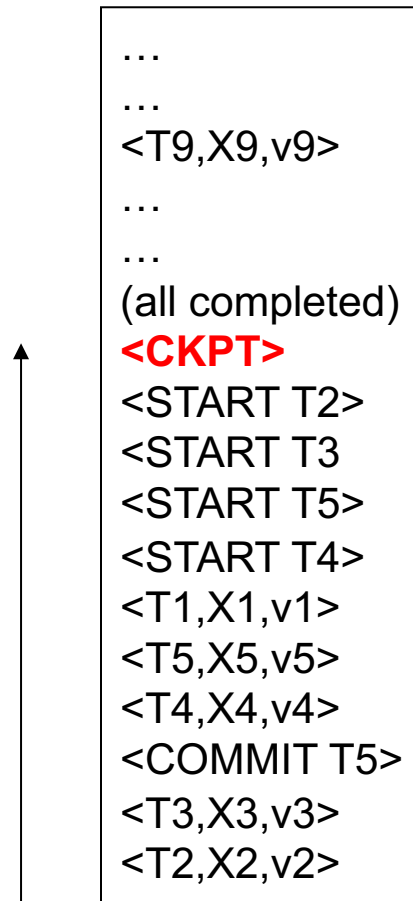
Checkpoint the database periodically

During a checkpoint:

- Stop accepting new transactions
- Wait until all current transactions complete
- Flush log to disk
- Write a **<CKPT>** log record, flush
- Resume transactions

Undo Recovery with Checkpointing

During recovery,
Can stop at first
<CKPT>



other transactions:
all have completed
no need to undo

transactions T2,T3,T4,T5

Nonquiescent Checkpointing

- Problem with checkpointing: database freezes during checkpoint
- Would like to checkpoint while database is operational
- Idea: nonquiescent checkpointing

Quiescent = being quiet, still, or at rest; inactive
Non-quiescent = allowing transactions to be active

