

Recovery

R&G Chapter 18

(slides adapted from content by J.Gehrke, J.Shanmugasundaram, and/or C.Koch)

Announcements

- Project 3 to be posted by Tuesday. Due mid-April
- Homework 6 to be posted Tonight. Due Monday.
- Project 2 Extra Credit: Schedule a 20 min meeting
 - <http://doodle.com/okennedy>

Recap A.C.I.D.

- **Atomicity**: All actions in a transaction happen, or none happen.
- **Consistency**: If the transaction maintains consistency, and the DB starts consistent, then the database ends consistent.
- **Isolation**: The execution of one transaction is isolated from all other transactions.
- **Durability**: If a transaction commits, its effects persist.

Recap: Write-Ahead Logging

- The write-ahead logging protocol
 - Force the log record for an update before the corresponding data page is written.
 - Guarantees Atomicity (REDO).
- Force all log records for a transaction before the transaction commits.
- Guarantees Durability (UNDO).

Recap: The Big Picture

The Log

Log Records

Prev LSN

XID

Record Type

Page ID

Length

Offset

Before-Image

After-Image

The DB (Disk)

Data Pages
(each with a PageLSN)

Master Record

The DB (Ram)

Transaction Table

XID

LastLSN

Status

Dirty Page Table

RecLSN

Flushed LSN

Checkpointing

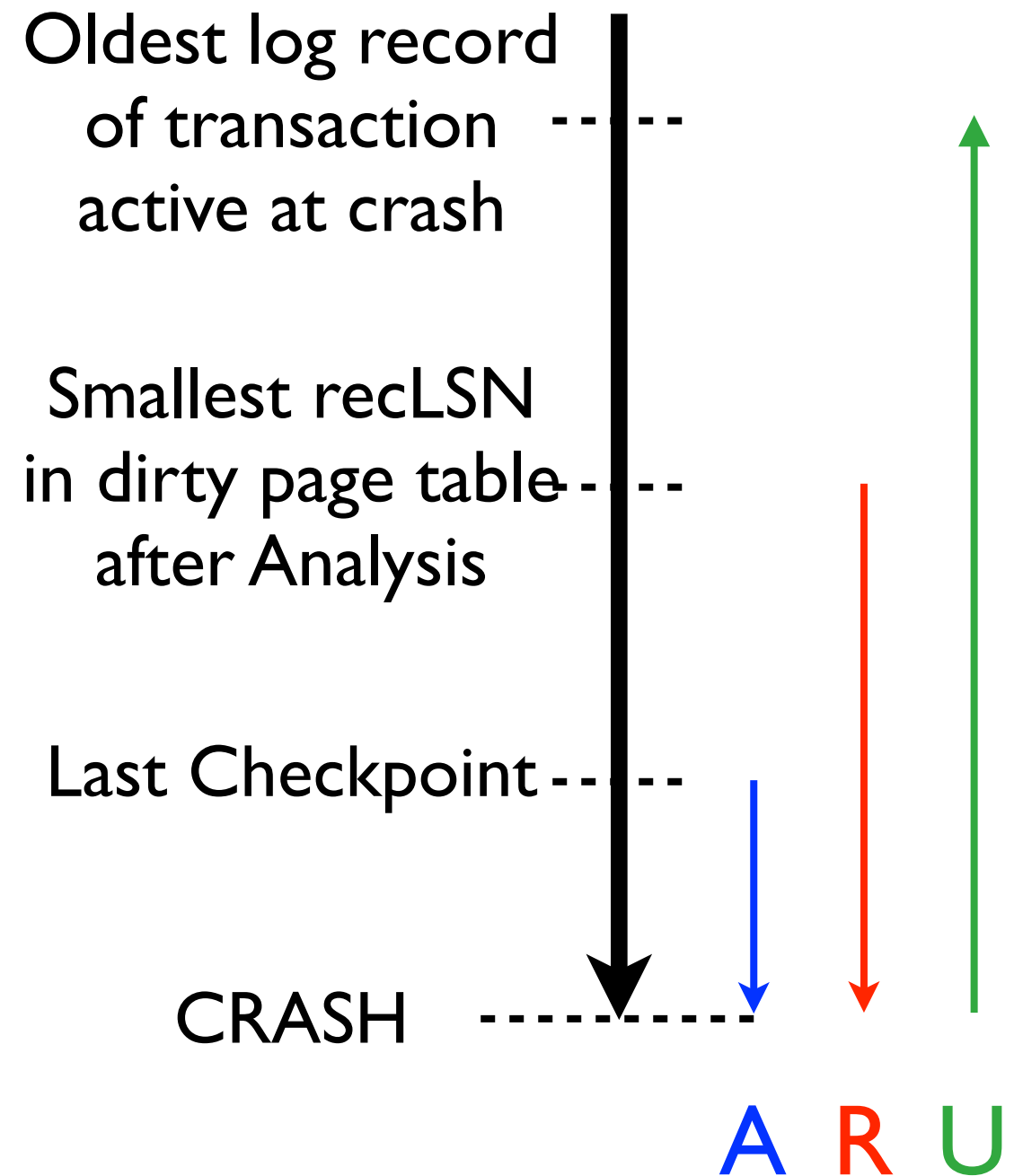
- Log grows indefinitely; Need a way to speed up the recovery process.
- Periodically, the DBMS creates a checkpoint.
- **begin_checkpoint** record indicates when the checkpoint began.
- **end_checkpoint** record contains the current transaction table and the dirty page table.

Checkpointing

- This is a ‘fuzzy’ checkpoint
- Other transactions continue to run, so the tables are only accurate as of the **begin_checkpoint** record.
- Checkpoints don’t force data to disk
 - Recovery process needs to go to oldest unwritten change as of the checkpoint.
- Store the LSN of the most recent checkpoint record in a safe place on disk (the master record).

Crash Recovery

- Start from checkpoint stored in master record.
- **Analysis**: Figure out which transactions committed since checkpoint.
- **Redo** all actions
- **Undo** effects of failed transactions.



Recovery: Analysis

- Reconstruct state at checkpoint
 - (via the end_checkpoint record)
- Scan log forward from checkpoint
 - End record: Remove transaction from transaction table
 - Other records: Add transaction to transaction table, set lastLSN=LSN, change transaction status on commit.
 - Update record: If P not in the dirty page table, add P to the DPT, set its recLSN=LSN

Recovery: REDO

- Repeat History to reconstruct state at crash
- Reapply **all** updates, even those of aborted transactions. Redo CLR as they appear.
- REDO actions by:
 - Reapplying logged action.
 - Set pageLSN to LSN.
 - Don't need to log anything (already logged)

Recovery: REDO

- Scan forward from the log record containing the smallest recLSN in the DPT.
- For each CLR or update log record LSN, REDO the action unless:
 - Affected page not in the DPT, or
 - Affected page is in DPT but has $\text{recLSN} > \text{LSN}$, or
 - $\text{pageLSN (in DB)} \geq \text{LSN}$.

Recovery: UNDO

- Define a priority queue:
 - $ToUndo := \{\text{all lastLSNs of undone xacts}\}$
- Choose largest LSN in ToUndo
 - If this LSN is the last CLR for an xact
 - Write an END record
 - If this LSN is a non-terminal CLR
 - Add the CLR's undonextLSN to ToUNDO

Recovery: UNDO

- Choose largest LSN in ToUndo
 - ...
- Otherwise, this LSN is an update.
 - Undo the update, write a CLR
 - Add the update's prevLSN to ToUndo if it exists.
- Repeat until ToUndo is empty

Recovery Example

<u>LSN</u>	<u>Log</u>
00	begin_checkpoint
05	end_checkpoint
10	update:T1 writes P5
20	update:T2 writes P3
30	T1 Abort
40	CLR Undo T1 LSN 10
45	T1 End
50	update:T3 writes P1
60	update:T2 writes P5

The diagram illustrates the recovery process after a crash. A point on the right is labeled 'PrevLSNs'. From this point, three dashed arrows point to log entries: one to 'update:T1 writes P5' (LSN 10), one to 'T1 Abort' (LSN 30), and one to 'CLR Undo T1 LSN 10' (LSN 40). Solid curved arrows show the sequence of operations: from 'update:T2 writes P3' (LSN 20) to 'T1 Abort' (LSN 30), from 'CLR Undo T1 LSN 10' (LSN 40) to 'T1 End' (LSN 45), and from 'update:T2 writes P5' (LSN 60) to 'T1 End' (LSN 45).

