Concurrency Implementation

P.J. McBrien

Imperial College London

Topic 27: REDO and UNDO logs

P.J. McBrien

Imperial College London

(日) (四) (三) (三) (三)

DBMS Architecture



æ

Recovery Manager (RM)

 $RM\ should\ protect\ the\ DBMS\ against\ failures$

æ

Recovery Manager (RM)

RM should protect the DBMS against failures

system failures

loss of volatile storage

- **1** committed transactions written to disc
- **2** uncommitted transactions not written to disc OR
- **3** sufficient information such that (1) and (2) may be met by a *recovery*

イロト 不得下 イヨト イヨト

Recovery Manager (RM)

RM should protect the DBMS against failures

system failures

loss of volatile storage

- **1** committed transactions written to disc
- **2** uncommitted transactions not written to disc OR
- **3** sufficient information such that (1) and (2) may be met by a *recovery*

media failures

loss of stable storage

1 committed data is held on multiple devices

イロト 不得下 イヨト イヨト

Before and After Images

before image

branch			
<u>sortcode</u>	bname	cash	
56	'Wimbledon'	94340.45	
34	'Goodge St'	8900.67	
67	'Strand'	34005.00	
	\Downarrow		
	$w_1[b_{56}]$		
	\Downarrow		
	branch		
<u>sortcode</u>	bname	cash	
56	'Wimbledon'	84340.45	
34	'Goodge St'	8900.67	
67	'Strand'	34005.00	
	after image		

```
BEGIN TRANSACTION

UPDATE branch

SET cash=cash-10000.00

WHERE sortcode=56

UPDATE branch

SET cash=cash+10000.00

WHERE sortcode=34

COMMIT TRANSACTION
```

- before image allows RM to undo w₁[b₅₆]
- after image allows RM to **redo** w₁[b₅₆]

<ロト <四ト < 回ト < 回ト

REDO and UNDO

Enhanced Data Manager Architecture



■ Need to cache log as well

æ

・ロト ・聞ト ・ヨト ・ヨト

Need to REDO



REDO required if committed transactions not in stable storage
must write all REDO to log before commit of transaction

3

Need to UNDO



- UNDO required if non-committed transactions in stable storage
- Must flush UNDO to log before corresponding write to data

3

Database Logs: REDO/UNDO Log



æ

Quiz 27.1: Contents of Data Disc After a Transaction with REDO/UNDO Log

	branch	
sortcode	bname	cash
56	'Wimbledon'	94340.45
34	'Goodge St'	8900.67
67	'Strand'	34005.00
	ANSACTION T1	
SET	cash=cash-10000.0	00
WHE	RE sortcode=56	
UPD	ATE branch	
SET	cash=cash+10000	.00
WHE	RE sortcode=34	
COMMIT	TRANSACTION T	1

What must the contents of the **branch** table on the data disc be after the transaction commits?



Quiz 27.2: Contents of REDO/UNDO Log Disc After a Transaction

Data Disc Before Transaction

branch			
sortcode	bname	cash	
56	'Wimbledon'	94340.45	
34	'Goodge St'	8900.67	
67	'Strand'	34005.00	

BEGIN TRANSACTION T1 UPDATE branch SET cash=cash-10000.00 WHERE sortcode=56

UPDATE branch SET cash=cash+10000.00 WHERE sortcode=34 COMMIT TRANSACTION T1 Data Disc At Commit Time

イロト イヨト イヨト イヨト

branch			
sortcode	bname	cash	
56	'Wimbledon'	94340.45	
34	'Goodge St'	18900.67	
67	'Strand'	34005.00	

What must be on the log disc after commit time?

А	В	С	D
REDO r_{56}	REDO r_{56}	UNDO r_{34}	REDO r_{56}
REDO r_{34}	UNDO r_{34}		
UNDO r_{56}			
UNDO r_{34}			

What must a complete REDO/UNDO log contain?