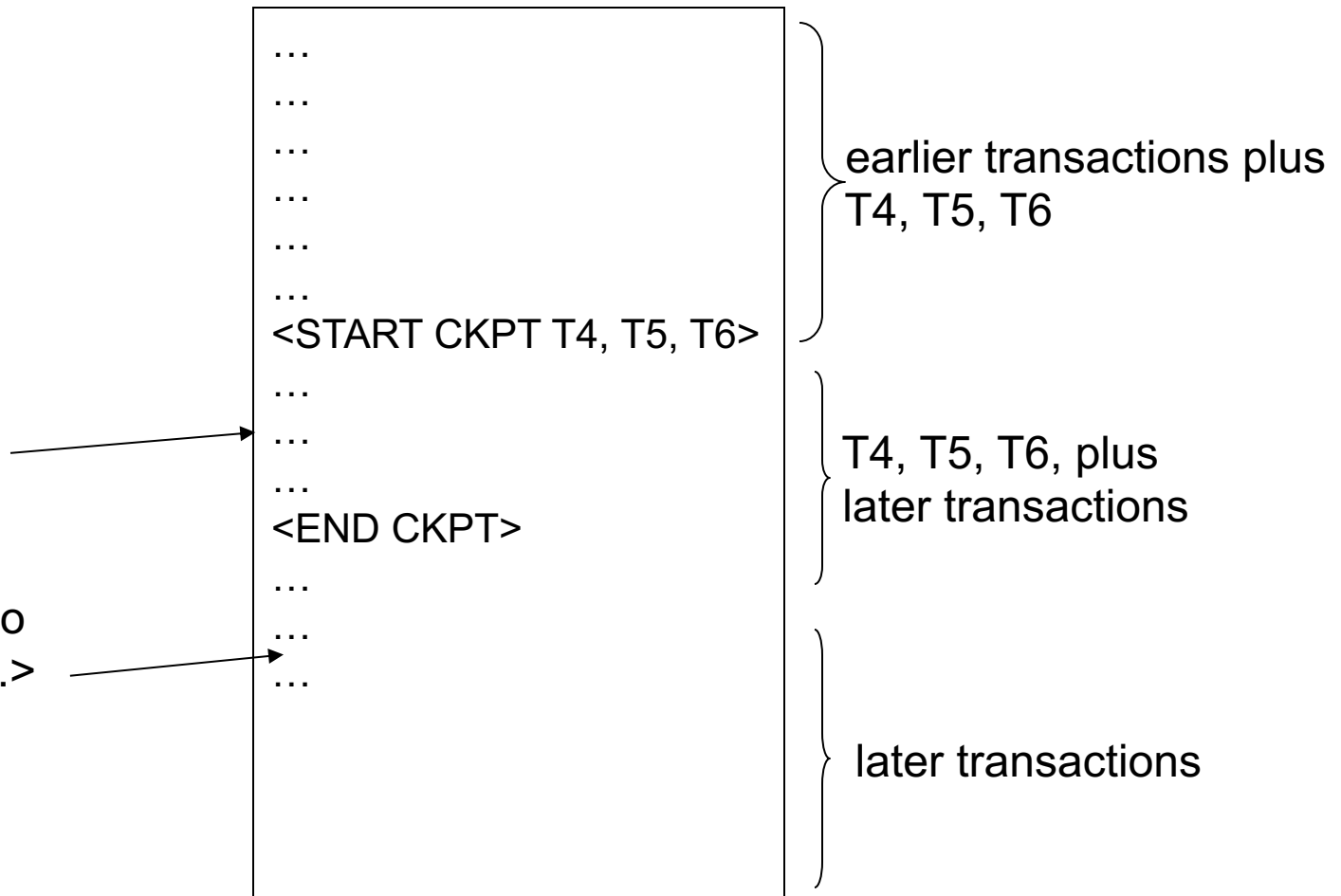
Nonquiescent Checkpointing

- Write a **<START CKPT(T1,...,Tk)>** where T_1, \dots, T_k are all active transactions.
- Flush log to disk
- Continue normal operation
- When all of T_1, \dots, T_k have completed, write **<END CKPT>**
- Flush log to disk

Undo with Nonquiescent Checkpointing

If we crash here:
Need to read
Back to start of
T4, T5, T6

If we crash here:
Need to read only to
<START CKPT T4..>



Undo with Nonquiescent Checkpointing

How far back
do we need
to read?

```
...  
...  
...  
<START CKPT T4, T5, T6>  
...  
...  
...  
...  
...  
<END CKPT>  
...  
...  
...  
...  
...
```

Undo with Nonquiescent Checkpointing

How far back
do we need
to read?

```
...  
...  
...  
<START CKPT T4, T5, T6>  
...  
...  
...  
...  
...  
<END CKPT>  
...  
...  
...  
...  
...
```

Here

Undo with Nonquiescent Checkpointing

How far back
now?

```
...  
...  
...  
<START CKPT T4, T5, T6>  
...  
...  
...  
...  
...  
<END CKPT>  
...  
...  
...  
<START CKPT T9, T10>  
...
```

Undo with Nonquiescent Checkpointing

How far back
now?

```
...  
...  
...  
<START CKPT T4, T5, T6>  
...  
...  
...  
...  
...  
<END CKPT>  
...  
...  
...  
<START CKPT T9, T10>  
...
```

Earliest of
<START T9>
<START T10>

Implementing ROLLBACK

- Recall: a transaction can end in COMMIT or ROLLBACK
- Idea: use the undo-log to implement ROLLBACK
- How ?
 - LSN = Log Sequence Number
 - Log entries for the same transaction are linked, using the LSN's
 - Read log in reverse, using LSN pointers