CRASH! Restart!

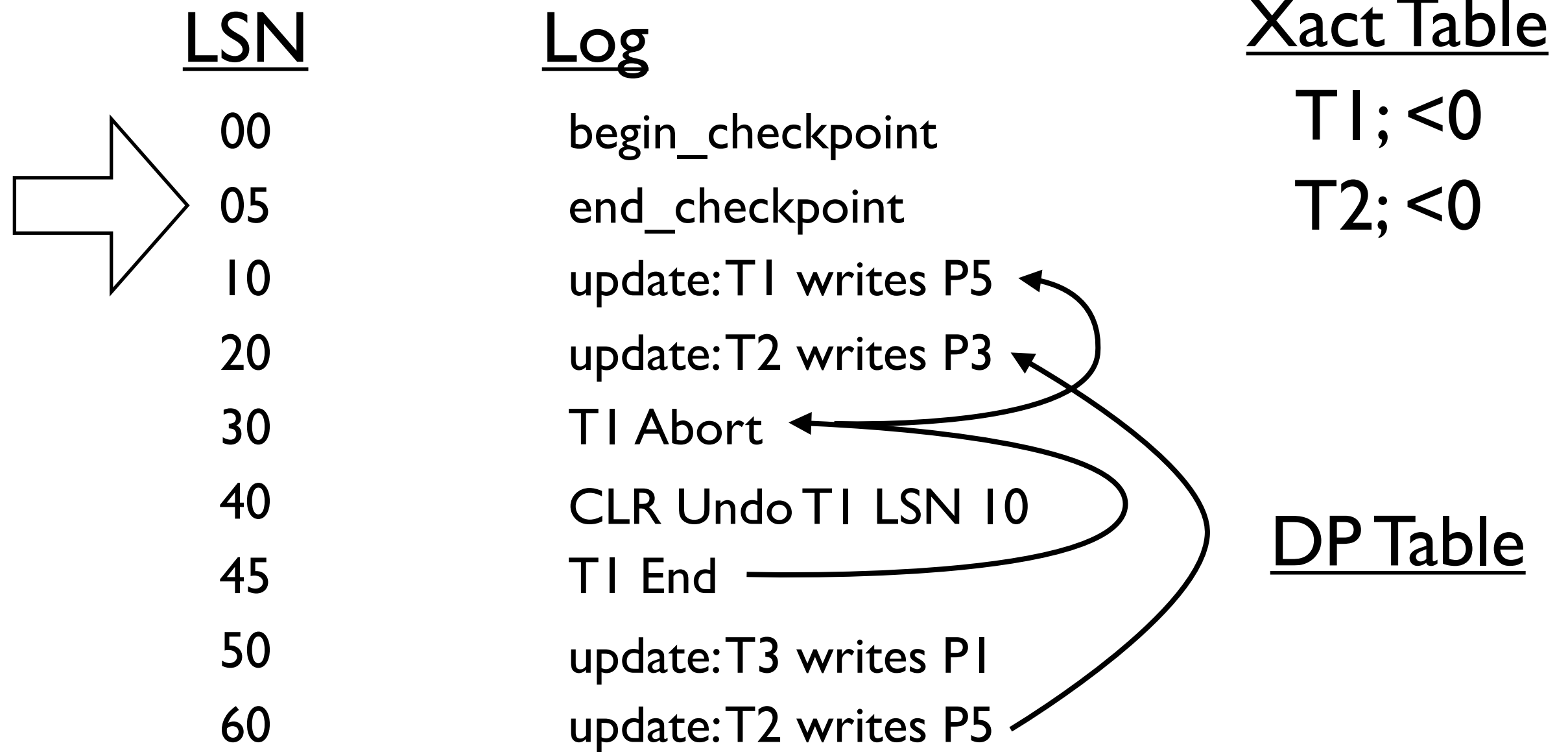
Analysis

<u>LSN</u>	<u>Log</u>	<u>Xact Table</u>
00	begin_checkpoint	T1; <0
05	end_checkpoint	T2; <0
10	update:T1 writes P5	
20	update:T2 writes P3	
30	T1 Abort	
40	CLR Undo T1 LSN 10	
45	T1 End	
50	update:T3 writes P1	
60	update:T2 writes P5	

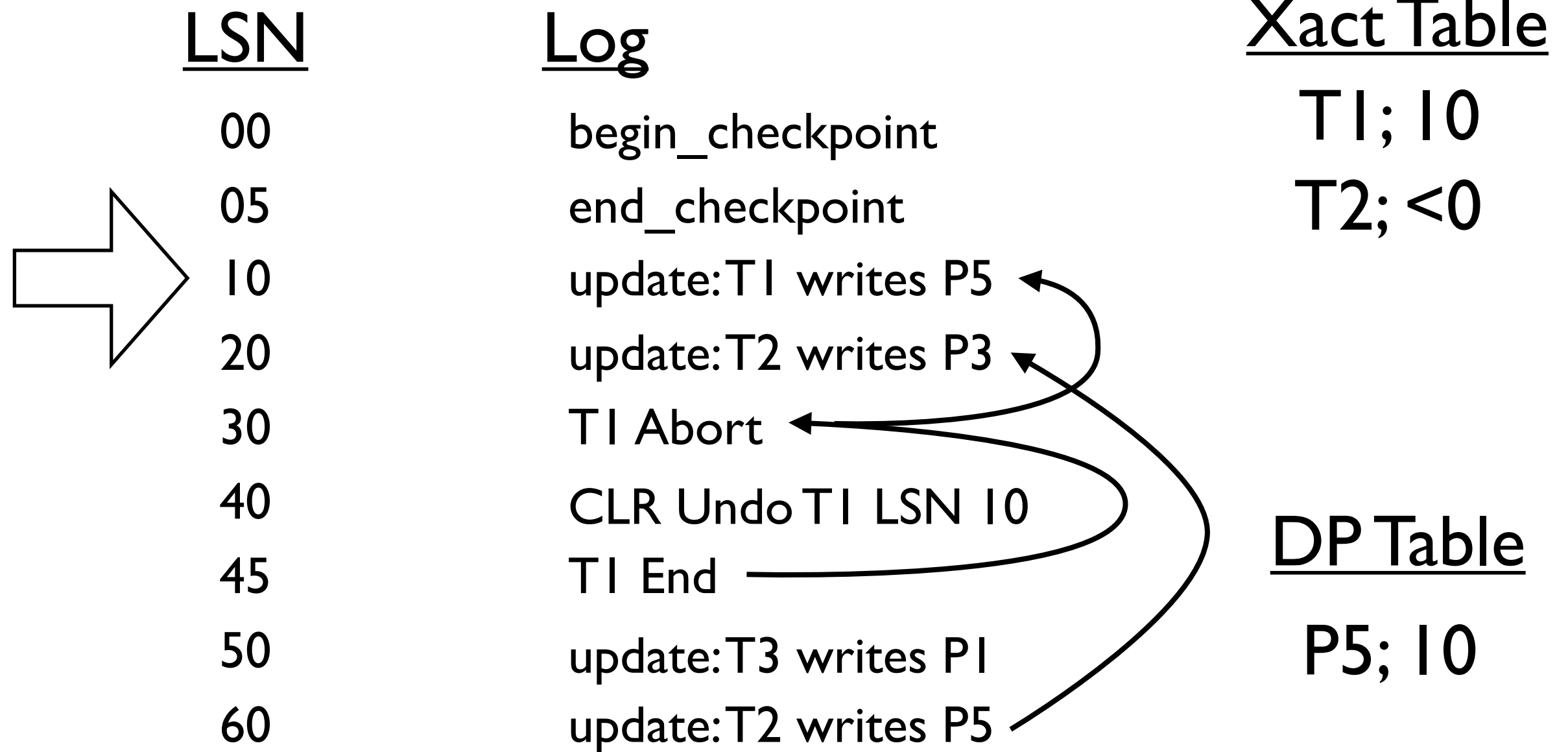
DP Table

```
graph LR; L10[LSN 10: update:T1 writes P5] --> DP[DP Table]; L20[LSN 20: update:T2 writes P3] --> DP; L60[LSN 60: update:T2 writes P5] --> DP; L45[LSN 45: T1 End] --> DP;
```