 $Must\ contain$

- REDO information for each update
- UNDO information for each update
- \blacksquare commit of each transaction

 $Might\ contain$

- \blacksquare begin of each transaction
 - can be inferred from first REDO/UNDO
 - presence useful to stop search of UNDO records
- abort of each transaction
 - can be inferred from lack of commit
 - presence useful to indicate UNDO already done

Rules for log and data updates

write ahead logging (WAL)

Redo rule

- \blacksquare commit \rightarrow flush log of transaction to disc
- never respond to scheduler before log written

Undo rule:

 \blacksquare flushing uncommitted data \rightarrow flush log of operations

 $\texttt{1} \mathsf{UNDO} \to \mathit{Scan \ back \ through \ the \ log}$

- \blacksquare Collect set of committed transactions C
- Collect set of incomplete transactions I
- Perform UNDO for any transaction in I

2 REDO \rightarrow Scan forward through the log

 \blacksquare Perform REDO for any transaction in C

(日) (四) (日) (日) (日)

- Collect set of committed transactions $C = \{v, y\}$
- Collect set of incomplete transactions I
- Perform UNDO for any transaction in I

2 REDO \rightarrow Scan forward through the log

 \blacksquare Perform REDO for any transaction in C

(日) (四) (日) (日) (日)

$$\begin{array}{l} \text{initial} \\ \text{state} \end{array} \Rightarrow w_v[o_1], c_v, w_x[o_2], w_y[o_1], c_y, w_z[o_2] \Rightarrow \end{array} \begin{array}{l} \text{final} \\ \text{state} \end{array}$$

- Collect set of committed transactions $C = \{v, y\}$
- Collect set of incomplete transactions $I = \{x, z\}$
- Perform UNDO for any transaction in I

2 REDO \rightarrow Scan forward through the log

 \blacksquare Perform REDO for any transaction in C

<ロト <四ト < 回ト < 回ト

$$\begin{array}{l} \begin{array}{l} \mbox{initial} \\ \mbox{state} \end{array} \Rightarrow w_v[o_1], c_v, w_x[o_2], w_y[o_1], c_y, w_z[o_2] \Rightarrow \end{array} \begin{array}{l} \begin{array}{l} \mbox{final} \\ \mbox{state} \end{array}$$

- Collect set of committed transactions $C = \{v, y\}$
- Collect set of incomplete transactions $I = \{x, z\}$
- Perform UNDO for any transaction in $I = w_z[o_2], w_x[o_2]$

2 REDO \rightarrow Scan forward through the log

 \blacksquare Perform REDO for any transaction in C

final state

イロト イ理ト イヨト イヨト

1 UNDO \rightarrow Scan back through the log

- Collect set of committed transactions $C = \{v, y\}$
- Collect set of incomplete transactions $I = \{x, z\}$
- Perform UNDO for any transaction in $I = w_z[o_2], w_x[o_2]$

2 REDO \rightarrow Scan forward through the log

Perform REDO for any transaction in $C = w_v[o_1], w_y[o_1]$

$$\begin{array}{l} \text{initial} \\ \text{state} \end{array} \Rightarrow w_v[o_1], c_v, w_x[o_2], w_y[o_1], c_y, w_z[o_2] \Rightarrow \end{array} \begin{array}{l} \text{final} \\ \text{state} \end{array}$$

- Collect set of committed transactions $C = \{v, y\}$
- Collect set of incomplete transactions $I = \{x, z\}$
- Perform UNDO for any transaction in $I = w_z[o_2], w_x[o_2]$

2 REDO \rightarrow Scan forward through the log

Perform REDO for any transaction in $C = w_v[o_1], w_y[o_1]$

$$\Rightarrow \text{UNDO}(w_{z}[o_{2}]), \text{UNDO}(w_{x}[o_{2}]), \text{REDO}(w_{v}[o_{1}]), \text{REDO}(w_{y}[o_{1}]) \Rightarrow \text{recovered} \text{state}$$