Implementing ROLLBACK

```
...  
<T9,X9,v9>  
...  
...  
(all completed)  
<CKPT>  
<START T2>  
<START T3  
<START T5>  
<START T4>  
<T1,X1,v1>  
<T5,X5,v5>  
<T2,X1,v2>  
<T4,X4,v4>  
<COMMIT T5>  
<T3,X3,v3>  
<T2,X2,v2>  
<T4 X4,v4>
```

LSN = Log Sequence Number

T2: ROLLBACK

Implementing ROLLBACK

```
...  
<T9,X9,v9>  
...  
...  
(all completed)  
<CKPT>  
<START T2>  
<START T3  
<START T5>  
<START T4>  
<T1,X1,v1>  
<T5,X5,v5>  
<T2,X1,v2>  
<T4,X4,v4>  
<COMMIT T5>  
<T3,X3,v3>  
<T2,X2,v2>  
<T4 X4,v4>
```

LSN = Log Sequence Number

T2: ROLLBACK

T2 has
pointer to
its last LSN



Implementing ROLLBACK

...
<T9,X9,v9>
...
...
(all completed)
<CKPT>
<START T2>
<START T3
<START T5>
<START T4>
<T1,X1,v1>
<T5,X5,v5>
<T2,X1,v2>
<T4,X4,v4>
<COMMIT T5>
<T3,X3,v3>
<T2,X2,v2>
<T4 X4,v4>

LSN = Log Sequence Number

T2: ROLLBACK

T2 has
pointer to
its last LSN

Implementing ROLLBACK

...
<T9,X9,v9>
...
...
(all completed)
<CKPT>
<START T2>
<START T3
<START T5>
<START T4>
<T1,X1,v1>
<T5,X5,v5>
<T2,X1,v2>
<T4,X4,v4>
<COMMIT T5>
<T3,X3,v3>
<T2,X2,v2>
<T4 X4,v4>

LSN = Log Sequence Number

T2: ROLLBACK

T2 has
pointer to
its last LSN

"REDO" Log

NO FORCE and NO STEAL

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

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: lost update

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 ?

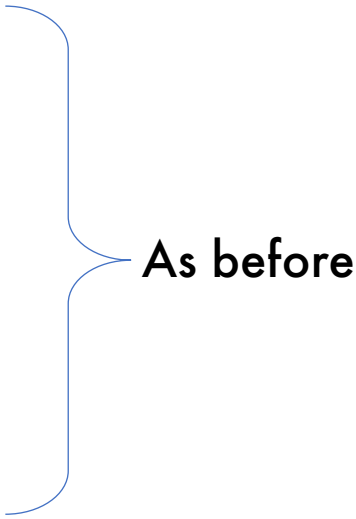
No: that's OK.

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




Redo Logging


Log records

- **<START T>**
 - transaction T has begun
 - **<COMMIT T>**
 - T has committed
 - **<ABORT T>**
 - T has aborted
 - **<T,X,v>**
 - T has updated element X, and its new value is v
 - *Idempotent, physical* log records
- 
- As before

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 ?