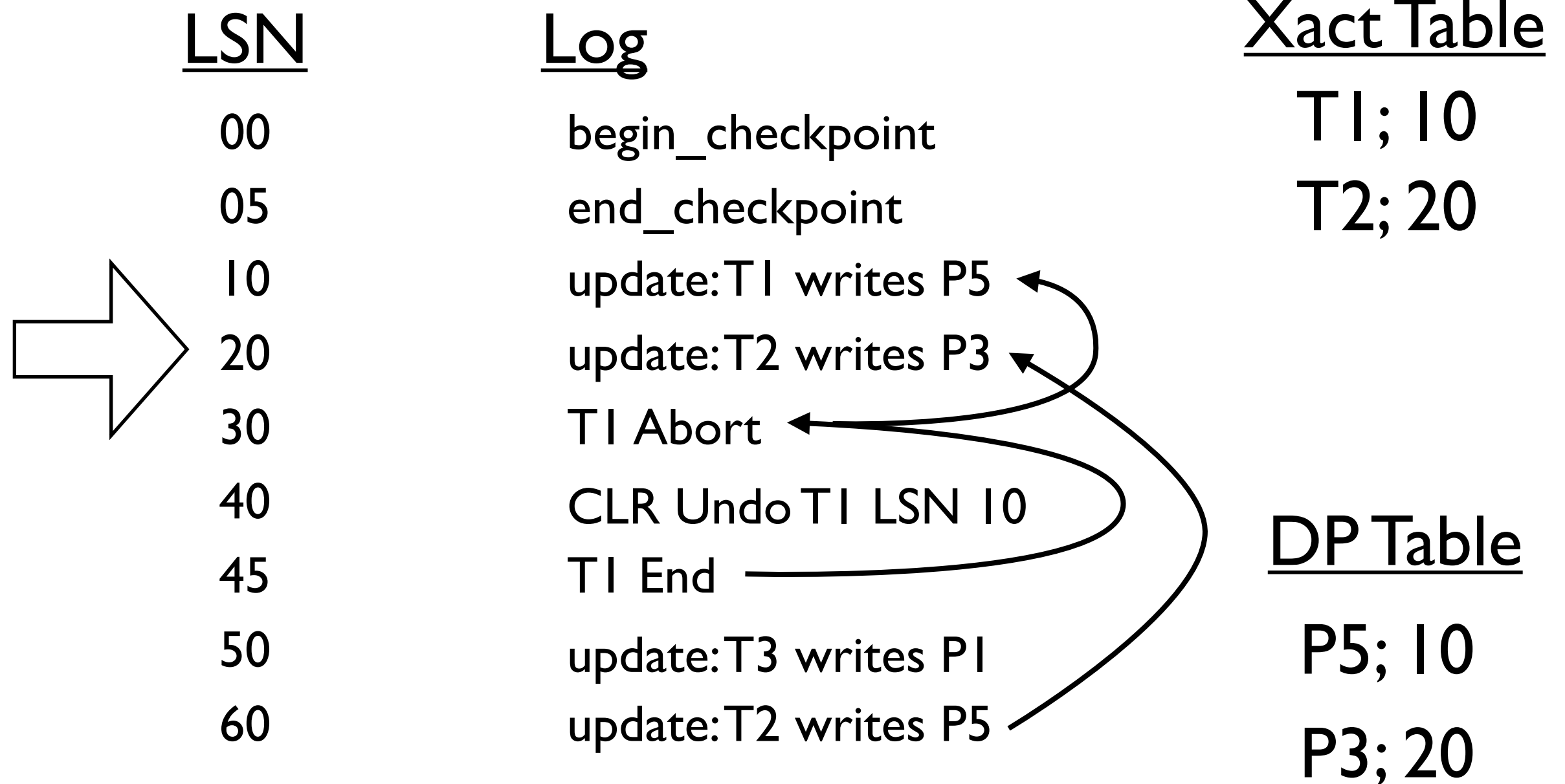
Analysis



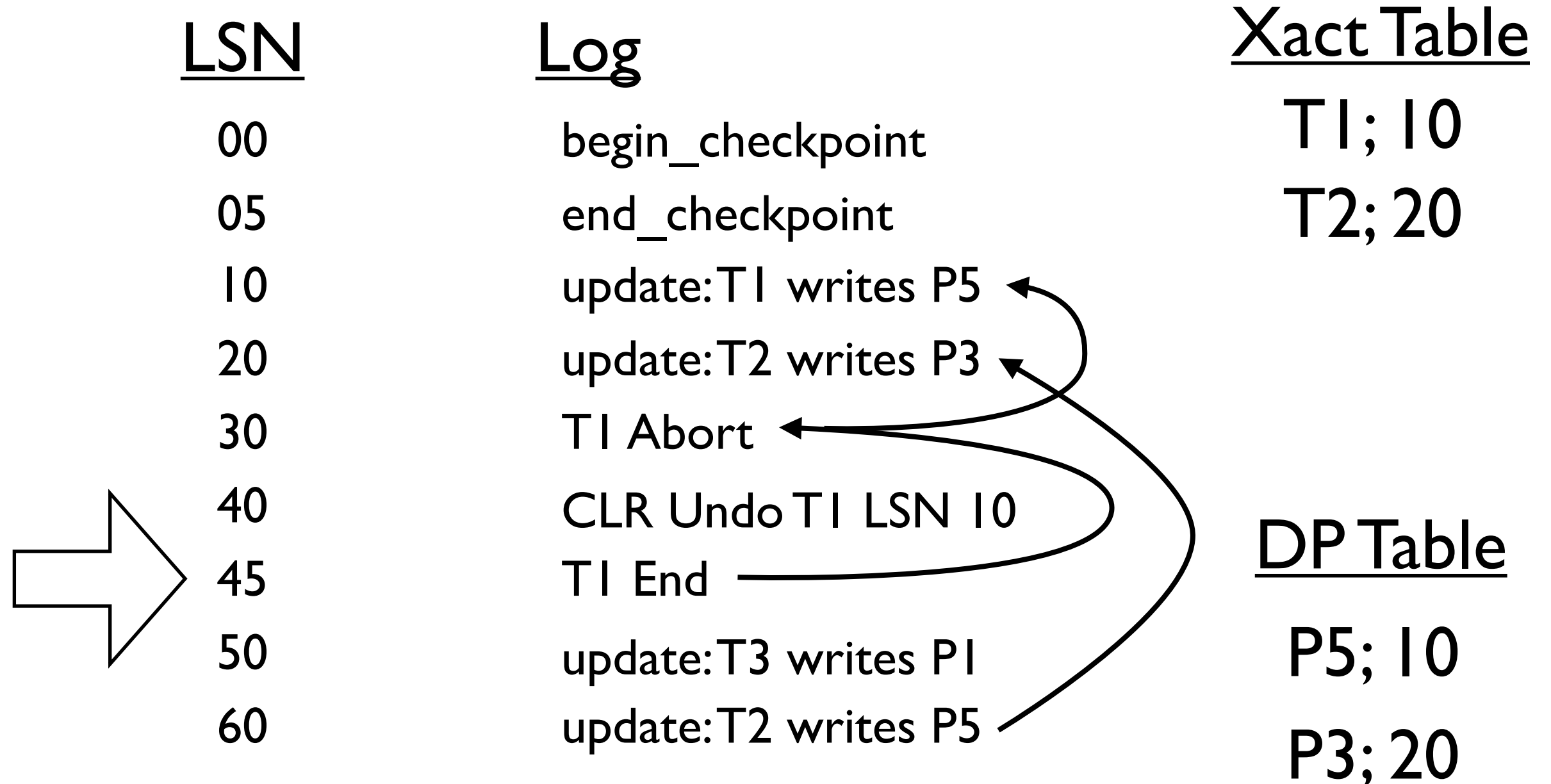
Analysis



Analysis



Analysis



Analysis

<u>LSN</u>	<u>Log</u>	<u>Xact Table</u>
00	begin_checkpoint	
05	end_checkpoint	T2; 20
10	update:T1 writes P5	
20	update:T2 writes P3	
30	T1 Abort	
40	CLR Undo T1 LSN 10	
45	T1 End	
50	update:T3 writes P1	
60	update:T2 writes P5	

