

## Database System Internals

# Optimistic Concurrency Control

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

February 19, 2020 CSE 444 - Winter 2020 1

1

## Announcements

- Quiz grades back this weekend on Gradescope
- Lab 3 part 1 due Tuesday
- HW 3 due date extended to Friday the 21st

February 19, 2020 CSE 444 - Winter 2020 2

2

## Pessimistic vs. Optimistic

- **Pessimistic CC** (locking)
  - Prevents unserializable schedules
  - Never abort for serializability (but may abort for deadlocks)
  - Best for workloads with high levels of contention
- **Optimistic CC** (timestamp, multi-version, validation)
  - Assume schedule will be serializable
  - Abort when conflicts detected
  - Best for workloads with low levels of contention

February 19, 2020 CSE 444 - Winter 2020 3

3

## Outline

- **Concurrency control by timestamps (18.8)**
- Concurrency control by validation (18.9)
- Snapshot Isolation

February 19, 2020 CSE 444 - Winter 2020 4

4

## Timestamps

- Each transaction receives unique timestamp  $TS(T)$

Could be:

- The system's clock
- A unique counter, incremented by the scheduler

February 19, 2020 CSE 444 - Winter 2020 5

5

## Timestamps

Main invariant:

The timestamp order defines the serialization order of the transaction

Will generate a schedule that is view-equivalent to a serial schedule, and recoverable

February 19, 2020 CSE 444 - Winter 2020 6

6

## Timestamps

With each element X, associate

- $RT(X)$  = the highest timestamp of any transaction U that read X
- $WT(X)$  = the highest timestamp of any transaction U that wrote X
- $C(X)$  = the commit bit: true when transaction with highest timestamp that **wrote** X committed

February 19, 2020 CSE 444 - Winter 2020 7

7

## Timestamps

With each element X, associate

- $RT(X)$  = the highest timestamp of any transaction U that read X
- $WT(X)$  = the highest timestamp of any transaction U that wrote X
- $C(X)$  = the commit bit: true when transaction with highest timestamp that **wrote** X committed

If transactions abort, we must **reset the timestamps**

February 19, 2020 CSE 444 - Winter 2020 8

8

## Main Idea

For any  $r_T(X)$  or  $w_T(X)$  request, check for conflicts:

- $w_U(X) \dots r_T(X)$
- $r_U(X) \dots w_T(X)$
- $w_U(X) \dots w_T(X)$

How do we check if Read too late?

Write too late?

February 19, 2020 CSE 444 - Winter 2020 9

9

## Main Idea

For any  $r_T(X)$  or  $w_T(X)$  request, check for conflicts:

- $w_U(X) \dots r_T(X)$
- $r_U(X) \dots w_T(X)$
- $w_U(X) \dots w_T(X)$

How do we check if Read too late?

Write too late?

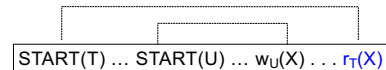
When T requests  $r_T(X)$ , need to check  $TS(U) \leq TS(T)$

February 19, 2020 CSE 444 - Winter 2020 10

10

## Read Too Late

- T wants to read X

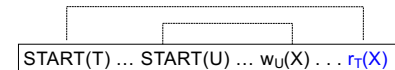


February 19, 2020 CSE 444 - Winter 2020 11

11

## Read Too Late

- T wants to read X



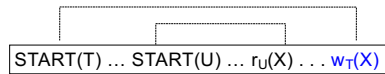
If  $WT(X) > TS(T)$  then need to rollback T !  
T tried to read **too late**

February 19, 2020 CSE 444 - Winter 2020 12

12

## Write Too Late

- T wants to write X



February 19, 2020 CSE 444 - Winter 2020 13

13

## Write Too Late

- T wants to write X



If  $RT(X) > TS(T)$  then need to rollback T !  
T tried to write **too late**

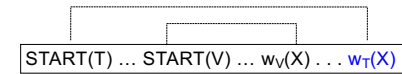
February 19, 2020 CSE 444 - Winter 2020 14

14

## Thomas' Rule

But... we can still handle it in one case:

- T wants to write X



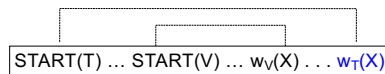
February 19, 2020 CSE 444 - Winter 2020 15

15

## Thomas' Rule

But we can still handle it:

- T wants to write X



If  $RT(X) \leq TS(T)$  and  $WT(X) > TS(T)$   
then don't write X at all !

