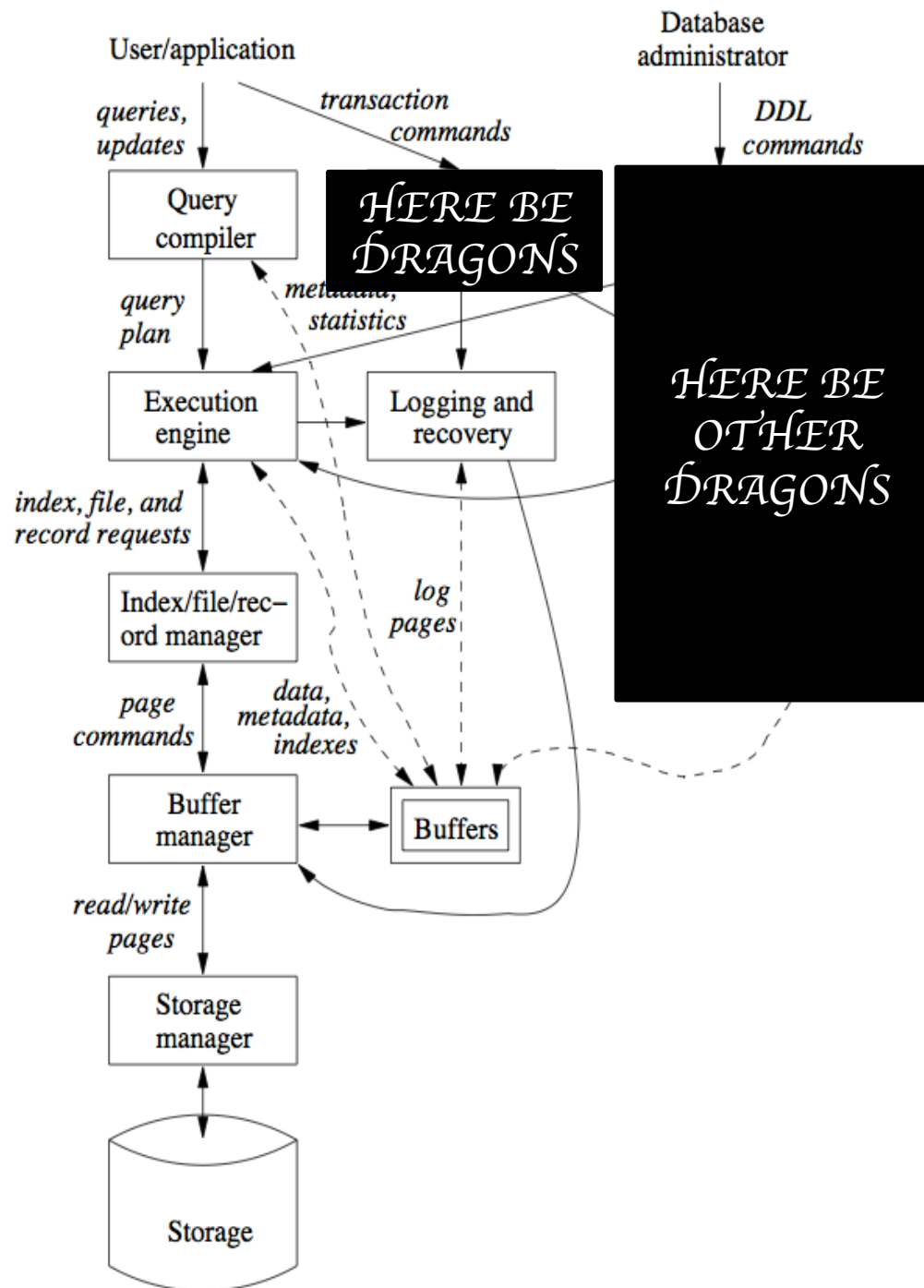


UNIVERSITY OF ILLINOIS
AT URBANA-CHAMPAIGN

CS411 - Logging and Recovery



illinois.edu



Review

- Why are we studying logging?
- What does it mean for a transaction to be “durable” or “aborted”?
- What was the name of the logging we learned last time?
- What parts of ACID does logging help ensure?