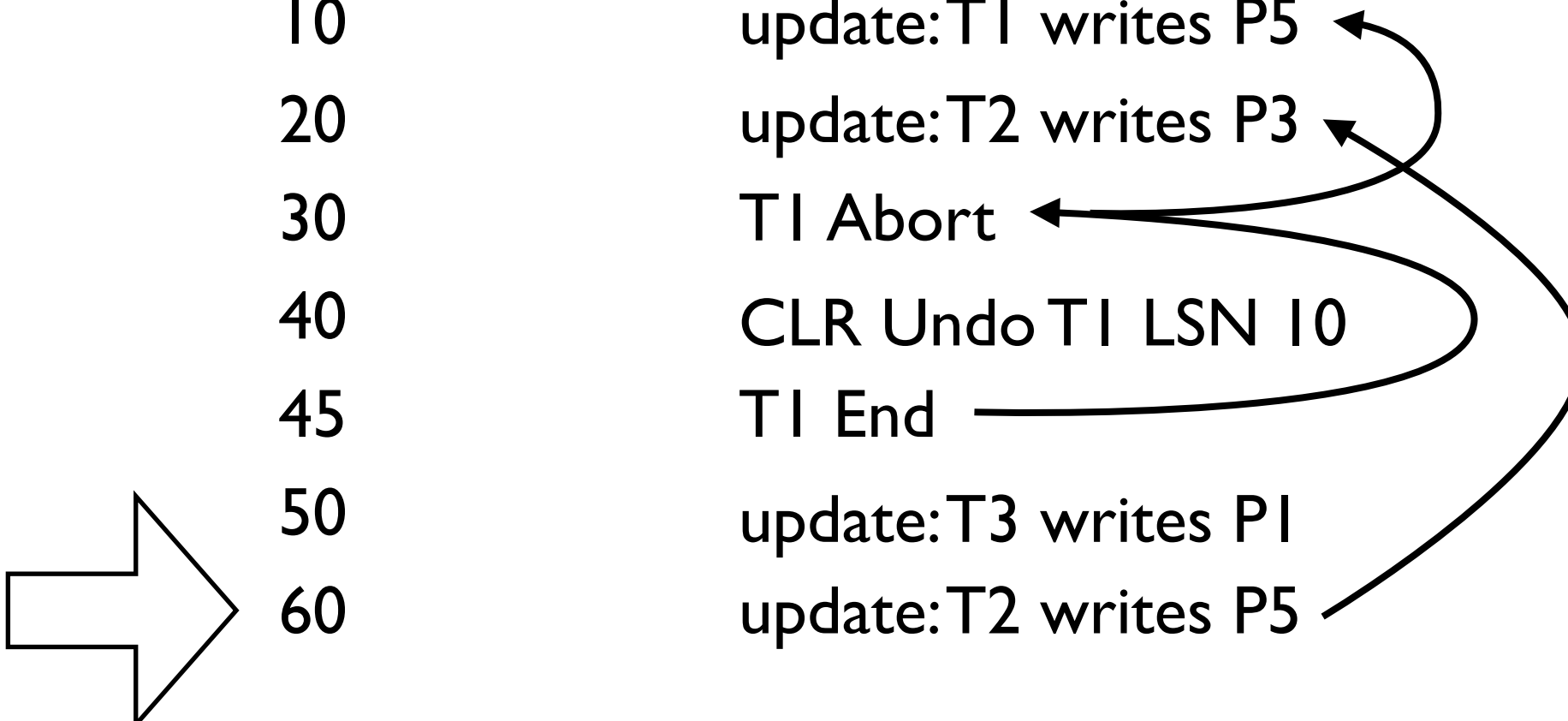
<u>DP Table</u>
P5; 10
P3; 20

Analysis

<u>LSN</u>	<u>Log</u>	<u>Xact Table</u>	<u>DP Table</u>
00	begin_checkpoint		
05	end_checkpoint	T2; 20	
10	update:T1 writes P5	T3; 50	
20	update:T2 writes P3		
30	T1 Abort		
40	CLR Undo T1 LSN 10		
45	T1 End		
50	update:T3 writes P1		P5; 10
60	update:T2 writes P5		P3; 20
			P1; 50

Analysis

<u>LSN</u>	<u>Log</u>	<u>Xact Table</u>	<u>DP Table</u>
00	begin_checkpoint		
05	end_checkpoint	T2; 60	
10	update:T1 writes P5	T3; 50	
20	update:T2 writes P3		
30	T1 Abort		
40	CLR Undo T1 LSN 10		
45	T1 End		
50	update:T3 writes P1		P5; 10
60	update:T2 writes P5		P3; 20
			P1; 50



