















What Is Recovery? (cont.)

- During recovery, the DBMS takes the steps needed to:
 - *redo* changes made by any committed txn, if there's a chance the changes didn't make it to disk
 - \rightarrow durability: the txn's changes are still there after the crash
 - → atomicity: *all* of its changes take effect
 - *undo* changes made by any txn that didn't commit, if there's a chance the changes made it to disk
 - → atomicity: *none* of its changes take effect
 - · also used when a transaction is rolled back
- In order for recovery to work, need to maintain enough state about txns to be able to redo or undo them.

 The <i>log</i> is a file that stores the store the stores the stores the store the stores the store the store	ne info	y . needed for recovery.
 It contains: update records, each of which summarizes a write records for transaction begin and commit 	LSN 100 150 225 350 400 470	record contents txn: 1; BEGIN txn: 1; item: D1; old: 3000; new: 2500 txn: 1; item: D2; old: 1000; new: 1500 txn: 2; BEGIN txn: 2; item: D3; old: 7200; new: 6780 txn: 1: item: D1: old: 2500; new: 2750
 It does <i>not</i> record reads. don't affect the state of the database aren't relevant to recover 	550 585 675	txn: 1; COMMIT txn: 2; item: D2; old: 1500; new: 1300 txn: 2; item: D3; old: 6780; new: 6760
 The log is append-only: reco and blocks of the log file are more efficient than non-se 	rds ar writte equer	re added at the end, n to disk sequentially. itial writes to the database files



• without the log records, the recovery system can't restore it





Undo-Redo Logging (cont.)	
 To ensure that it can undo/redo txns as needed, undo-redo logging follows the WAL policy. 	
 In addition, it does the following when a transaction commits: 1. writes the commit log record to the in-memory log buffer 2. forces to disk all dirty log records (dirty = not yet written to disk) 	
• It does <i>not</i> force the dirty database pages to disk.	
At recovery, it performs two passes:	
 first, a <i>backward pass</i> to undo uncommitted transactions then, a <i>forward pass</i> to redo committed transactions 	



• He	ere's how it would work on our	earlier example	
LSN	record contents	backward pass	forward pass
100	txn: 1; BEGIN	skip	skip
150	txn: 1; item: D1; old: 3000; new: 2500	skip	redo: D1 = 250
225	txn: 1; item: D2; old: 1000; new: 1500	skip	redo: D2 = 150
350	txn: 2; BEGIN	skip	skip
400	txn: 2; item: D3; old: 7200; new: 6780	undo: D3 = 7200	skip
470	txn: 1; item: D1; old: 2500; new: 2750	skip	redo: D1 = 275
500	txn: 1; item: D2; old: 1500; new: 2100	skip	redo: D2 = 210
550	txn: 1; COMMIT	add to commit list	skip
585	txn: 2; item: D2; old: 1500; new: 1300	undo: D2 = 1500	skip
675	txn: 2; item: D3; old: 6780; new: 6760	undo: D3 = 6780	skip

- Recovery restores the database to a consistent state that reflects:
 - all of the updates by txn 1 (which committed before the crash)
 - none of the updates by txn 2 (which did not commit)

LSN	record contents	backward pass	forward pass
100	txn: 1; BEGIN	skip	l skip
150	txn: 1; item: D1; old: 3000; new: 2500	skip	redo: D1 = 2500
225	txn: 1; item: D2; old: 1000; new: 1500	skip	redo: D2 = 1500
350	txn: 2; BEGIN	skip	skip
400	txn: 2; item: D3; old: 7200; new: 6780	undo: D3 = 7200	skip
470	txn: 1; item: D1; old: 2500; new: 2750	skip	redo: D1 = 2750
500	txn: 1; item: D2; old: 1500; new: 2100	skip	redo: D2 = 2100
550	txn: 1; COMMIT	add to commit list	skip
585	txn: 2; item: D2; old: 1500; new: 1300	undo: D2 = 1500	skip
675	txn: 2; item: D3; old: 6780; new: 6760	undo: D3 = 6780	skip

1) Scanning backward at the start of recovery provides the info needed for undo / redo decisions.