Why does this work?

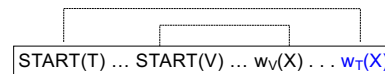
February 19, 2020 CSE 444 - Winter 2020 16

16

## Thomas' Rule

But we can still handle it:

- T wants to write X



If  $RT(X) \leq TS(T)$  and  $WT(X) > TS(T)$   
then don't write X at all !

Why does this work?

View-serializable:  
V will have overwritten T!

February 19, 2020 CSE 444 - Winter 2020 17

17

## View-Serializability

- By using Thomas' rule we do obtain a view-serializable schedule

February 19, 2020 CSE 444 - Winter 2020 18

18

## Summary So Far

Only for transactions that do not abort  
Otherwise, may result in non-recoverable schedule

Transaction wants to **READ** element X  
If  $WT(X) > TS(T)$  then ROLLBACK  
Else READ and update  $RT(X)$  to larger of  $TS(T)$  or  $RT(X)$

Transaction wants to **WRITE** element X  
If  $RT(X) > TS(T)$  then ROLLBACK  
Else if  $WT(X) > TS(T)$  ignore write & continue (Thomas Write Rule)  
Otherwise, WRITE and update  $WT(X) = TS(T)$

February 19, 2020 CSE 444 - Winter 2020 19

19

## Ensuring Recoverable Schedules

Recall:

- Schedule avoids cascading aborts if whenever a transaction reads an element, then the transaction that wrote it must have **already committed**
- Use the commit bit  $C(X)$  to keep track if the transaction that **last wrote** X has committed (just a read will not change the commit bit)

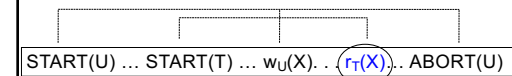
February 19, 2020 CSE 444 - Winter 2020 20

20

## Ensuring Recoverable Schedules

Read dirty data:

- T wants to read X, and  $WT(X) < TS(T)$
- Seems OK, but...



If  $C(X)=\text{false}$ , T needs to wait for it to become true

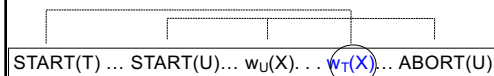
February 19, 2020 CSE 444 - Winter 2020 21

21

## Ensuring Recoverable Schedules

Thomas' rule needs to be revised:

- T wants to write X, and  $WT(X) > TS(T)$
- Seems OK not to write at all, but ...



If  $C(X)=\text{false}$ , T needs to wait for it to become true

February 19, 2020 CSE 444 - Winter 2020 22

22

## Timestamp-based Scheduling

- When a transaction T requests  $r_T(X)$  or  $w_T(X)$ , the scheduler examines  $RT(X)$ ,  $WT(X)$ ,  $C(X)$ , and decides one of:

- To grant the request, or
- To rollback T (and restart with later timestamp)
- To delay T until  $C(X) = \text{true}$

February 19, 2020 CSE 444 - Winter 2020 23

23

## Timestamp-based Scheduling

RULES including commit bit

- There are 4 long rules in Sec. 18.8.4
- You should be able to derive them yourself, based on the previous slides
- Make sure you understand them !