Redo

LSN

Log

Xact Table

00

begin_checkpoint

05

end_checkpoint

10

update:T1 writes P5

20

update:T2 writes P3

30

T1 Abort

40

CLR Undo T1 LSN 10

45

T1 End

50

update:T3 writes P1

60

update:T2 writes P5

T2; 60

T3; 50

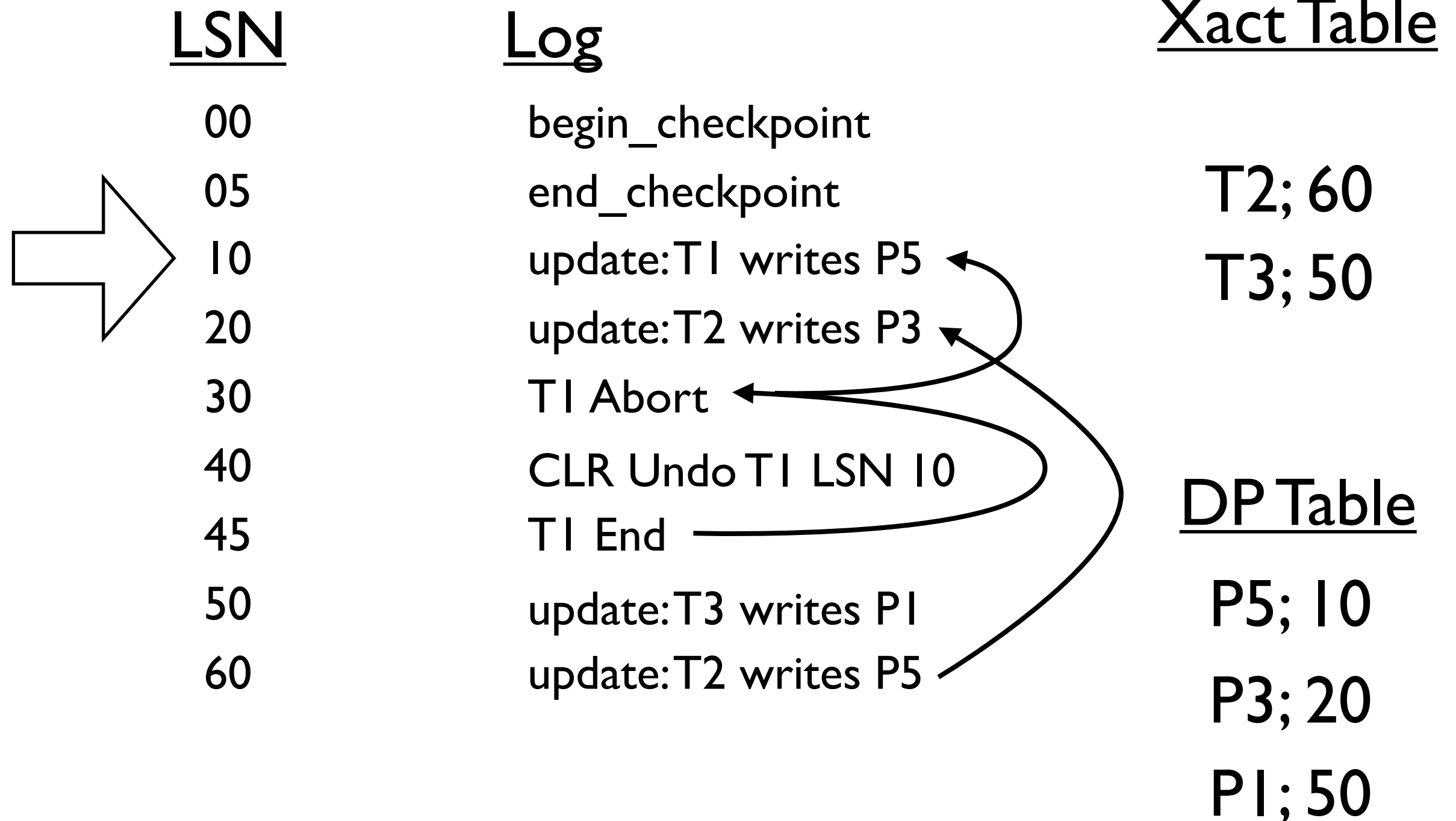
DP Table

P5; 10

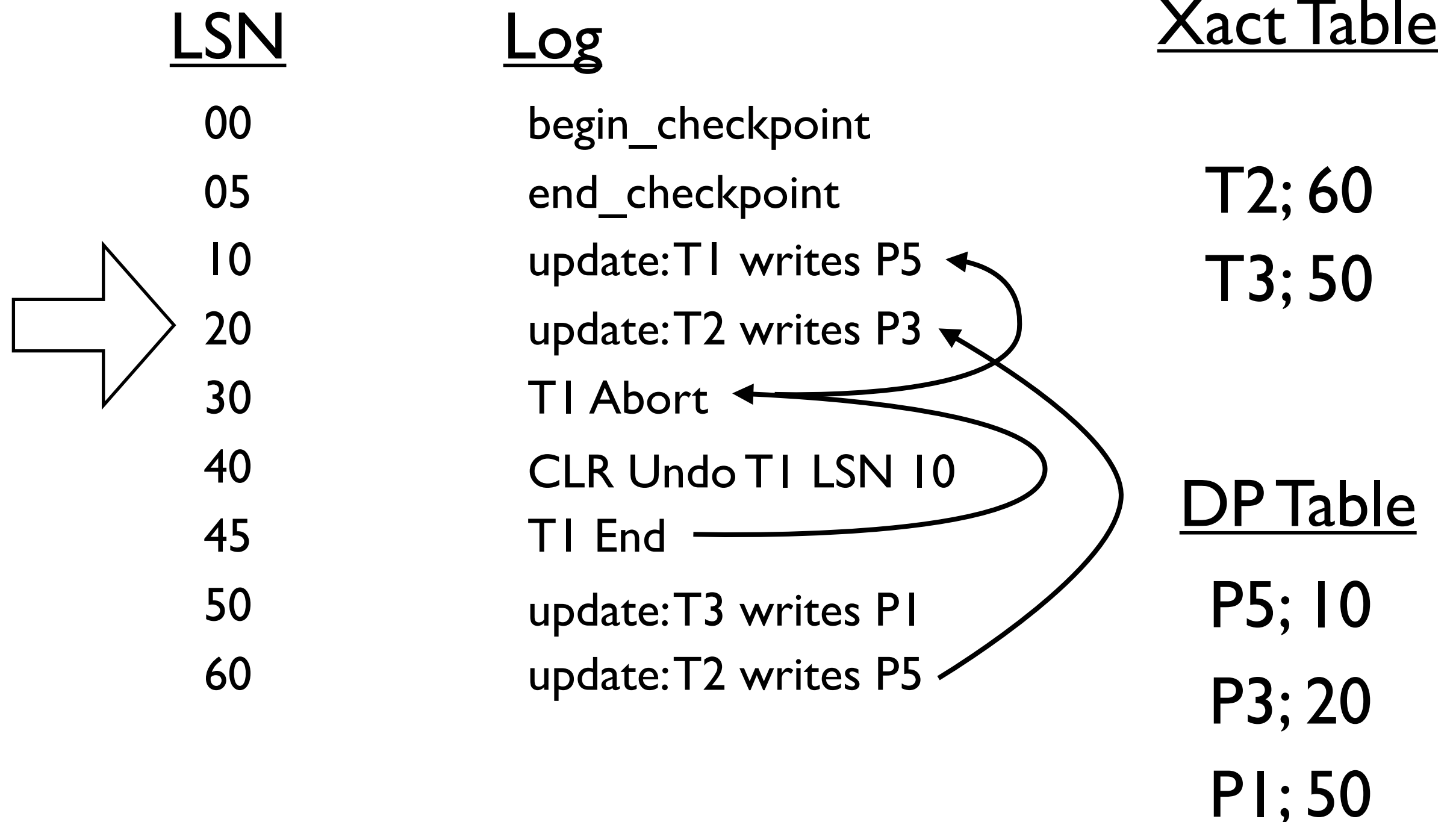
P3; 20

P1; 50

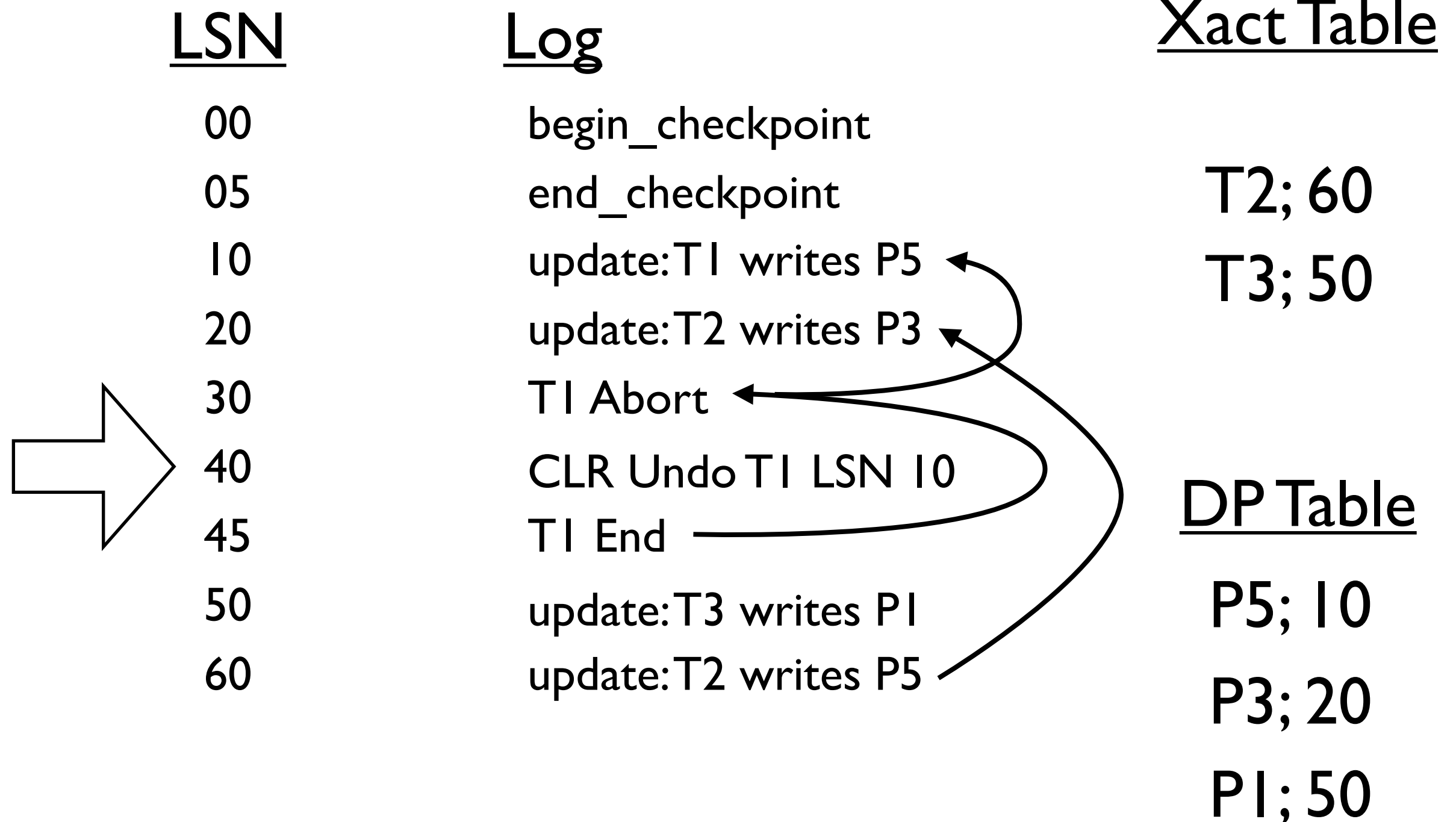
Redo



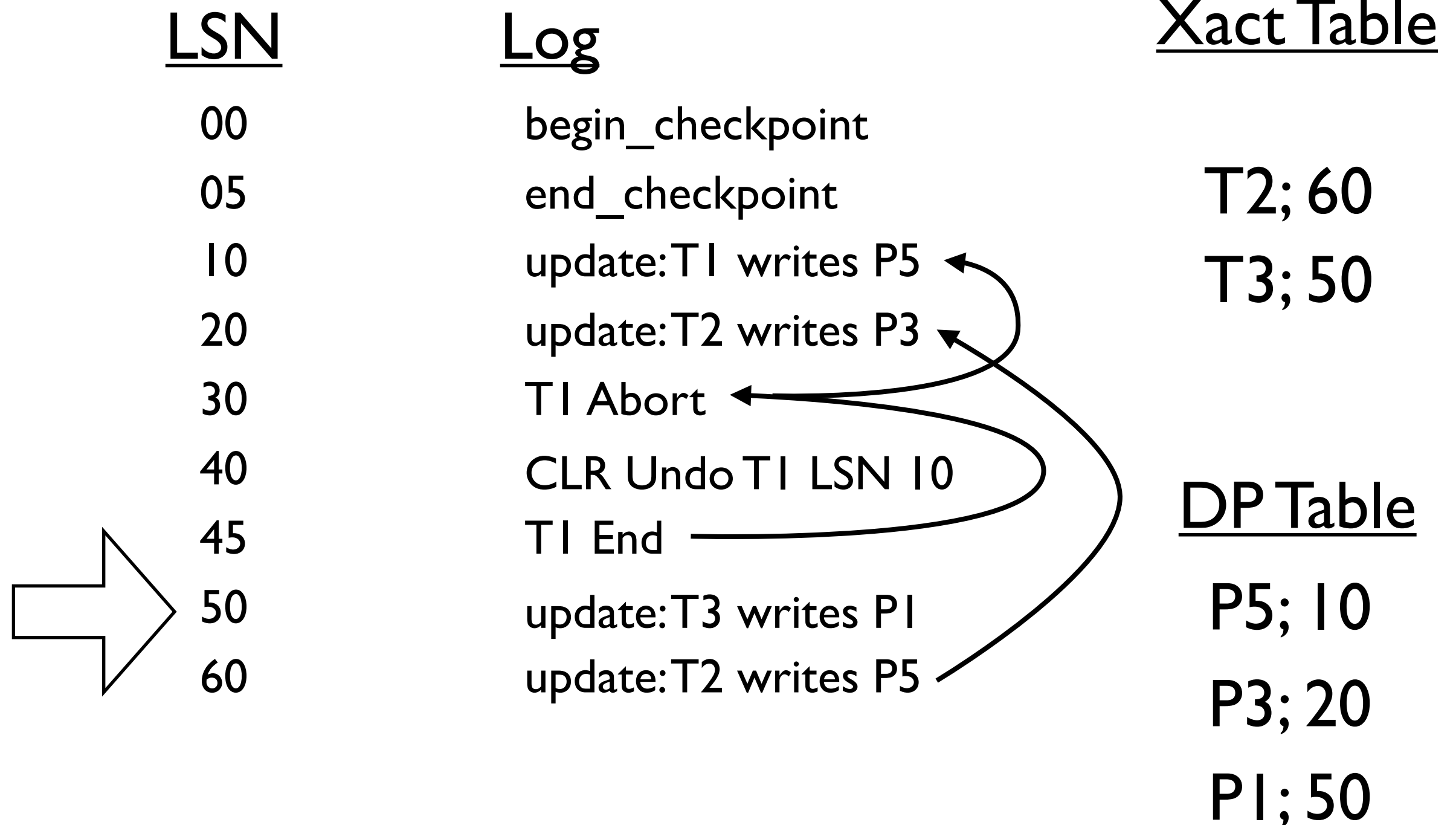