• when we see an update, we already know whether the txn has committed!

LSN	record contents	backward pass	forward pass
100	txn: 1; BEGIN	skip	skip
150	txn: 1; item: D1; old: 3000; new: 2500	skip	redo: D1 = 2500
225	txn: 1; item: D2; old: 1000; new: 1500	skip	redo: D2 = 1500
350	txn: 2; BEGIN	skip	skip
400	txn: 2; item: D3; old: 7200; new: 6780	undo: D3 = 7200	skip
470	txn: 1; item: D1; old: 2500; new: 2750	skip	redo: D1 = 2750
500	txn: 1; item: D2; old: 1500; new: 2100	skip	redo: D2 = 2100
550	txn: 1; COMMIT	add to commit list	skip
585	txn: 2; item: D2; old: 1500; new: 1300	undo: D2 = 1500	skip
675	txn: 2; item: D3; old: 6780; new: 6760	undo: D3 = 6780	skip
2) Tc •	ensure the correct values are put all redos after all undos (c	on disk after rec consider D2 abo	covery, we: ve)





Storing LSNs with Data Elements

When a data element is updated, the DBMS:

txn: 1; item: D1; new: "bar"; old: "foo"; olsn: 0

log file

100

150

record

txn: 1;

- stores the LSN of the update log record with the data element
 known as the *datum LSN*
- stores the old LSN of the data element in the log record

	data eleme	nts (value / d	atum LSN)
contents	D1	D2	D3
BEGIN	"foo" / 0	"oh" / 0	"moo" / 0

"bar"/ 150





During recovery, there are three LSNs to consider		on-disk datum LSNs: D4: 0, D5: 0, <mark>D6: 1100</mark> , D7: 93	
for each update record:	LSN	record contents	
1) the record LSN : the one	700	txn: 3; BEGIN	
for the update record itself	770	txn: 3; item: D5; old: "foo"; new: "bar"; olsn: 0	
2) the on-disk <i>datum LSN</i>	825	txn: 4; BEGIN	
for the data item	850	txn: 4; item: D4; old: 9000; new: 8500; olsn: 0	
in the database file	900	txn: 4; item: D6; old: 5.7; new: 8.9; olsn: 0	
3) the olsn : the old datum LSN	930	txn: 3; item: D7; old: "zoo"; new: "cat"; olsn: 0	
for the data item	980	txn: 4; COMMIT	
 the one associated with it when the update was 	1000	txn: 3; item: D4; old: 8500; new: 7300; olsn: 850	
originally requested	1100	txn: 3; item: D6; old: 8.9; new: 4.1: olsn: 900	



	Which updates will b	e undone?	
• (datum LSNs: D4: 0 D5: 0 D	6: 1100	D7: 930
LSN	record contents	backward pass	forward pass
700	txn: 3; BEGIN		
770	txn: 3; item: D5; old: "foo"; new: "bar"; olsn: 0		
825	txn: 4; BEGIN		
850	txn: 4; item: D4; old: 9000; new: 8500; olsn: 0		
900	txn: 4; item: D6; old: 5.7; new: 8.9; olsn: 0		
930	txn: 3; item: D7; old: "zoo"; new: "cat"; olsn: 0		
980	txn: 4; COMMIT		
1000	txn: 3; item: D4; old: 8500; new: 7300; olsn: 850		
1100	txn: 3; item: D6; old: 8.9; new: 4.1; olsn: 900		

	Which updates will b	e undone?	
• (datum LSNs: D4: 0 D5: 0 D	6: 1100, <mark>900</mark>	D7: 930, <mark>0</mark>
LSN	record contents	backward pass	forward pass
700	txn: 3; BEGIN	skip	
770	txn: 3; item: D5; old: "foo"; new: "bar"; olsn: 0	0 != 770 don't undo	
825	txn: 4; BEGIN	skip	
850	txn: 4; item: D4; old: 9000; new: 8500; olsn: 0	skip	
900	txn: 4; item: D6; old: 5.7; new: 8.9; olsn: 0	skip	
930	txn: 3; item: D7; old: "zoo"; new: "cat"; olsn: 0	930 == 930 undo: D7 = "zoo" datum LSN = 0	
980	txn: 4; COMMIT	add to commit list	
1000	txn: 3; item: D4; old: 8500; new: 7300; olsn: 850	0 != 1000 don't undo	
1100	txn: 3; item: D6; old: 8.9; new: 4.1; olsn: 900	1100 == 1100 undo: D6 = 8.9 datum LSN = 900	