<ロト <四ト < 回ト < 回ト

\mathbf{Log}	
LOG	b_4
LOG	b_1
UNDO	$w_1[b_{56}, cash=94340.45]$
REDO	$w_1[b_{56}, cash=84340.45]$
LOG	b_2
UNDO	$w_2[b_{34}, cash=10900.67]$
REDO	$w_2[b_{34}, cash=8900.67]$
UNDO	$w_2[b_{67}, cash=34005.00]$
REDO	$w_2[b_{67}, cash=36005.25]$
LOG	b_7
LOG	c_2
UNDO	$w_1[b_{34}, cash=8900.67]$
REDO	$w_1[b_{34}, cash=18900.67]$
UNDO	$w_7[b_{67}, cash=36005.25]$
REDO	$w_7[b_{67}, cash=37005.25]$
LOG	c_7
LOG	CA

Disc	Before	Reco	very

	branch	
<u>sortcode</u>	bname	cash
56	'Wimbledon'	84340.45
34	'Goodge St'	18900.67
67	'Strand'	34005.00

Disc After Recovery

branch			
<u>sortcode</u>	bname	cash	
56	'Wimbledon'		
34	'Goodge St'		
67	'Strand'		

イロト イヨト イヨト イヨト

\mathbf{Log}	
LOG	b_4
LOG	b_1
UNDO	$w_1[b_{56}, cash=94340.45]$
REDO	$w_1[b_{56}, cash=84340.45]$
LOG	b_2
UNDO	$w_2[b_{34}, cash=10900.67]$
REDO	$w_2[b_{34}, cash=8900.67]$
UNDO	$w_2[b_{67}, cash=34005.00]$
REDO	$w_2[b_{67}, cash=36005.25]$
LOG	b_7
LOG	c_2
UNDO	$w_1[b_{34}, cash=8900.67]$
REDO	$w_1[b_{34}, cash=18900.67]$
UNDO	$w_7[b_{67}, cash=36005.25]$
REDO	$w_7[b_{67}, cash=37005.25]$
LOG	C7
LOG	C4

\mathbf{Disc}	Before	Reco	very

branch			
<u>sortcode</u>	bname	cash	
56	'Wimbledon'	84340.45	
34	'Goodge St'	18900.67	
67	'Strand'	34005.00	

Disc After Recovery

	branch	
<u>sortcode</u>	bname	cash
56	'Wimbledon'	
34	'Goodge St'	8900.67
67	'Strand'	

イロト イヨト イヨト イヨト

\mathbf{Log}	
LOG	b_4
LOG	b_1
UNDO	$w_1[b_{56}, cash=94340.45]$
REDO	$w_1[b_{56}, cash=84340.45]$
LOG	b_2
UNDO	$w_2[b_{34}, cash=10900.67]$
REDO	$w_2[b_{34}, cash=8900.67]$
UNDO	$w_2[b_{67}, cash=34005.00]$
REDO	$w_2[b_{67}, cash=36005.25]$
LOG	b_7
LOG	c_2
UNDO	$w_1[b_{34}, cash=8900.67]$
REDO	$w_1[b_{34}, cash=18900.67]$
UNDO	$w_7[b_{67}, cash=36005.25]$
REDO	$w_7[b_{67}, cash=37005.25]$
LOG	C7
LOG	CA

	branch	
<u>sortcode</u>	bname	cash
56	'Wimbledon'	84340.45
34	'Goodge St'	18900.67
67	'Strand'	34005.00

Disc After Recovery

	branch	
<u>sortcode</u>	bname	cash
56	'Wimbledon'	94340.45
34	'Goodge St'	8900.67
67	'Strand'	

イロト イヨト イヨト イヨト

\mathbf{Log}	
LOG	b_4
LOG	b_1
UNDO	$w_1[b_{56}, cash=94340.45]$
REDO	$w_1[b_{56}, cash=84340.45]$
LOG	b_2
UNDO	$w_2[b_{34}, cash=10900.67]$
REDO	$w_2[b_{34}, cash=8900.67]$
UNDO	$w_2[b_{67}, cash=34005.00]$
REDO	$w_2[b_{67}, cash=36005.25]$
LOG	b_7
LOG	c_2
UNDO	$w_1[b_{34}, cash=8900.67]$
REDO	$w_1[b_{34}, cash=18900.67]$
UNDO	$w_7[b_{67}, cash=36005.25]$
REDO	$w_7[b_{67}, cash=37005.25]$
LOG	C7
LOG	C4