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle T_1, B, 3 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle T_1, C, 7 \rangle$

$\langle T_2, D, 5 \rangle$

$\langle \text{COMMIT } T_2 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_3, A, 2 \rangle$

$\langle \text{START } T_4 \rangle$

$\langle T_4, A, 3 \rangle$

$\langle T_3, D, 5 \rangle$

$\langle \text{COMMIT } T_4 \rangle$

A=2

B=5

C=6

D=8



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle T_1, B, 3 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle T_1, C, 7 \rangle$

$\langle T_2, D, 5 \rangle$

$\langle \text{COMMIT } T_2 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_3, A, 2 \rangle$

$\langle \text{START } T_4 \rangle$

$\langle T_4, A, 3 \rangle$

$\langle T_3, D, 5 \rangle$

$\langle \text{COMMIT } T_4 \rangle$

A=5

B=3

C=7

D=5



Undo Logging

- Rules:
 1. Add $\langle T, X, v \rangle$ to log ***before*** new value of X is written to disk
 2. Add $\langle \text{COMMIT } T \rangle$ ***after all*** results for T have been written to disk



Example

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>
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>
FLUSH LOG						
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	
						<COMMIT T>
FLUSH LOG						



Recovery

- If we see $\langle \text{START } T \rangle$ without corresponding $\langle \text{COMMIT } T \rangle$, must undo the transaction
- When the system fails, scan log backward
 - If we see $\langle \text{COMMIT } F \rangle$, ignore transaction F
 - For every other transaction T , when we see $\langle T, X, v \rangle$, write v to X



Undo Logging

- Two problems
 1. If we want to recover, we have to read the entire log file
 2. Must write all changes to disk before committing



Checkpointing

- Periodically:
 - stop accepting new transactions.
 - Wait for all active transactions to commit or abort
 - Flush the log
 - Write <CKPT> into the log
 - Resume transactions
- Recovery only needs to read to <CKPT>



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle T_2, C, 15 \rangle$

$\langle T_1, D, 20 \rangle$

$\langle \text{COMMIT } T_1 \rangle$

$\langle \text{COMMIT } T_2 \rangle$

$\langle \text{CKPT} \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_3, E, 25 \rangle$

$\langle T_3, F, 30 \rangle$



Checkpointing

- Problem
 - We have to halt all incoming transactions
 - Active transactions might take a long time to complete
 - Our entire DBMS is just waiting!



Nonquiescent Checkpointing

- “*quiescent* - adj. marked by inactivity or repose”
- Merriam-Webster
- “*quiescent* - adj. being lazy”
- Ryan Cunnigham



Nonquiescent Checkpointing

1. Write $\langle \text{START CKPT } (T_1, T_2, \dots, T_n) \rangle$ for all active transactions
2. When all of T_1, T_2, \dots, T_n commit or abort, write $\langle \text{END CKPT} \rangle$



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle \text{START CKPT}(T_1, T_2) \rangle$

$\langle T_2, C, 15 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_1, D, 20 \rangle$

$\langle \text{COMMIT } T_1 \rangle$

$\langle T_3, E, 25 \rangle$

$\langle \text{COMMIT } T_2 \rangle$

$\langle \text{END CKPT} \rangle$

$\langle T_3, F, 30 \rangle$



Recovery

- Scan log backward:
 - If we meet $\langle \text{END CKPT} \rangle$ first, continue backward to the next $\langle \text{START CKPT} \rangle$
 - If we meet $\langle \text{START CKPT} (T_1, T_2, \dots, T_n) \rangle$ first, continue backward
 - only need to undo T_1, T_2, \dots, T_n



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle \text{START CKPT}(T_1, T_2) \rangle$

$\langle T_2, C, 15 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_1, D, 20 \rangle$

$\langle \text{COMMIT } T_1 \rangle$

$\langle T_3, E, 25 \rangle$

$\langle \text{COMMIT } T_2 \rangle$

$\langle \text{END CKPT} \rangle$

$\langle T_3, F, 30 \rangle$



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle T_1, C, 2 \rangle$

$\langle \text{START CKPT}(T_1, T_2) \rangle$

$\langle T_2, C, 15 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_1, D, 20 \rangle$