READING ASSIGNMENT:  
Garcia-Molina et al. 18.8.4

February 19, 2020 CSE 444 - Winter 2020 24

24

## Timestamp-based Scheduling (sec. 18.8.4)

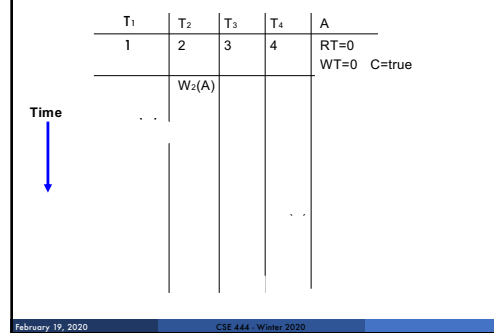
Transaction wants to READ element X  
 If  $WT(X) > TS(T)$  then ROLLBACK  
 Else if  $C(X) = \text{false}$ , then WAIT  
 Else READ and update  $RT(X)$  to larger of  $TS(T)$  or  $RT(X)$

Transaction wants to WRITE element X  
 If  $RT(X) > TS(T)$  then ROLLBACK  
 Else if  $WT(X) > TS(T)$   
 Then If  $C(X) = \text{false}$  then WAIT  
 else IGNORE write (Thomas Write Rule)  
 Otherwise, WRITE, and update  $WT(X)=TS(T)$ ,  $C(X)=\text{false}$

February 19, 2020 CSE 444 - Winter 2020 25

25

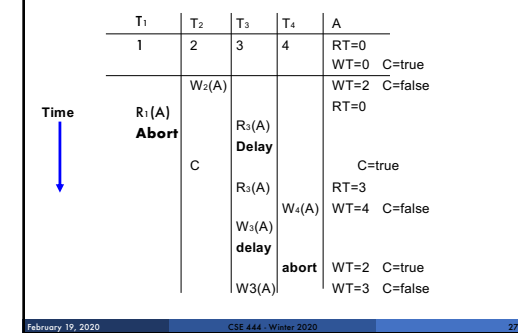
## Basic Timestamps with Commit Bit



February 19, 2020 CSE 444 - Winter 2020 26

26

## Basic Timestamps with Commit Bit



February 19, 2020 CSE 444 - Winter 2020 27

27

## Summary of Timestamp-based Scheduling

- View-serializable
- Avoids cascading aborts (hence: recoverable)
- Does NOT handle phantoms
  - These need to be handled separately, e.g. predicate locks

February 19, 2020 CSE 444 - Winter 2020 28

28

## Multiversion Timestamp

- When transaction T requests  $r(X)$  but  $WT(X) > TS(T)$ , then T must rollback

- Idea: keep multiple versions of X:  
 $X_t, X_{t-1}, X_{t-2}, \dots$   
 $TS(X_t) > TS(X_{t-1}) > TS(X_{t-2}) > \dots$

February 19, 2020 CSE 444 - Winter 2020 29

29

## Details

- When  $w_t(X)$  occurs,  
 if the write is legal then  
 create a new version, denoted  $X_t$  where  $t = TS(T)$

February 19, 2020 CSE 444 - Winter 2020 30

30

## Details

- When  $w_T(X)$  occurs,  
if the write is legal then  
create a **new version**, denoted  $X_t$  where  $t = TS(T)$
- When  $r_T(X)$  occurs,  
find **most recent version**  $X_t$  such that  $t \leq TS(T)$   
Notes:
  - $WT(X_t) = t$  and it never changes for that version
  - $RT(X_t)$  must still be maintained to check legality of writes
- Can delete  $X_t$  if we have a later version  $X_{t+1}$  and all active transactions  $T$  have  $TS(T) > t+1$

February 19, 2020

CSE 444 - Winter 2020

31

31

## Example (in class)

Four versions of X:  $X_3$   $X_9$   $X_{12}$   $X_{18}$ 

$R_6(X)$  -- Read  $X_3$   
 $W_{21}(X)$  -- Check read timestamp of  $X_{18}$   
 $R_{15}(X)$  -- Read  $X_{12}$   
 $W_5(X)$  -- Check read timestamp of  $X_3$

When can we delete  $X_3$ ?