\mathbf{Disc}	Before	Reco	very

branch			
<u>sortcode</u>	bname	cash	
56	'Wimbledon'	84340.45	
34	'Goodge St'	18900.67	
67	'Strand'	34005.00	

Disc After Recovery

	branch	
<u>sortcode</u>	bname	cash
56	'Wimbledon'	94340.45
34	'Goodge St'	8900.67
67	'Strand'	

イロト イヨト イヨト イヨト

\mathbf{Log}	
LOG	b_4
LOG	b_1
UNDO	$w_1[b_{56}, cash=94340.45]$
REDO	$w_1[b_{56}, cash=84340.45]$
LOG	b_2
UNDO	$w_2[b_{34}, cash=10900.67]$
REDO	$w_2[b_{34}, cash=8900.67]$
UNDO	$w_2[b_{67}, cash=34005.00]$
REDO	$w_2[b_{67}, cash=36005.25]$
LOG	b ₇
LOG	c_2
UNDO	$w_1[b_{34}, cash=8900.67]$
REDO	$w_1[b_{34}, cash=18900.67]$
UNDO	$w_7[b_{67}, cash=36005.25]$
REDO	$w_7[b_{67}, cash=37005.25]$
LOG	C7
LOG	C4

\mathbf{Disc}	Before	Recovery	
		branch	

branch	
bname	cash
'Wimbledon'	84340.45
'Goodge St'	18900.67
'Strand'	34005.00
	branch bname 'Wimbledon' 'Goodge St' 'Strand'

Disc After Recovery

	branch	
<u>sortcode</u>	bname	cash
56	'Wimbledon'	94340.45
34	'Goodge St'	8900.67
67	'Strand'	36005.25

イロト イヨト イヨト イヨト

\mathbf{Log}	
LOG	b_4
LOG	b_1
UNDO	$w_1[b_{56}, cash=94340.45]$
REDO	$w_1[b_{56}, cash=84340.45]$
LOG	b_2
UNDO	$w_2[b_{34}, cash=10900.67]$
REDO	$w_2[b_{34}, cash=8900.67]$
UNDO	$w_2[b_{67}, cash=34005.00]$
REDO	$w_2[b_{67}, cash=36005.25]$
LOG	b ₇
LOG	c_2
UNDO	$w_1[b_{34}, cash=8900.67]$
REDO	$w_1[b_{34}, cash=18900.67]$
UNDO	$w_7[b_{67}, cash=36005.25]$
REDO	$w_7[b_{67}, cash=37005.25]$
LOG	C7
LOG	c_4

\mathbf{Disc}	Before	Recovery

	branch	
<u>sortcode</u>	bname	cash
56	'Wimbledon'	84340.45
34	'Goodge St'	18900.67
67	'Strand'	34005.00

Disc After Recovery

	branch	
<u>sortcode</u>	bname	cash
56	'Wimbledon'	94340.45
34	'Goodge St'	8900.67
67	'Strand'	37005.25

イロト イヨト イヨト イヨト

Omitting the REDO Log

If no REDO records kept

must flush committed transactions to data disc

- $C = \emptyset, D = \emptyset$
- **2** Scan the log backwards from the end.
- **3** commit entry \rightarrow add to C
- **4** undo entry for member of $C \to \text{add}$ object to D without making changes to the data.
- **5** perform undo entry for object not of member D

Omitting the UNDO Log

If no UNDO records kept

transaction must never write uncommitted data

- add fix command between RM and CM to stop CM flushing data
- commit is followed by flush or **unfix** of fixed objects