$\langle \text{COMMIT } T_1 \rangle$

$\langle T_3, E, 25 \rangle$

$\langle T_3, F, 30 \rangle$



Undo Logging

- Two problems
 1. ~~If we want to recover, we have to read the entire log file~~
 2. Must write all changes to disk before committing



Redo Logging

- Main idea:
 - Log what we're going to change ***before*** we change it
 - Anything not logged was never written
 - After a crash, repeat all the committed transactions



Redo Logging

1. Before modifying X to value v on the disk, log $\langle T, X, v \rangle$
2. Before modifying X on the disk, log $\langle \text{COMMIT } T \rangle$



Example

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,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 T>
FLUSH LOG						
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	



Recovery

1. Scan log to identify all committed transactions
2. Scan log from beginning for each $\langle T, X, v \rangle$
 - If T is not committed, ignore
 - If T is committed, write v into X



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle T_1, B, 3 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle T_1, C, 7 \rangle$

$\langle T_2, D, 5 \rangle$

$\langle \text{COMMIT } T_2 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_3, A, 2 \rangle$

$\langle \text{START } T_4 \rangle$

$\langle T_4, A, 3 \rangle$

$\langle T_3, D, 5 \rangle$

$\langle \text{COMMIT } T_4 \rangle$

A=2

B=5

C=6

D=8



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle T_1, B, 3 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle T_1, C, 7 \rangle$

$\langle T_2, D, 5 \rangle$

$\langle \text{COMMIT } T_2 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_3, A, 2 \rangle$

$\langle \text{START } T_4 \rangle$

$\langle T_4, A, 3 \rangle$

$\langle T_3, D, 5 \rangle$

$\langle \text{COMMIT } T_4 \rangle$

A=3

B=5

C=6

D=5



Nonquiescent Checkpointing

1. Write $\langle \text{START CKPT } (T_1, T_2, \dots, T_n) \rangle$ for all active transactions
2. Write all modified elements for committed transactions to the disk
3. Write $\langle \text{END CKPT} \rangle$



Example

<START T_1 >

< $T_1, A, 5$ >

<START T_2 >

<COMMIT T_1 >

< $T_2, B, 10$ >

<START CKPT (T_2)>

< $T_2, C, 15$ >

<START T_3 >

<END CKPT>

<COMMIT T_2 >

<COMMIT T_3 >



Recovery

- Scan log backward:
 - If we meet $\langle \text{END CKPT} \rangle$ first, continue backward to the next $\langle \text{START CKPT } (T_1, T_2, \dots, T_n) \rangle$
 - Redo all committed transactions in T_1, T_2, \dots, T_n or started after checkpoint
 - If we meet $\langle \text{START CKPT } (T_1, T_2, \dots, T_n) \rangle$ first, this checkpoint is bunk. Continue looking for $\langle \text{END CKPT} \rangle$



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle \text{COMMIT } T_1 \rangle$

$\langle T_2, B, 10 \rangle$

$\langle \text{START CKPT } (T_2) \rangle$

$\langle T_2, C, 15 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_3, D, 20 \rangle$

$\langle \text{END CKPT} \rangle$

$\langle \text{COMMIT } T_2 \rangle$

$\langle \text{COMMIT } T_3 \rangle$



Example

<START T_1 >

< T_1 ,A,5>

<START T_2 >

<COMMIT T_1 >

< T_2 ,B,10>

<START CKPT (T_2)>

< T_2 ,C,15>

<START T_3 >

< T_3 ,D,20>

<END CKPT>

<COMMIT T_2 >

<COMMIT T_3 >



Comparison

Undo Logging

- Data must be written to disk after transaction finishes
- More Disk I/O
- Less Memory

Redo Logging

- Data must remain in memory until transaction finishes
- Less Disk I/O
- More memory



Undo/Redo Logging

- Best of both worlds:
 - Keep logs such that we can either undo or redo transactions
 - When recovering, undo all uncommitted transactions and redo all committed transactions



Undo/Redo Logging

1. Whenever a transaction modifies a value, write $\langle T, X, v, w \rangle$ to the log where
 - v is the old value
 - w is the new value
 - We don't care whether or not the changes have been made on disk



