

MVCC: MULTIVERSION CONCURRENCY CONTROL

83

83

Multiversion Schemes

- Provide consistent version of database to transaction
 - Multiversion Timestamp Ordering
 - Multiversion Two-Phase Locking
- Use timestamps to label versions
- **Write**: Create new version of the data item
- **Read(A)**: Return appropriate version of A based on transaction timestamp
- **Reads** never have to wait!

84

84

Multiversion Timestamp Ordering

- Data item A has a sequence of versions $\langle A_1, A_2, \dots, A_m \rangle$
- Each version A_k contains three data fields:
 - Content: the value of version A_k .
 - W-timestamp(A_k): TS of transaction that created (wrote) version A_k
 - R-timestamp(A_k): Largest TS of transaction that read version A_k

85

85

Multiversion Timestamp Ordering

- When T_i creates a new version A_k of A
 - W-timestamp(A_k) := TS(T_i)
 - R-timestamp(A_k) := TS(T_i)
- When T_j reads A_k :
 - if TS(T_j) > R-timestamp(A_k)
R-timestamp(A_k) := TS(T_j)

86

86

Multiversion Timestamp Ordering

- Suppose T issues $\text{read}(A)$
- Versions $A_1, \dots, A_k, \dots, A_m$ ordered by $\text{W-TS}(A_i)$
- Pick A_k for some k such that
 - $\text{W-TS}(A_k) \leq \text{TS}(T)$
 - $\text{W-TS}(A_{k+1}) > \text{TS}(T)$
- Then return A_k

87

87

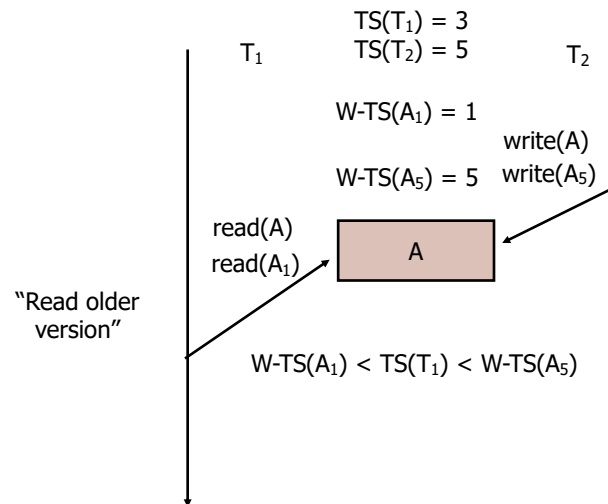
Multiversion Timestamp Ordering

- Suppose T issues $\text{write}(A)$
- Versions $A_1, \dots, A_k, \dots, A_m$ ordered by $\text{W-TS}(A_i)$
- Pick A_k for some k such that
 - $\text{W-TS}(A_k) \leq \text{TS}(T)$
 - $\text{W-TS}(A_{k+1}) > \text{TS}(T)$
- If $\text{TS}(T) < \text{R-timestamp}(A_k)$, then T is rolled back
- If $\text{TS}(T) = \text{W-timestamp}(A_k)$, then update A_k
- Otherwise create new version of A