February 19, 2020

CSE 444 - Winter 2020

32

32

## Example w/ Basic Timestamps

	$T_1$	$T_2$	$T_3$	$T_4$	A
Timestamps:	150	200	175	225	RT=0 WT=0
$R_1(A)$					RT=150
$W_1(A)$					WT=150
		$R_2(A)$			RT=200
		$W_2(A)$			WT=200
			$R_3(A)$		RT=225
			<b>Abort</b>		
				$R_4(A)$	RT=225

February 19, 2020

CSE 444 - Winter 2020

33

33

## Example w/ Multiversion

$T_1$	$T_2$	$T_3$	$T_4$	$A_0$	$A_{150}$	$A_{200}$
150	200	175	225			
$R_1(A)$				RT=150		
$W_1(A)$					Create RT=200	Create
	$R_2(A)$				RT=200	
	$W_2(A)$					
		$R_3(A)$				
		$W_3(A)$				
		<b>abort</b>				
			$R_4(A)$			RT=225

February 19, 2020

CSE 444 - Winter 2020

34

34

## Example w/ Multiversion

$T_1$	$T_2$	$T_3$	$T_4$	$A_0$	$A_{150}$	$A_{200}$
150	200	175	225			
$R_1(A)$				RT=150		
$W_1(A)$					Create RT=200	Create
	$R_2(A)$				RT=200	
	$W_2(A)$					
		$R_3(A)$				
		$W_3(A)$				
		<b>abort</b>				
			$R_4(A)$			RT=225

February 19, 2020

CSE 444 - Winter 2020

35

35

## Second Example w/ Multiversion

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$A_0$	$A_{150}$	$A_{200}$	$A_{175}$	$A_{225}$	$A_{250}$
1	2	3	4	5						
			$W_4(A)$							

February 19, 2020

CSE 444 - Winter 2020

36

36

T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	T <sub>4</sub>	T <sub>5</sub>	A <sub>0</sub>	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>
1	2	3	4	5						
W <sub>1</sub> (A)			W <sub>4</sub> (A)				Create		Create	
	R <sub>2</sub> (A)						RT=2			
	W <sub>2</sub> (A)						RT=3			
	abort									
		R <sub>3</sub> (A)								
				R <sub>5</sub> (A)					RT=5	
				W <sub>5</sub> (A)					RT=5	Create
			R <sub>4</sub> (A)							
R <sub>1</sub> (A)							RT=3			
C						X				
		C								

X means that we can delete this version

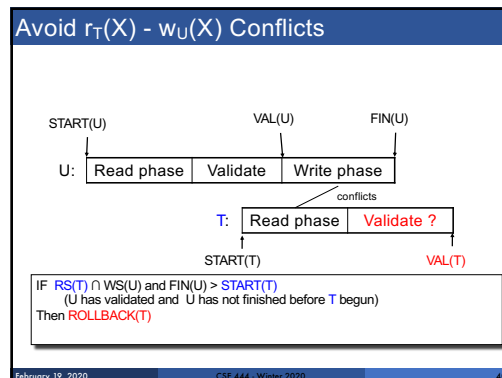
37

Outline										
<ul style="list-style-type: none"> <li>Concurrency control by timestamps (18.8)</li> <li>Concurrency control by validation (18.9)</li> <li>Snapshot Isolation</li> </ul>										

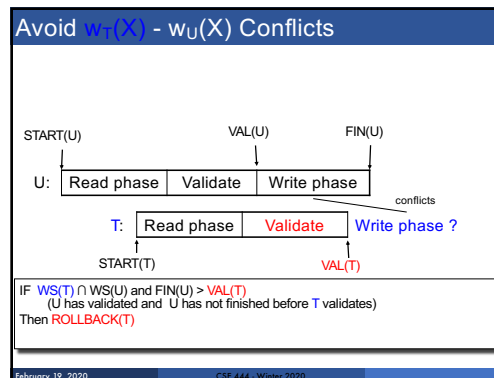
38

Concurrency Control by Validation										
<ul style="list-style-type: none"> <li>Each transaction T defines:             <ul style="list-style-type: none"> <li>Read set <math>RS(T)</math> = the elements it reads</li> <li>Write set <math>WS(T)</math> = the elements it writes</li> </ul> </li> <li>Each transaction T has three phases:             <ul style="list-style-type: none"> <li>Read phase; time = <math>START(T)</math></li> <li>Validate phase (may need to rollback); time = <math>VAL(T)</math></li> <li>Write phase; time = <math>FIN(T)</math></li> </ul> </li> </ul> <p>Main invariant: the serialization order is <math>VAL(T)</math></p>										