Omitting UNDO and REDO

atomic commit \rightarrow out of place updating

3

・ロト ・聞ト ・ヨト ・ヨト

Quiz 27.3: Contents of Disc Before Commit if no UNDO log

	branch (T)	branch (T)	branch (1)	branch (1)	branch (1) branch (3)
sortcode	sortcode bname	sortcode bname cash	sortcode bname cash sort	sortcode bname cash sortcode	sortcode bname cash sortcode bname
56	56 'Wimbledon'	56 'Wimbledon' 94340.45	56 'Wimbledon' 94340.45	56 'Wimbledon' 94340.45 56	56 'Wimbledon' 94340.45 56 'Wimbledon'
34	34 'Goodge St'	34 'Goodge St' 8900.67	34 'Goodge St' 8900.67	34 'Goodge St' 8900.67 34	34 'Goodge St' 8900.67 34 'Goodge St'
67	67 'Strand'	67 'Strand' 34005.00	67 'Strand' 34005.00	67 'Strand' 34005.00 67	67 'Strand' 34005.00 67 'Strand'
		<u> </u>			
	branch (2)	branch (2)	branch (2)	branch (2)	branch (2) branch (4)
sortcode	sortcode bname	sortcode bname cash	sortcode bname cash sort	sortcode bname cash sortcode	sortcode bname cash sortcode bname
56	56 'Wimbledon'	56 'Wimbledon' 94340.45	56 'Wimbledon' 94340.45	56 'Wimbledon' 94340.45 56	56 'Wimbledon' 94340.45 56 'Wimbledon'
34	34 'Goodge St'	34 'Goodge St' 18900.67	34 'Goodge St' 18900.67	34 'Goodge St' 18900.67 34	34 'Goodge St' 18900.67 34 'Goodge St'
67	67 'Strand'	67 'Strand' 34005.00	67 'Strand' 34005.00	67 'Strand' 34005.00 67	67 'Strand' 34005.00 67 'Strand'

What must the contents of the **branch** table on disc be before the transaction commits?



э

Quiz 27.4: Contents of Disc After Commit if no REDO log

	branch (T)	branch (T)	branch (1)	branch (1)	branch (1) branch (3)
sortcode	sortcode bname	sortcode bname cash	sortcode bname cash sort	sortcode bname cash sortcode	sortcode bname cash sortcode bname
56	56 'Wimbledon'	56 'Wimbledon' 94340.45	56 'Wimbledon' 94340.45	56 'Wimbledon' 94340.45 56	56 'Wimbledon' 94340.45 56 'Wimbledon'
34	34 'Goodge St'	34 'Goodge St' 8900.67	34 'Goodge St' 8900.67	34 'Goodge St' 8900.67 34	34 'Goodge St' 8900.67 34 'Goodge St'
67	67 'Strand'	67 'Strand' 34005.00	67 'Strand' 34005.00	67 'Strand' 34005.00 67	67 'Strand' 34005.00 67 'Strand'
		<u> </u>			
	branch (2)	branch (2)	branch (2)	branch (2)	branch (2) branch (4)
sortcode	sortcode bname	sortcode bname cash	sortcode bname cash sort	sortcode bname cash sortcode	sortcode bname cash sortcode bname
56	56 'Wimbledon'	56 'Wimbledon' 94340.45	56 'Wimbledon' 94340.45	56 'Wimbledon' 94340.45 56	56 'Wimbledon' 94340.45 56 'Wimbledon'
34	34 'Goodge St'	34 'Goodge St' 18900.67	34 'Goodge St' 18900.67	34 'Goodge St' 18900.67 34	34 'Goodge St' 18900.67 34 'Goodge St'
67	67 'Strand'	67 'Strand' 34005.00	67 'Strand' 34005.00	67 'Strand' 34005.00 67	67 'Strand' 34005.00 67 'Strand'

What must the contents of the branch table on disc be after the transaction commits?