Redo



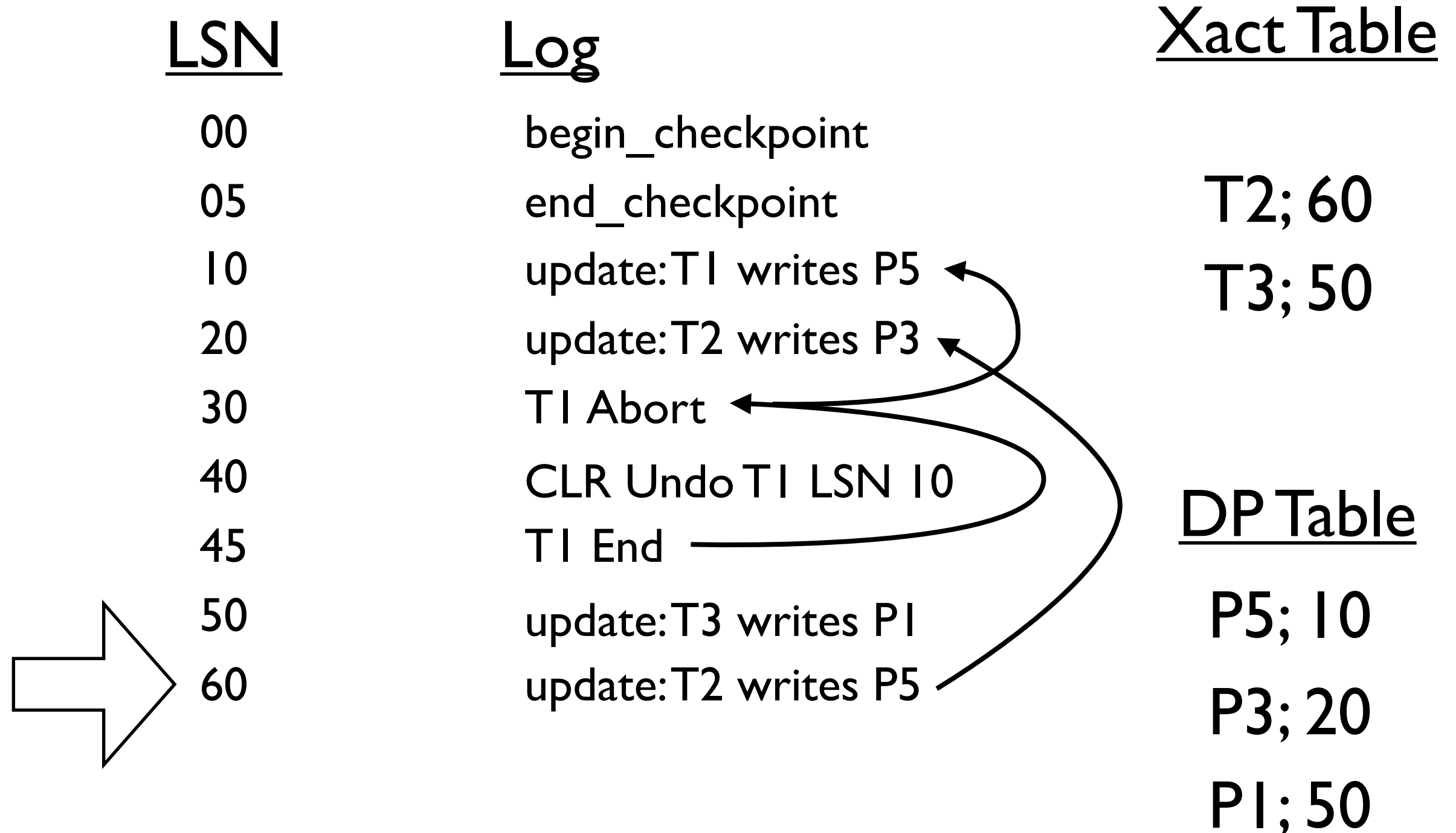
Redo



Redo



Redo



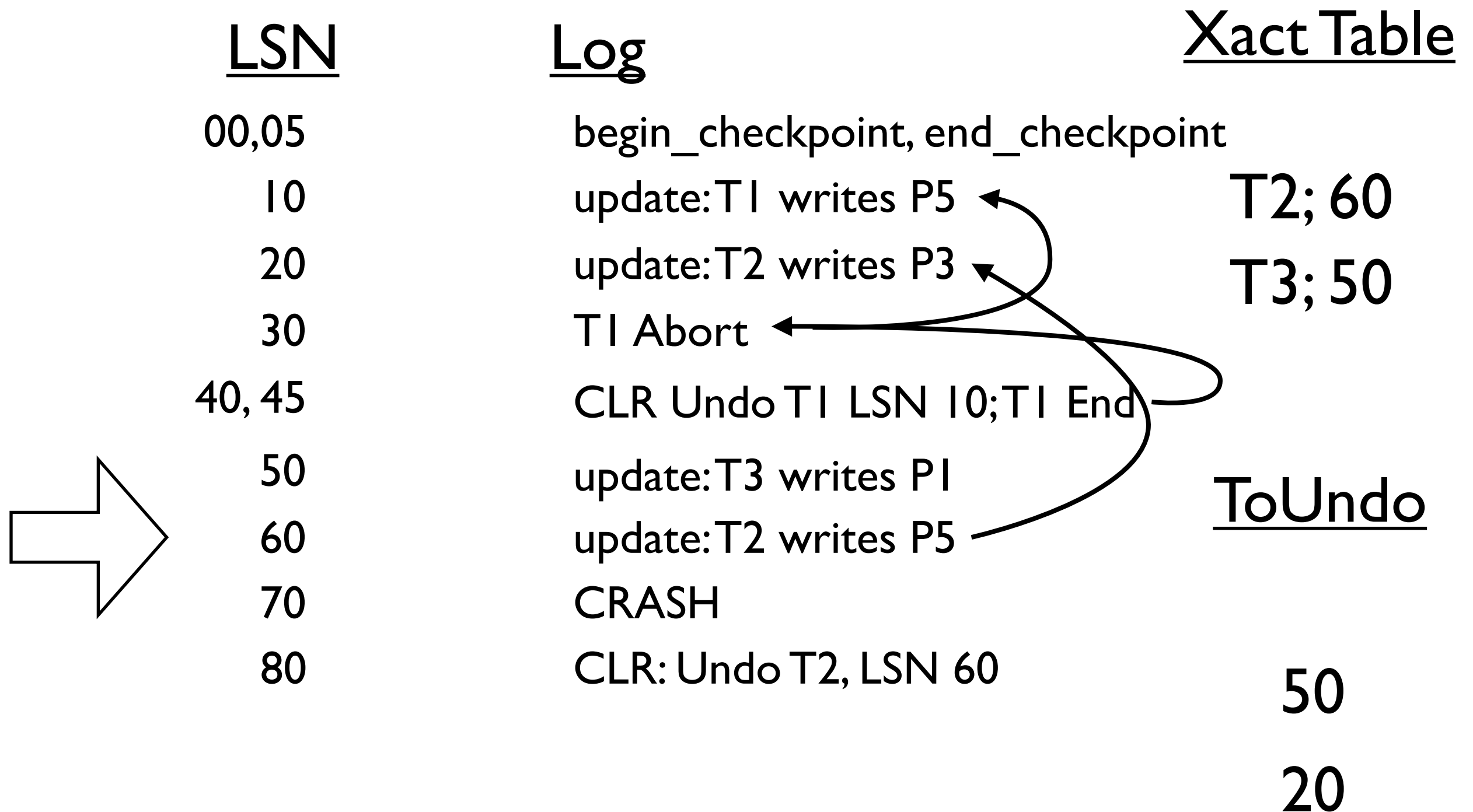
Undo

<u>LSN</u>	<u>Log</u>	<u>Xact Table</u>
00,05	begin_checkpoint, end_checkpoint	
10	update:T1 writes P5	T2; 60
20	update:T2 writes P3	T3; 50
30	T1 Abort	
40, 45	CLR Undo T1 LSN 10;T1 End	
50	update:T3 writes P1	
60	update:T2 writes P5	<u>ToUndo</u>
70	CRASH	60
		50

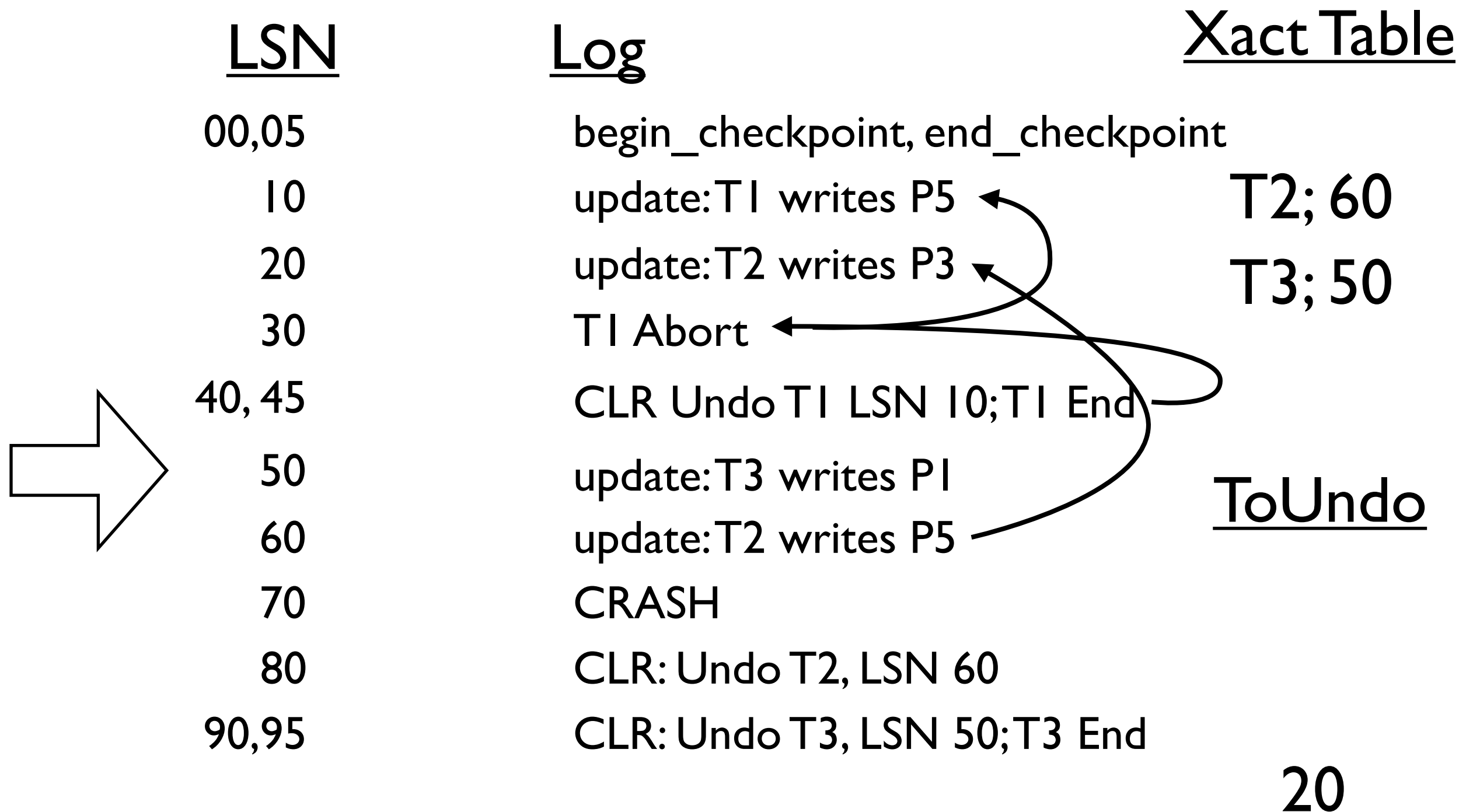
The diagram illustrates the undo process flow. Arrows indicate the following connections:

- From Log entry 'update:T1 writes P5' (LSN 10) to Xact Table entry 'T2; 60'.
- From Log entry 'update:T2 writes P3' (LSN 20) to Xact Table entry 'T3; 50'.
- From Log entry 'CLR Undo T1 LSN 10;T1 End' (LSN 40, 45) to Xact Table entry 'T2; 60'.
- From Log entry 'update:T2 writes P5' (LSN 60) to ToUndo entry '60'.