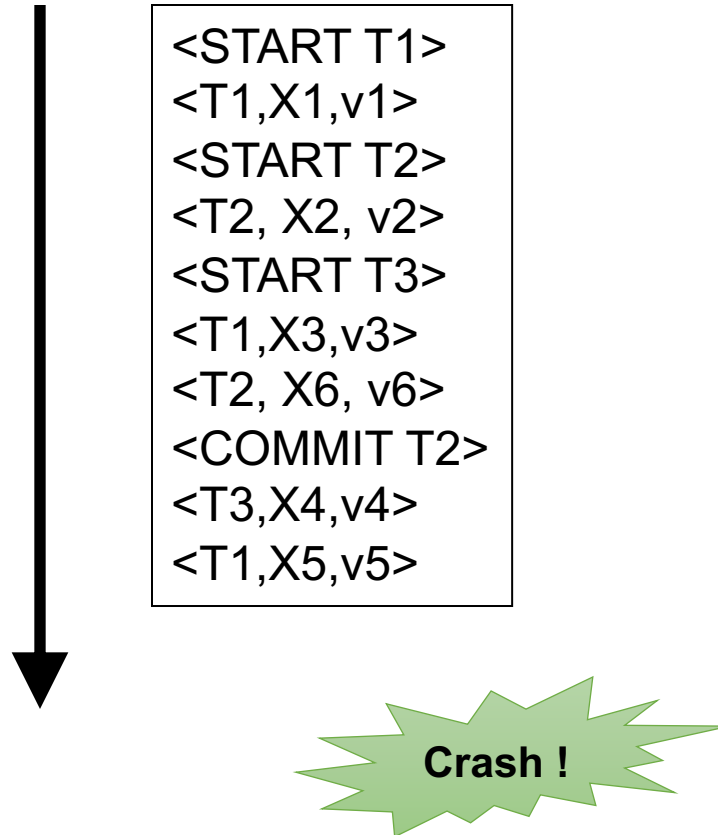
We **REDO** by setting A=16 and B=16

Recovery with Redo Log

After system's crash, run recovery manager

- Step 1. Decide for each transaction T whether it is committed or not
 - <START T>....<COMMIT T>.... = yes
 - <START T>....<ABORT T>..... = no
 - <START T>..... = no
- Step 2. Read log from the beginning, redo all updates of committed transactions

Recovery with Redo Log



Recovery with Redo Log



```
<START T1>  
<T1,X1,v1>  
<START T2>  
<T2, X2, v2>  
<START T3>  
<T1,X3,v3>  
<T2, X6, v6>  
<COMMIT T2>  
<T3,X4,v4>  
<T1,X5,v5>
```

Step 1: find committed TXNs



Recovery with Redo Log



```
<START T1>  
<T1,X1,v1>  
<START T2>  
<T2, X2, v2>  
<START T3>  
<T1,X3,v3>  
<T2, X6, v6>  
<COMMIT T2>  
<T3,X4,v4>  
<T1,X5,v5>
```

Step 1: find committed TXNs

T2 (only)



Recovery with Redo Log



```
<START T1>  
<T1,X1,v1>  
<START T2>  
<T2, X2, v2>  
<START T3>  
<T1,X3,v3>  
<T2, X6, v6>  
<COMMIT T2>  
<T3,X4,v4>  
<T1,X5,v5>
```

Step 1: find committed TXNs

T2 (only)

Step 2: redo their actions



Recovery with Redo Log



```
<START T1>  
<T1,X1,v1>  
<START T2>  
<T2, X2, v2>  
<START T3>  
<T1,X3,v3>  
<T2, X6, v6>  
<COMMIT T2>  
<T3,X4,v4>  
<T1,X5,v5>
```

Step 1: find committed TXNs

T2 (only)

Step 2: redo their actions

X2 = v2



Recovery with Redo Log



```
<START T1>  
<T1,X1,v1>  
<START T2>  
<T2, X2, v2>  
<START T3>  
<T1,X3,v3>  
<T2, X6, v6>  
<COMMIT T2>  
<T3,X4,v4>  
<T1,X5,v5>
```

Step 1: find committed TXNs

T2 (only)

Step 2: redo their actions

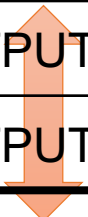
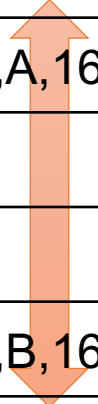
X2 = v2

X6 = v6



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

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

NO-STEAL

Redo-Logging Rules

R1: If T modifies X, then both $\langle T, X, v \rangle$ and $\langle \text{COMMIT } T \rangle$ must be written to disk before $\text{OUTPUT}(X)$

- Hence: OUTPUTs are done late

NO-STEAL

Nonquiescent Checkpointing

- Write $\langle \text{START CKPT}(T_1, \dots, T_k) \rangle$
where T_1, \dots, T_k are all active txn's
- Begin flush to disk all **blocks of committed transactions** (*dirty blocks*)
- Meantime, continue normal operation
- When **all blocks have been written**, write $\langle \text{END CKPT} \rangle$