	Which updates will b	pe redone?	
•	datum LSNs: D4: 0 D5: 0 D	6: 1100, <mark>900</mark>	D7: 930, <mark>0</mark>
LSN	record contents	backward pass	forward pass
700	txn: 3; BEGIN	skip	
770	txn: 3; item: D5; old: "foo"; new: "bar"; olsn: 0	0 != 770 don't undo	
825	txn: 4; BEGIN	skip	
850	txn: 4; item: D4; old: 9000; new: 8500; olsn: 0	skip	
900	txn: 4; item: D6; old: 5.7; new: 8.9; olsn: 0	skip	
930	txn: 3; item: D7; old: "zoo"; new: "cat"; olsn: 0	930 == 930 undo: D7 = "zoo" datum LSN = 0	
980	txn: 4; COMMIT	add to commit list	
1000	txn: 3; item: D4; old: 8500; new: 7300; olsn: 850	0 != 1000 don't undo	
1100	txn: 3; item: D6; old: 8.9; new: 4.1; olsn: 900	1100 == 1100 undo: D6 = 8.9 datum LSN = 900	



	Which updates will t	be redone?	
• (datum LSNs: D4: ⁄0, 850 D5: 0 D	6: 1100, <mark>900</mark>	D7: 930, <mark>0</mark>
LSN	record contents	backward pass	forward pass
700	txn: 3; BEGIN	skip	skip
770	txn: 3; item: D5; old: "foo"; new: "bar"; olsn: 0	0 != 770 don't undo	skip
825	txn: 4; BEGIN	skip	skip
850	txn: 4; item: D4; old: 9000; new: 8500; olsn: 0	skip	0 == 0 redo: D4 = 8500 datum LSN = 850
900	txn: 4; item: D6; old: 5.7; new: 8.9; olsn: 0	skip	900 != 0 don't redo
930	txn: 3; item: D7; old: "zoo"; new: "cat"; olsn: 0	930 == 930 undo: D7 = "zoo" datum LSN = 0	skip
980	txn: 4; COMMIT	add to commit list	skip
1000	txn: 3; item: D4; old: 8500; new: 7300; olsn: 850	0 != 1000 don't undo	skip
1100	txn: 3; item: D6; old: 8.9; new: 4.1; olsn: 900	1100 == 1100 undo: D6 = 8.9 datum LSN = 900	skip

Checkpoints As a DBMS runs, the log gets longer and longer. thus, recovery could end up taking a very long time! To avoid long recoveries, periodically perform a *checkpoint*. force data and log records to disk to create a consistent on-disk database state during recovery, don't need to consider operations that preceded this consistent state

Static Checkpoints

- Stop activity and wait for a consistent state.
 - 1) prohibit new transactions from starting and wait until all current transactions have aborted or committed.
- Once there is a consistent state:
 - force all dirty log records to disk (dirty = not yet written to disk)
 - 3) force all dirty database pages to disk
 - 4) write a *checkpoint record* to the log
 - · these steps must be performed in the specified order!
- When performing recovery, go back to the most recent checkpoint record.
- Problem with this approach?