Example

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>
FLUSH LOG						
OUTPUT(A)	16	16	16	16	8	
						<COMMIT T>
OUTPUT(B)	16	16	16	16	16	



Recovery

1. Redo all committed transactions from last to first
2. Undo all uncommitted transactions in reverse



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 4, 5 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle \text{COMMIT } T_1 \rangle$

$\langle T_2, B, 9, 10 \rangle$

$\langle T_2, C, 14, 15 \rangle$

$\langle \text{START } T_4 \rangle$

$\langle T_3, D, 19, 20 \rangle$

$\langle T_4, E, 30, 31 \rangle$

$\langle \text{COMMIT } T_3 \rangle$

$\langle T_4, F, 43, 44 \rangle$



Example

<START T_1 >

< T_1 ,A,4,5>

<START T_2 >

<START T_3 >

<COMMIT T_1 >

< T_2 ,B,9,10>

< T_2 ,C,14,15>

<START T_4 >

< T_3 ,D,19,20>

< T_4 ,E,30,31>

<COMMIT T_3 >

< T_4 ,F,43,44>



Nonquiescent Checkpointing

1. Write $\langle \text{START CKPT } (T_1, T_2, \dots, T_n) \rangle$ for all active transactions
2. Write ***all*** modified elements to the disk
3. Write $\langle \text{END CKPT} \rangle$



Recovery

- Identify matching <START CKPT> and <END CKPT>
 - Undo all uncommitted transactions that were
 - active at <START CKPT>
 - started after START CKPT
 - Redo all committed transactions that committed after START CKPT



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 4, 5 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle \text{COMMIT } T_1 \rangle$

$\langle T_2, B, 9, 10 \rangle$

$\langle \text{START CKPT } (T_2) \rangle$

$\langle T_2, C, 14, 15 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_3, D, 19, 20 \rangle$

$\langle \text{END CKPT} \rangle$



Example

<START T_1 >

< T_1 ,A,4,5>

<START T_2 >

<COMMIT T_1 >

< T_2 ,B,9,10>

<START CKPT (T_2)>

< T_2 ,C,14,15>

<START T_3 >

< T_3 ,D,19,20>

<END CKPT>



Example

<START T_1 >

< T_1 ,A,4,5>

<START T_2 >

<COMMIT T_1 >

< T_2 ,B,9,10>

<START CKPT (T_2)>

< T_2 ,C,14,15>

<START T_3 >

< T_3 ,D,19,20>

<END CKPT>



Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 4, 5 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle \text{COMMIT } T_1 \rangle$

$\langle T_2, B, 9, 10 \rangle$

$\langle \text{START CKPT } (T_2) \rangle$

$\langle T_2, C, 14, 15 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_3, D, 19, 20 \rangle$