END CKPT different from Undo!
Does not mean that T_1, \dots, T_k are complete

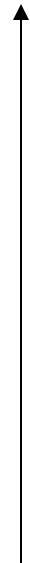
Nonquiescent Checkpointing

```
...  
<START T1>  
...  
<COMMIT T1>  
...  
<START T4>  
...  
<START CKPT T4, T5, T6>  
...  
...  
<COMMIT T4>  
...  
<END CKPT>  
...  
...  
<COMMIT T5>  
...
```

Nonquiescent Checkpointing

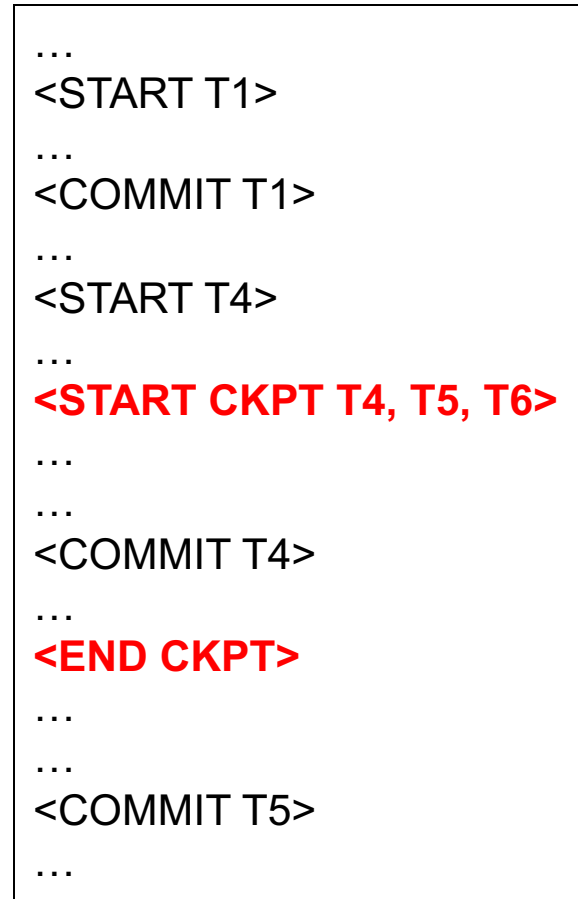
Step 1: look for
The last
<END CKPT> and it's
<START CKPT>

```
...  
<START T1>  
...  
<COMMIT T1>  
...  
<START T4>  
...  
<START CKPT T4, T5, T6>  
...  
...  
<COMMIT T4>  
...  
<END CKPT>  
...  
...  
<COMMIT T5>  
...
```



Nonquiescent Checkpointing

Step 1: look for
The last
<END CKPT> and it's
<START CKPT>



Step 2: redo...
...from where?

Nonquiescent Checkpointing

Step 1: look for
The last
<END CKPT> and it's
<START CKPT>

We know
all pages
of T1 are
on disk

```
...  
<START T1>  
...  
<COMMIT T1>  
...  
<START T4>  
...  
<START CKPT T4, T5, T6>  
...  
...  
<COMMIT T4>  
...  
<END CKPT>  
...  
...  
<COMMIT T5>  
...
```

Step 2: redo...
...from where?

Nonquiescent Checkpointing

Step 1: look for
The last
<END CKPT> and it's
<START CKPT>

We know
all pages
of T1 are
on disk

```
...  
<START T1>  
...  
<COMMIT T1>  
...  
<START T4>  
...  
<START CKPT T4, T5, T6>  
...  
...  
<COMMIT T4>  
...  
<END CKPT>  
...  
...  
<COMMIT T5>  
...
```

Step 2: redo...
...from where?

What about
pages of T4?