Dynamic Checkpoints Don't stop and wait for a consistent state. Steps: prevent all update operations force all dirty log records to disk force all dirty database pages to disk write a *checkpoint record* to the log include a list of all active txns When performing recovery: backward pass: go back until you've seen the start records of all uncommitted txns in the most recent checkpoint record forward pass: begin from the log record that comes after the most recent checkpoint record. why? note: if all txns in the checkpoint record are on the commit list, we stop the backward pass at the checkpoint record

•	Initial datum LSNs: D4: 110 D	5: 140,0 D6: 80	C
.SN	record contents	backward pass	forward pass
100	txn: 1; BEGIN		
110	txn: 1; item: D4 ; old: 20; new: 15; olsn: 0		
120	txn: 2; BEGIN	stop here	
130	txn: 1; COMMIT	add to commit list	
140	txn: 2; item: D5; old: 12; new: 13; olsn: 0	undo: D5 = 12 datum LSN = 0	
150	CHECKPOINT (active txns = 2)	note active txns	
160	txn: 2; item: D4 ; old: 15; new: 50; olsn: 110	don't undo	start here skip
170	txn: 3; BEGIN	skip	skip
180	txn: 3; item: D6; old: 6; new: 8; olsn: 80	don't undo	skip

Undo-Only Logging Only store the info. needed to undo txns. update records include only the old value Like undo-redo logging, undo-only logging follows WAL. In addition, all database pages changed by a transaction must be forced to disk before allowing the transaction to commit. Why? At transaction commit: force all dirty log records to disk force database pages changed by the txn to disk write the commit log record force the commit log record to disk During recovery, the system only performs the backward pass.

Redo-Only Logging

- Only store the info. needed to redo txns.
 - update records include only the new value
- Like the other two schemes, redo-only logging follows WAL.
- In addition, all database pages changed by a txn are held in memory until it commits and its commit record is forced to disk.
- At transaction commit:
 - 1. write the commit log record
 - 2. force all dirty log records to disk

(changed database pages are allowed to go to disk anytime after this)

- If a transaction aborts, none of its changes can be on disk.
- During recovery, perform the backward pass to build the commit list (no undos). Then perform the forward pass as in undo-redo.



Comparing the Three Logging Schemes (cont.)

- Redo-only:
 - + smaller logs than undo-redo
 - +/ recovery: more complex than undo-only, less than undo-redo
 - must be able to cache all changes until the txn commits
 - · limits the size of transactions
 - constrains the replacement policy of the cache
 - + forces only log records to disk at commit

• Undo-redo:

- larger logs
- more complex recovery
- + forces only log records to disk at commit
- + don't need to retain all data in the cache until commit

Previous the previous of the previous

Review Problem

- Recall the three logging schemes:
 undo-redo, undo-only, redo-only
- What type of logging is being used to create the log at right?

<u>txn 1</u>			
writes	75	for	D1
writes	90	for	D3
<u>txn 2</u>			
writes	25	for	D2
writes	60	for	D3

LSN	record contents
100	txn: 1; BEGIN
210	txn: 1; item: D1; old: 45
300	txn: 2; BEGIN
420	txn: 2; item: D2; old: 80
500	txn: 2; item: D3; old: 30
525	txn: 2; COMMIT
570	txn: 1; item: D3; old: 60

 Review Problem Recall the three logging schemes: undo-redo, undo-only, redo-only What type of logging is being 		<u>txn 1</u> writes 75 for D1 writes 90 for D3 <u>txn 2</u> writes 25 for D2 writes 60 for D3
used to create the log at right?	LSN	record contents
undo-only	100	txn: 1; BEGIN
	210	txn: 1; item: D1; old: 45; new: 75
 Io make the rest of the problem 	300	txn: 2; BEGIN
easier, add the new values to	420	txn: 2; item: D2; old: 80; new: 25
the log	500	txn: 2; item: D3; old: 30; new: 60
	525	txn: 2; COMMIT
	570	txn: 1; item: D3; old: 60; new: 90



Review Problem			
 Recall the three logging schemes: undo-redo, undo-only, redo-only 		<u>txn 1</u> writes 75 for D1 writes 90 for D3 <u>txn 2</u>	
At the start of recovery, what are		writes 25 for D2 writes 60 for D3	
the possible on-disk values	I.SN	record contents	
under <i>redo-only</i> ?	100	txn: 1; BEGIN	
 does pin values in memory → can't go to disk until commit at commit, unpins values but does not force them to disk 		txn: 1; item: D1; old: 45; new: 75	
		txn: 2; BEGIN	
		txn: 2; item: D2; old: 80; new: 25	
		txn: 2; item: D3; old: 30; new: 60	
→ older values are still 525		txn: 2; COMMIT	
possible	570	txn: 1; item: D3; old: 60; new: 90	
<u>in-memory</u> <u>possible on-dis</u> D1: D2: D3:	<u>sk</u>		









Distributed Atomicity (cont.)