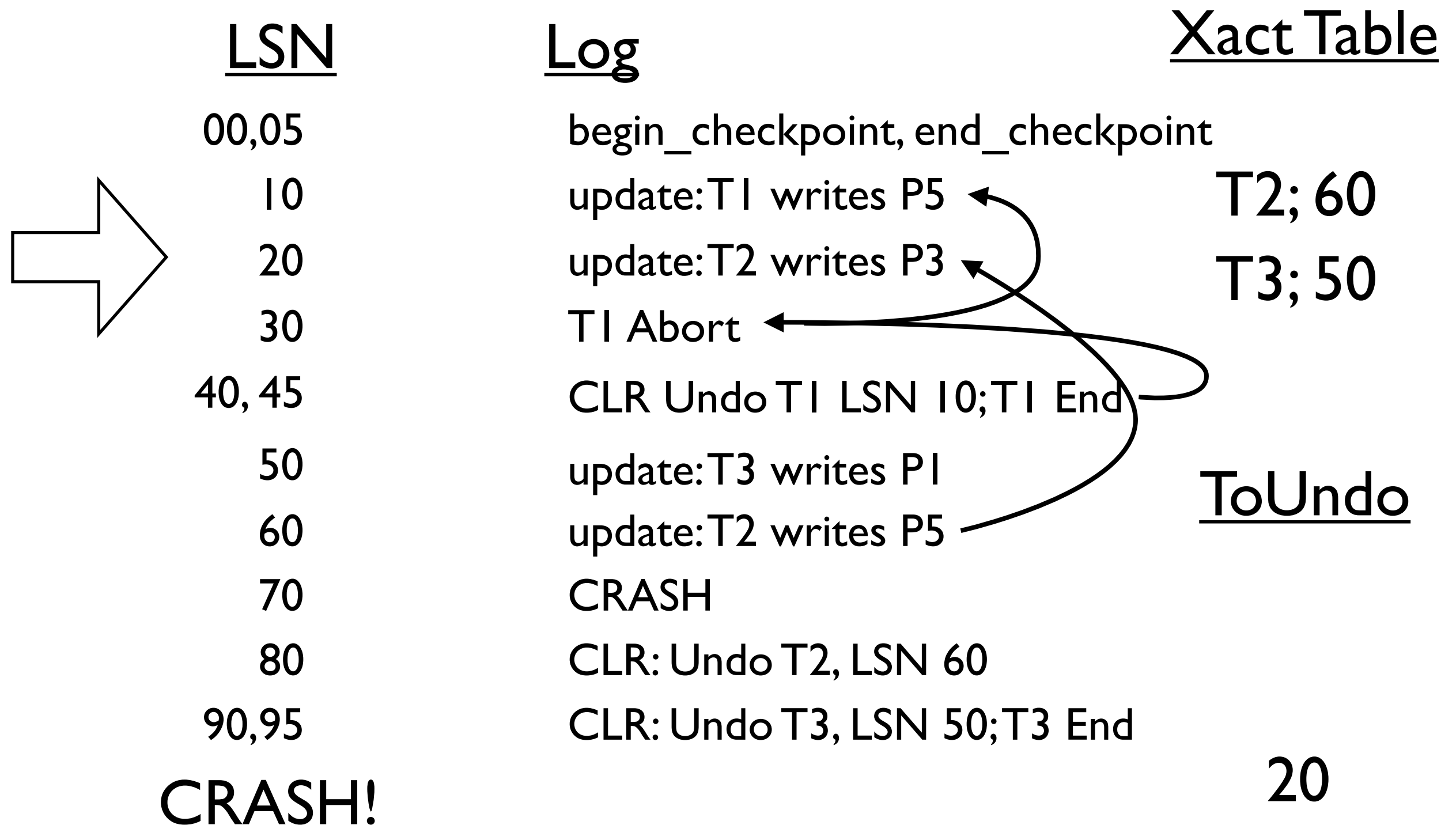
Undo



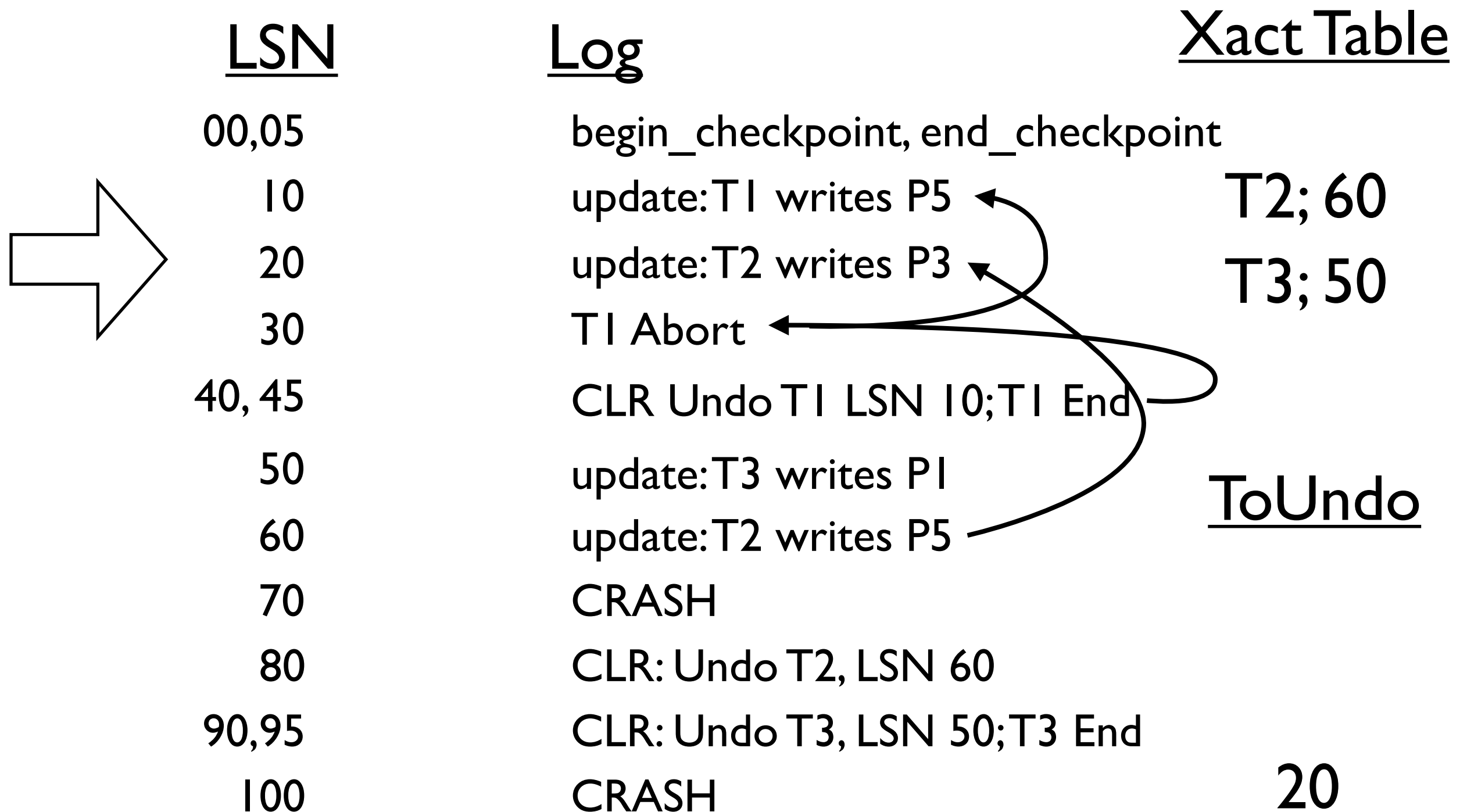
Undo



Undo



Undo



Undo

	<u>LSN</u>	<u>Log</u>	<u>Xact Table</u>
	00,05	begin_checkpoint, end_checkpoint	
	10	update:T1 writes P5	T2; 60
	20	update:T2 writes P3	T3; 50
	30	T1 Abort	
	40, 45	CLR Undo T1 LSN 10;T1 End	
	50	update:T3 writes P1	
	60	update:T2 writes P5	<u>ToUndo</u>
	70	CRASH	
	80	CLR: Undo T2, LSN 60	
	90,95	CLR: Undo T3, LSN 50;T3 End	
	100	CRASH	
	110	CLR: Undo T2, LSN 20;T2 End	

Other Crash Issues

- What happens if the system crashes during analysis?
- What happens if the system crashes during redo?
- How do you limit the amount of work in REDO?
- How do you limit the amount of work in UNDO?