$\langle \text{END CKPT} \rangle$

$\langle \text{COMMIT } T_2 \rangle$



Example

<START T_1 >

< T_1 ,A,4,5>

<START T_2 >

<COMMIT T_1 >

< T_2 ,B,9,10>

<START CKPT (T_2)>

< T_2 ,C,14,15>

<START T_3 >

< T_3 ,D,19,20>

<END CKPT>

<COMMIT T_2 >



Example

<START T_1 >

< T_1 ,A,4,5>

<START T_2 >

<COMMIT T_1 >

< T_2 ,B,9,10>

<START CKPT (T_2)>

< T_2 ,C,14,15>

<START T_3 >

< T_3 ,D,19,20>

<END CKPT>

<COMMIT T_2 >



Example

<START T_1 >

< T_1 ,A,4,5>

<START T_2 >

<COMMIT T_1 >

< T_2 ,B,9,10>

<START CKPT (T_2)>

< T_2 ,C,14,15>

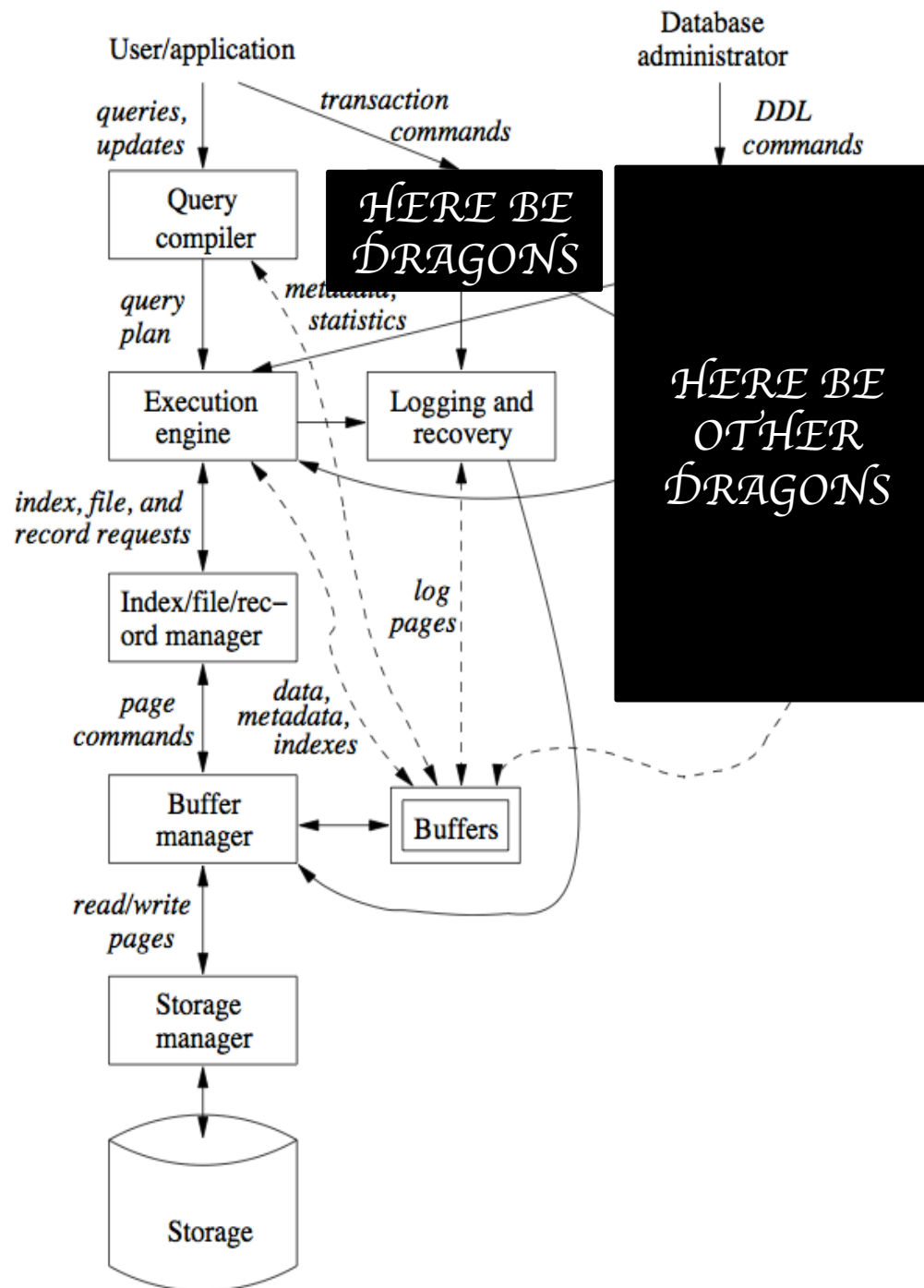
<START T_3 >

< T_3 ,D,19,20>

<END CKPT>

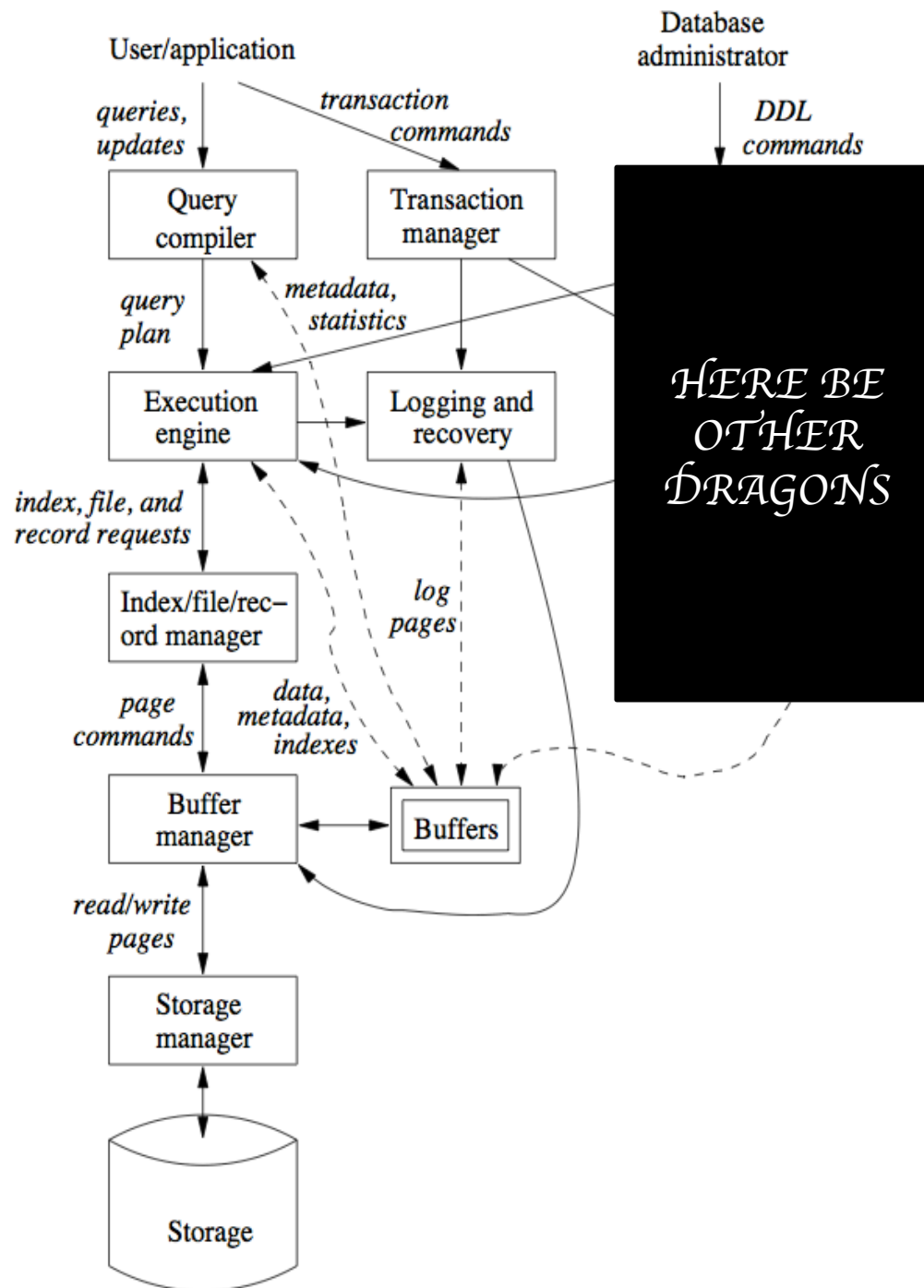
<COMMIT T_2 >



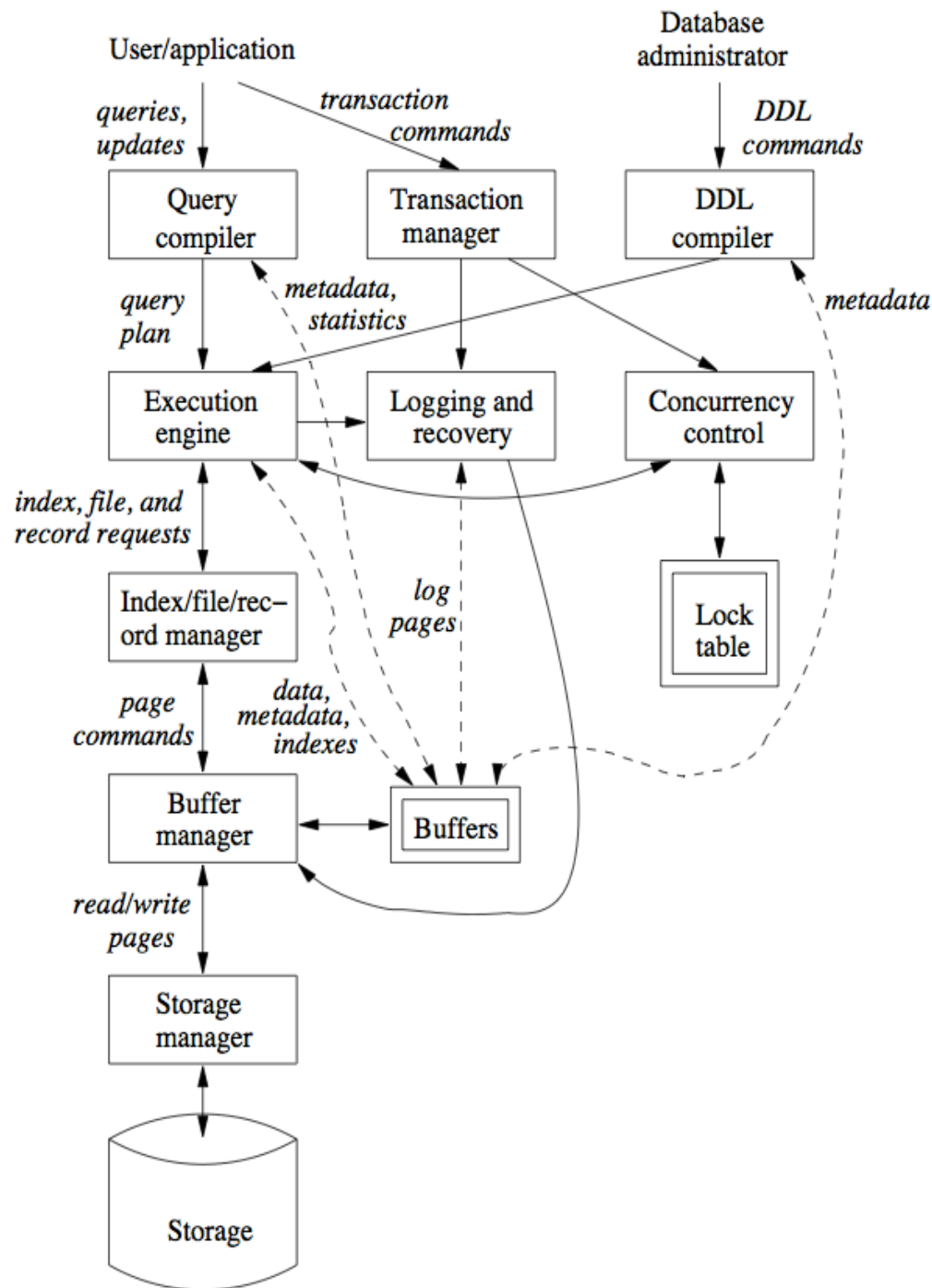




illinois.edu







Concurrency Control

- Interaction among transactions can cause the database state to be inconsistent
- Need to regulate the order of individual steps of transactions
- Assure that transactions preserve consistency when executing simultaneously



Schedule

- Time-ordered sequence of the important actions taken by one or more transactions
- We only care about READ and WRITE actions
 - The same database primitives from logging



Example

T_1	T_2
READ(A, t)	READ(A, s)
$t := t + 100$	$s := s * 2$
WRITE(A, t)	WRITE(A, s)
READ(B, t)	READ(B, s)
$t := t + 100$	$s := s * 2$
WRITE(B, t)	WRITE(B, s)