- Example of what could go wrong:
 - a subtxn at one of the sites deadlocks and is aborted
 - before the coordinator of the txn finds out about this, it tells the other sites to commit, and they do so
- · Another example:
 - · the coordinator notifies the other sites that it's time to commit
 - · most of the sites commit their subtxns
 - · one of the sites crashes before committing

Fue-Phase Commit (2PC) A protocol for deciding whether to commit a distributed txn. Basic idea: coordinator asks sites if they're ready to commit if a site is ready, it: prepares its subtxn – putting it in the ready state tells the coordinator it's ready if all sites say they're ready, all subtxns are committed otherwise, all subtxns are aborted (i.e., rolled back) Preparing a subtxn means ensuring it can be either committed or rolled back – even after a failure. need to at least... some logging schemes need additional steps After saying it's ready, a site *must wait* to be told what to do next.

2PC Phase I: Prepare

- When it's time to commit a distributed txn T, the coordinator:
 - force-writes a prepare record for T to its own log
 - sends a prepare message to each participating site
- If a site can commit its subtxn, it:
 - takes the steps needed to put its txn in the ready state
 - force-writes a ready record for T to its log
 - sends a *ready message* for T to the coordinator and waits
- If a site needs to abort its subtxn, it:
 - force-writes a *do-not-commit record* for T to its log
 - sends a *do-not-commit message* for T to the coordinator
 - can it abort the subtxn now?
- Note: we always log a message before sending it to others.
 - · allows the decision to send the message to survive a crash







Recovery When Using 2PC (cont.)

- Case 4: the last log record for T is <u>from before 2PC began</u> (e.g., an update record).
 - undo the subtxn's updates as needed
 - this works in both of the possible situations:
 - 2PC has already completed without hearing from this site *why*?
 - 2PC is still be going on *why*?
- Case 5: the last log record for T is a <u>ready</u> record.
 - contact the coordinator (or another site) to determine T's fate
 - why can the site still commit or abort T as needed?
 - if it can't reach another site, it must block until it can reach one!



or Fails? (cont.)			
 Case 3: a site sent a ready message for T but didn't hear back poll the other sites to determine T's fate 			
conclusion/action			
???			
???			
???			
can't know T's fate unless coordinator recovers. why?			

 What type of logging is being used to create the log at right? 	original values: D1=1000, D2=3000	
	LSN	record contents
	100	txn: 1; BEGIN
 At the start of recovery, what are the possible on-disk values? 	150	txn: 1; item: D1; new: 2500
	350	txn: 2; BEGIN
	400	txn: 2; item: D2; new: 6780
	470	txn: 1; item: D1; new: 2750
	550	txn: 1; COMMIT
	585	txn: 2; item: D1; new: 1300

Extra Practice • What if the DBMS were using original values: undo-only logging instead? D1=1000, D2=3000 LSN record contents 100 txn: 1; BEGIN • At the start of recovery, what are 150 txn: 1; item: D1; new: 2500 the possible on-disk values? 350 txn: 2; BEGIN txn: 2; item: D2; new: 6780 400 470 txn: 1; item: D1; new: 2750 550 txn: 1; COMMIT txn: 2; item: D1; new: 1300 585 possible on-disk in-memory D1: 1000 D2: 3000

Extra Practice			
 What if the DBMS were using undo-redo logging instead? 	original values: D1=1000, D2=3000		
 At the start of recovery, what are the possible on-disk values? 	LSN record contents 100 txn: 1; BEGIN 150 txn: 1; item: D1; new: 2500 350 txn: 2; BEGIN 400 txn: 2; item: D2; new: 6780 470 txn: 1; item: D1; new: 2750 550 txn: 1; COMMIT 585 txn: 2; item: D1; new: 1300		
in-memory possible on-dis D1: 1000 D2: 3000	<u>sk</u>		