Nonquiescent Checkpointing

Step 1: look for
The last
<END CKPT> and it's
<START CKPT>

We know
all pages
of T1 are
on disk

```
...  
<START T1>  
...  
<COMMIT T1>  
...  
<START T4>  
...  
<START CKPT T4, T5, T6>  
...  
...  
<COMMIT T4>  
...  
<END CKPT>  
...  
...  
<COMMIT T5>  
...
```

Step 2: redo...
...from where?

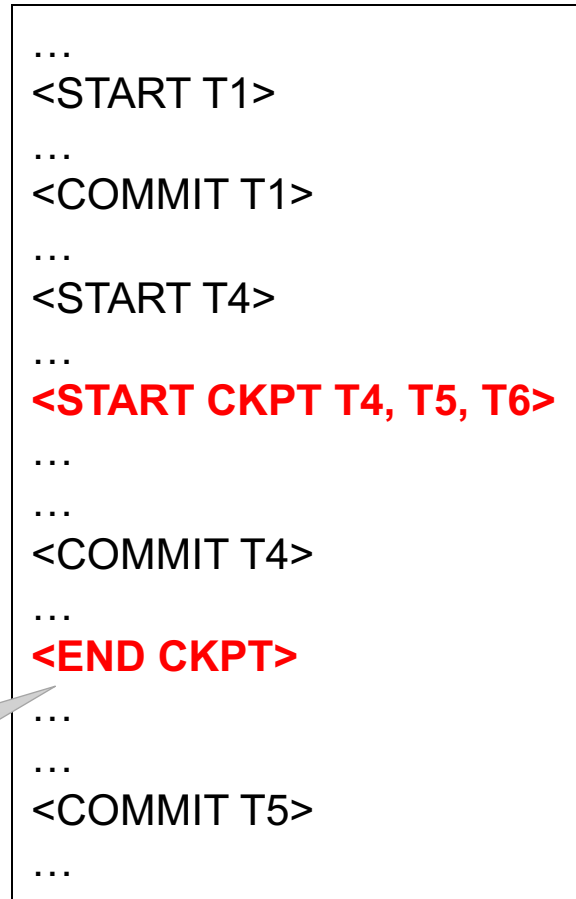
What about
pages of T4?

NO!
They may
be in BP

Nonquiescent Checkpointing

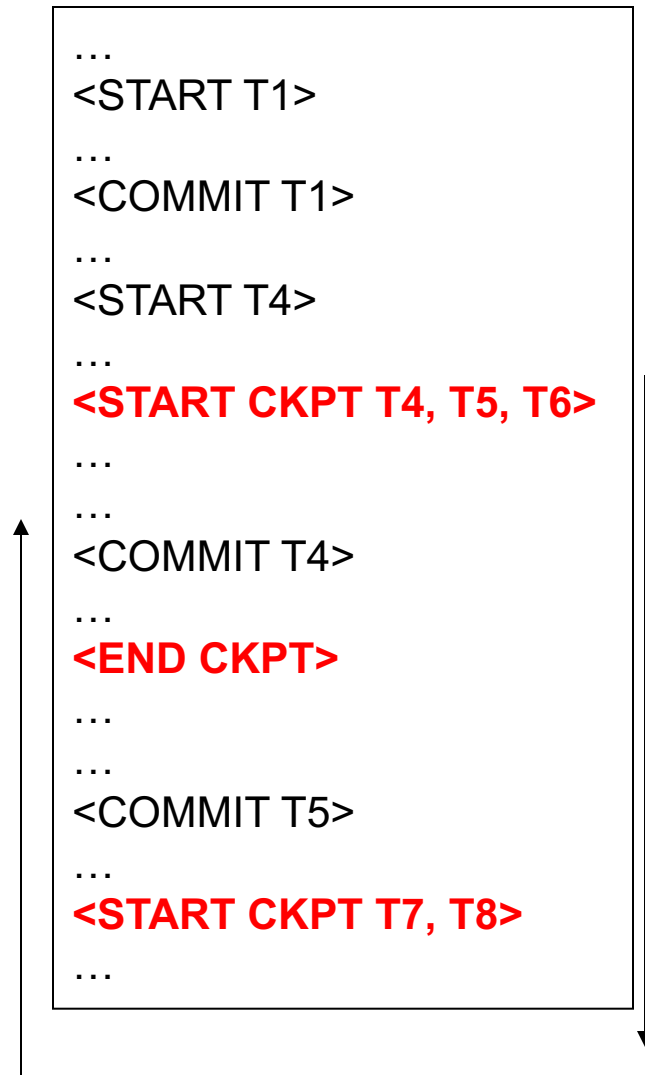
Step 1: look for
The last
<END CKPT> and it's
<START CKPT>