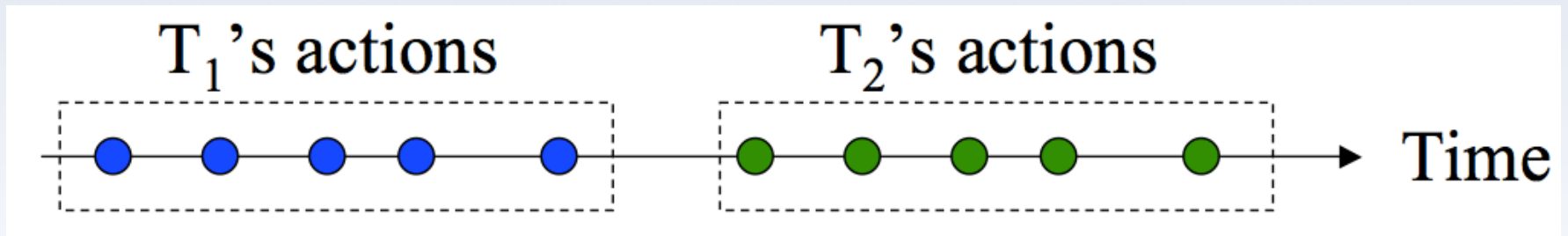
Example

- Our running example will have the constraint that $A=B$



Serial Schedule

- Execute one transaction first, then the other



Example

T ₁	T ₂	A	B
<div> READ(A, t) t := t + 100 WRITE(A, t) READ(B, t) t := t + 100 WRITE(B, t) </div>		25	25
		125	
			125
	<div> READ(A, s) s := s * 2 WRITE(A, s) READ(B, s) s := s * 2 WRITE(B, s) </div>		
		250	
			250



Example

T_1	T_2	A	B
		25	25
	READ(A, s)		
	$s := s * 2$		
	WRITE(A, s)	50	
	READ(B, s)		
	$s := s * 2$		
	WRITE(B, s)		50
READ(A, t)			
$t := t + 100$			
WRITE(A, t)		150	
READ(B, t)			
$t := t + 100$			
WRITE(B, t)			150



Concurrent Execution

- Improves throughput
 - CPU actions and I/O actions can happen in parallel
- Reduces average waiting time
 - Short transactions don't have to wait for long transactions to complete
- Just like with operating systems



Serializable Schedule

- A schedule is *serializable* if it has the same effect as some serial schedule



Example

T_1	T_2	A	B
READ(A, t)		25	25
$t := t + 100$			
WRITE(A, t)	READ(A, s)	125	
	$s := s * 2$		
	WRITE(A, s)	250	
READ(B, t)			
$t := t + 100$			
WRITE(B, t)			125
	READ(B, s)		
	$s := s * 2$		
	WRITE(B, s)		250



Example

T_1	T_2	A	B
		25	25
READ(A, t)			
$t := t + 100$			
WRITE(A, t)		125	
	READ(A, s)		
	$s := s * 2$		
	WRITE(A, s)	250	
	READ(B, s)		
	$s := s * 2$		50
	WRITE(B, s)		
READ(B, t)			
$t := t + 100$			
WRITE(B, t)			150



Notation

- We don't care about values
- Just ***what*** is accessed, ***who*** accessed it, and ***how*** it is accessed
 - Read: $r_i(X)$ transaction i reads element X
 - Write: $w_i(X)$ transaction i writes element X



Example

T1: $r_1(A)$; $w_1(A)$; $r_1(B)$; $w_1(B)$;

T2: $r_2(A)$; $w_2(A)$; $r_2(B)$; $w_2(B)$;

S: $r_1(A)$; $w_1(A)$; $r_2(A)$; $w_2(A)$; $r_1(B)$; $w_1(B)$; $r_2(B)$; $w_2(B)$;



Conflicting swaps

- A pair of consecutive actions conflict if changing order changes the behavior of a transaction:
 1. two actions from the same transaction
 2. $r_i(X), w_j(X)$ or $w_i(X), r_j(X)$
 3. $w_i(X), w_j(X)$



Conflict Serializable

- Making a series of nonconflicting swaps to a schedule, we can produce a serial schedule



Example

r1(A); w1(A); r2(A); w2(A); r1(B); w1(B); r2(B); w2(B);
r1(A); w1(A); r2(A); r1(B); w2(A); w1(B); r2(B); w2(B);
r1(A); w1(A); r1(B); r2(A); w2(A); w1(B); r2(B); w2(B);
r1(A); w1(A); r1(B); r2(A); w1(B); w2(A); r2(B); w2(B);
r1(A); w1(A); r1(B); w1(B); r2(A); w2(A); r2(B); w2(B);



Next Week...

- We'll investigate these notions more deeply
 - How to tell if a schedule is conflict-serializable
 - Mechanisms for enforcing serializability