39



40



41

Outline										
<ul style="list-style-type: none"> <li>Concurrency control by timestamps (18.8)</li> <li>Concurrency control by validation (18.9)</li> <li>Snapshot Isolation             <ul style="list-style-type: none"> <li>Not in the book, but good overview in Wikipedia</li> </ul> </li> </ul>										

42

### Snapshot Isolation

- A type of multiversion concurrency control algorithm
- Provides yet another level of isolation
- Very efficient, and very popular
  - Oracle, PostgreSQL, SQL Server 2005
- Prevents many classical anomalies BUT...
- Not serializable (!), yet ORACLE and PostgreSQL use it even for SERIALIZABLE transactions!
  - But "serializable snapshot isolation" now in PostgreSQL

February 19, 2020 CSE 444 - Winter 2020 43

43

### Snapshot Isolation Overview

- Each transactions receives a timestamp  $TS(T)$
- Transaction T sees snapshot at time  $TS(T)$  of the database
- Write/write conflicts resolved by "first committer wins" rule
  - Loser gets aborted
- Read/write conflicts are ignored

February 19, 2020 CSE 444 - Winter 2020 44

44

### Snapshot Isolation Details

- Multiversion concurrency control:
  - Versions of X:  $X_{t1}, X_{t2}, X_{t3}, \dots$
- When T reads X, return  $X_{TS(T)}$ .
- When T writes X (to avoid lost update):
  - If latest version of X is  $TS(T)$  then **proceed**
  - Else if  $C(X) = \text{true}$  then **abort**
  - Else if  $C(X) = \text{false}$  then **wait**
- When T commits, write its updates to disk

February 19, 2020 CSE 444 - Winter 2020 45

45

### What Works and What Not

- No dirty reads (Why ?)
  - Start each snapshot with consistent state
- No inconsistent reads (Why ?)
  - Two reads by the same transaction will read same snapshot
- No lost updates ("first committer wins")
- Moreover: no reads are ever delayed
- However: read-write conflicts not caught !

February 19, 2020 CSE 444 - Winter 2020 46

46

### Write Skew

```

T1:  READ(X);
    if X >= 50
    then Y = -50; WRITE(Y)
    COMMIT

T2:  READ(Y);
    if Y >= 50
    then X = -50; WRITE(X)
    COMMIT
  
```

In our notation:

$R_1(X), R_2(Y), W_1(Y), W_2(X), C_1, C_2$

Starting with  $X=50, Y=50$ , we end with  $X=-50, Y=-50$ .  
Non-serializable !!!

February 19, 2020 CSE 444 - Winter 2020 47

47

### Write Skews Can Be Serious

- Acidicland had two viceroys, Delta and Rho
- Budget had two registers:  $taxes$ , and  $spendYng$
- They had high taxes and low spending...

```

Delta:  READ(taxes);
    if taxes = 'High'
    then ( spendYng = 'Raise';
          WRITE(spendYng) )
    COMMIT

Rho:    READ(spendYng);
    if spendYng = 'Low'
    then ( taxes = 'Cut';
          WRITE(taxes) )
    COMMIT
  
```

... and they ran a deficit ever since.

February 19, 2020 CSE 444 - Winter 2020 48

48



## Discussion: Tradeoffs

- **Pessimistic CC: Locks**
  - Great when there are many conflicts
  - Poor when there are few conflicts
- **Optimistic CC: Timestamps, Validation, SI**
  - Poor when there are many conflicts (rollbacks)
  - Great when there are few conflicts
- **Compromise**
  - READ ONLY transactions → timestamps
  - READ/WRITE transactions → locks

February 19, 2020 CSE 444 - Winter 2020 49

49

## Commercial Systems

Always check documentation!

- **DB2**: Strict 2PL
- **SQL Server**:
  - Strict 2PL for standard 4 levels of isolation
  - Multiversion concurrency control for snapshot isolation
- **PostgreSQL**: SI; recently: serializable SI (!)
- **Oracle**: SI

February 19, 2020 50

50