э

Topic 28: Checkpoints

P.J. McBrien

Imperial College London

イロト イロト イヨト イヨト 一回

Checkpointing



- Forces the database into some known state
- Recovery limited to only look back to checkpoint (or a 'bit' before!)
 - speeds the recovery operation
 - limits the size of log
- The more consistent this known state
 - the easier it is to recover
 - the longer it takes to perform the checkpoint

Commit Consistent Checkpoint

Generating a Commit Consistent Checkpoint

- **1** Stop accepting new transactions
- **2** Finish existing transactions.
- **3** Flush all dirty data cache objects to disc.
- 4 Write a checkpoint to stable log.
- \blacksquare recovery now only needs to scan back to cp in log \checkmark
- \blacksquare possible long hold-up at checkpoint \bigstar

イロト イヨト イヨト

Cache Consistent Checkpoint

Generating a Cache Consistent Checkpoint

- **1** Suspend all transactions
- **2** Flush all dirty cache objects to disc
- **3** Write a checkpoint + active transactions to stable log

Recovery from Cache Consistent Checkpoint records

- \blacksquare perform UNDOs of non-committed transactions back to cp
- **2** perform UNDO of non-committed transactions before cp if they were active at cp
- **3** perform REDOs of committed transactions after cp
- could still have delay whilst flushing cached objects

Worksheet: Cache Consistent Checkpoint

b_7
$w_7[b_{67}, cash=34005.25]$
$w_7[b_{67}, cash=37005.25]$
b_2
$w_2[b_{34}, cash=10900.67]$
$w_2[b_{34}, cash=8900.67]$
b_6
$w_6[a_{101}, rate=5.25]$
$w_6[a_{101}, rate=6.00]$
b_1
$w_1[b_{56}, cash=94340.45]$
$w_1[b_{56}, cash=84340.45]$
a_7
$cp\{1, 2, 6\}$
:

	÷
UNDO	$w_6[a_{119}, rate=5.50]$
REDO	$w_6[a_{119}, rate=6.00]$
LOG	c_6
UNDO	$w_2[b_{67}, cash=34005.00]$
REDO	$w_2[b_{67}, cash=36005.25]$
LOG	b_8
LOG	c_2
UNDO	$w_1[b_{34}, cash=8900.67]$
REDO	$w_1[b_{34}, cash=18900.67]$
LOG	b_9
UNDO	$w_9[b_{67}, cash=36005.00]$
REDO	$w_9[b_{67}, cash=20000.00]$
LOG	<i>C</i> 9

イロト イヨト イヨト イヨト

Fuzzy Checkpointing

Generating a Fuzzy Checkpoint

- **1** Suspend all transactions
- **2** Flush any dirty cache objects to disc not flushed in previous cp
- **3** Write a checkpoint + active transactions to stable log

Recovery from Fuzzy Checkpoint records

Recovery works like cache consistent checkpoint, but working with penultimate cp

- \blacksquare perform UNDOs of non-committed transactions back to penultimate cp
- **2** perform UNDO of non-committed transactions before penultimate cp if they were active at cp
- \blacksquare perform REDOs of committed transactions after penultimate cp

Media Failures: Mirroring (RAID-1)



- Keep more than one active copy of data and log
- Writes sent to both
- Read from either

æ

Media Failures: Dumping



- 'tape' might also be a external file server, removable HD, etc.
- To use normal OS backup procedure
 - DBMS must not be still running
 - raw partition must not be used

・ロト ・四ト ・ヨト ・ヨト

Checkpoints and Dumps

- Dump must do a checkpoint
- Restore involves:
 - **1** copy tape to disc
 - 2 undo transactions active at the archive time
 - 3 redo transactions that committed after the archive
- \blacksquare commit consistent checkpoint obvious choice

3

イロト イヨト イヨト

Media Failures: Archive Database



- mirror log, but only have one active database
- periodically archive updates onto archive database
- failure of active database disc involves restore of archive database using logs

・ロト ・四ト ・ヨト ・ヨト

THE END

- Content of the course is what has been presented in the lectures
- Revise by reviewing worksheets and courseworks
- 2011 to 2019 exam papers good for revision
- Revision exercises will be made available on the course homepage