We know
all pages
of T1 are
on disk

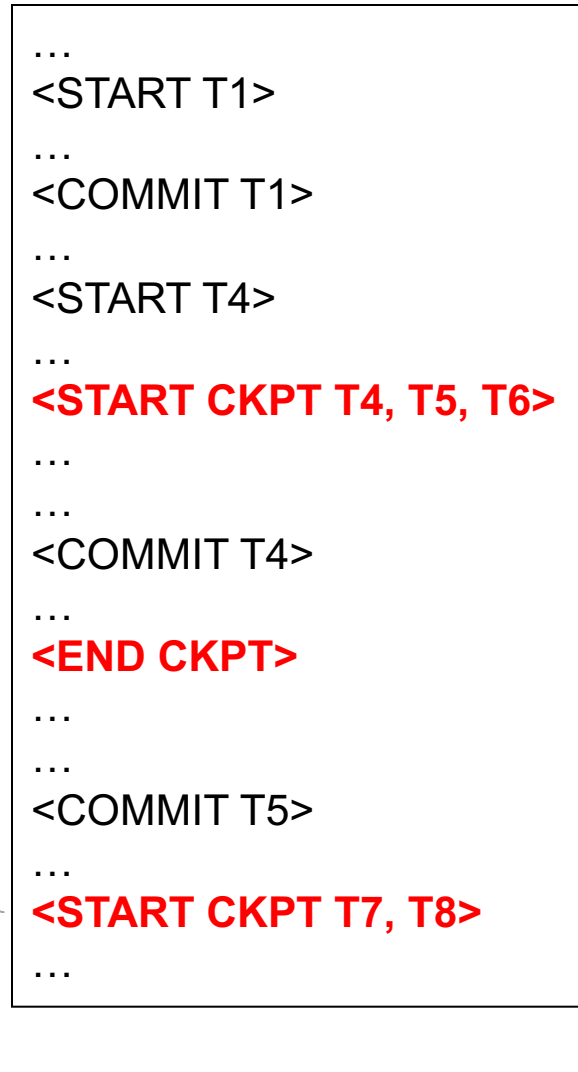


Step 2: redo
from the
earliest
start of
T4, T5, T6
ignoring
transactions
committed
earlier

Nonquiescent Checkpointing



Nonquiescent Checkpointing



Can we Use this?

Nonquiescent Checkpointing

```
...  
<START T1>  
...  
<COMMIT T1>  
...  
<START T4>  
...  
<START CKPT T4, T5, T6>  
...  
...  
<COMMIT T4>  
...  
<END CKPT>  
...  
...  
<COMMIT T5>  
...  
<START CKPT T7, T8>  
...
```

Can we Use this?

Cannot use!

Comparison Undo/Redo Logging

■ Undo:

- OUTPUT must be done early
- Later reads must reload page; Inefficient

Steal/Force

■ Redo

- OUTPUT must be done late
- Buffer manager cannot evict page; Inflexible

No-Steal/No-Force

■ Undo/redo logging (next)

- Both efficient and flexible

Steal/No-Force

“UNDO-REDO” Log

NO FORCE and STEAL

Undo/Redo Logging

Log records, only one change

- $\langle T, X, u, v \rangle =$ T has updated element X, its old value was u, and its new value is v

Undo/Redo-Logging Rule

UR1: If T modifies X, then $\langle T, X, u, v \rangle$ must be written to disk before $\text{OUTPUT}(X)$

Note: we are free to OUTPUT early or late relative to $\langle \text{COMMIT } T \rangle$

Action	T	Mem A	Mem B	Disk A	Disk B	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,8,16>
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,8,16>
OUTPUT(A)	16	16	16	16	8	
						<COMMIT T>
OUTPUT(B)	16	16	16	16	16	

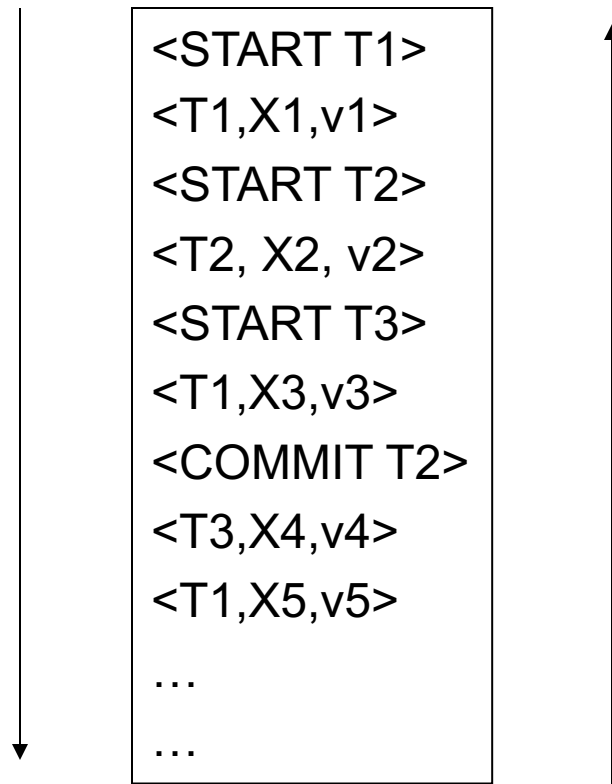
Can OUTPUT whenever we want: before/after COMMIT

Recovery with Undo/Redo Log

After system's crash, run recovery manager

- **Redo all committed transaction, top-down**
- **Undo all uncommitted transactions, bottom-up**

Recovery with Undo/Redo Log



ARIES

- ARIES pieces together several techniques into a comprehensive algorithm
- Developed at IBM Almaden, by Mohan
- IBM botched the patent, so everyone uses it now
- Several variations, e.g. for distributed transactions

Two basic types of log records for update operations

▪ Physical log records

- Position on a particular page where update occurred
- Both before and after image for undo/redo logs
- Benefits: Idempotent & updates are fast to redo/undo

▪ Logical log records

- Record only high-level information about the operation
- Benefit: Smaller log
- BUT difficult to implement because crashes can occur in the middle of an operation

ARIES Recovery Manager

Log entries:

- **<START T>** – when T begins
- **Update: <T,X,u,v>**
 - T updates X, old value=u, new value=v
 - Logical description of the change
- **<COMMIT T>** or **<ABORT T>** then **<END>**
- **<CLR>** – 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 w.r.t commits

LSN = Log Sequence Number

- **LSN** = identifier of a log entry
 - Log entries belonging to the same TXN are linked with extra entry for previous LSN

- Each page contains a **pageLSN**:
 - LSN of log record for latest update to that page

ARIES Data Structures

▪ Active Transactions Table

- Lists all active TXN's
- For each TXN: **lastLSN** = its most recent update LSN

▪ 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

Data Structures

$W_{T100}(P7)$

$W_{T200}(P5)$

$W_{T200}(P6)$

$W_{T100}(P5)$

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