88

88

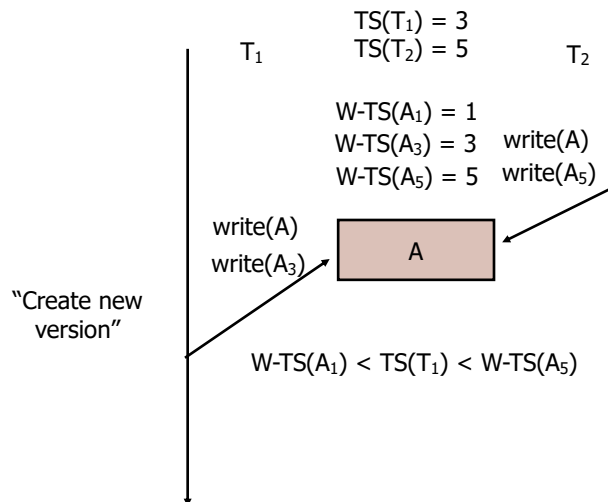
Multiversion Timestamp Ordering



89

89

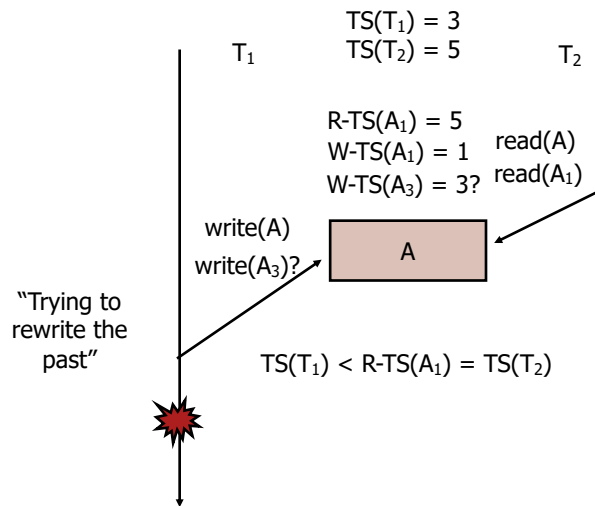
Multiversion Timestamp Ordering



90

90

Multiversion Timestamp Ordering



Multiversion Timestamp Ordering

- Reads always succeed!
- Write by T is rejected if already read by "future" transaction T_j
 - Future in serialization order
 - Written by T or "older" transaction
- Protocol guarantees serializability

92

Multiversion Two-Phase Locking

- Global version counter VC for commits
- Type 1: Update transactions
 - Acquire read and write locks
 - Keep locks until commit
 - Write: Create new version of A
 - Data item has single timestamp $TS(A)$
- Type 2: Read-only transactions
 - $TS(T)$ = current VC when it starts
 - T reads "newest" A_k with $TS(A_k) \leq TS(T)$

93

93

Multiversion 2PL: Update Transaction

- Read(Q)
 - Acquire read lock
 - Read latest version
- Write(Q)
 - Acquire write lock
 - Create new version A_{m+1} with $TS(A_{m+1}) = \infty$
- Commit Transaction
 - $TS(A_k) := VC$ for every variable written
 - $VC := VC + 1$

94

94

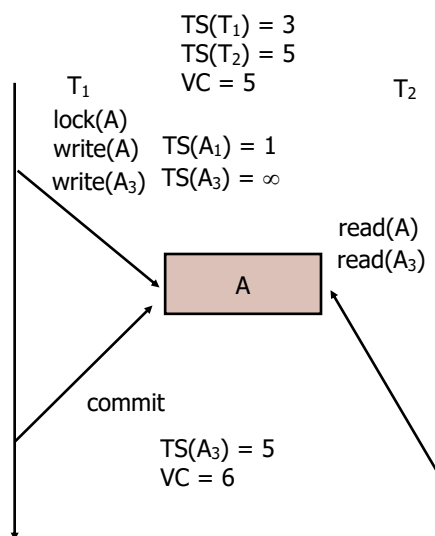
Multiversion 2PL: Read-Only Transaction

- Read-only transaction T_j
 - T_j reads "newest" A_k with $TS(A_k) \leq TS(T_j)$
 - Start before T : see the value before T updates
 - Start after T : see the values updated by T
- Only serializable schedules are produced

95

95

Multiversion 2PL



96

96

MVCC: Implementation Issues

- Creation of multiple versions increases storage overhead
 - Extra records
 - Version information
- Versions can be garbage collected

97

97

SNAPSHOT ISOLATION

98

98

Snapshot Isolation

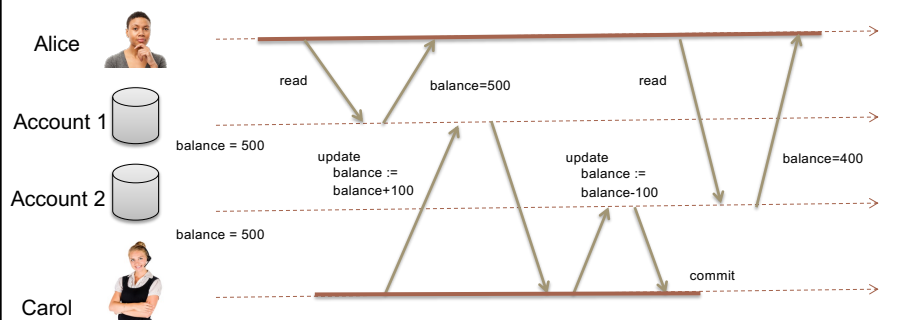
- Motivation: Large data query conflicts with OLTP transactions
- 2PL: Read locks block OLTP
- MVCC: Reads force rollback of OLTP

99

99

Weaker Isolation Levels

- Read Committed:
 - Read Skew: Reading different versions



100

100

Snapshot Isolation

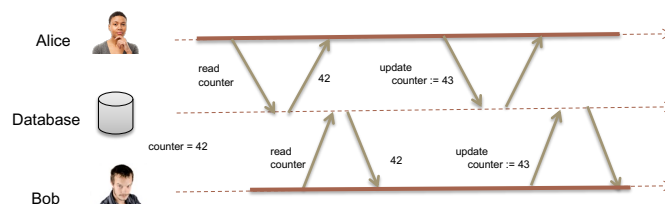
- Motivation: Large data query conflicts with OLTP transactions
- Solution 1: Give logical “snapshot” of database state to read only transactions
 - RW transactions use normal locking
 - Multiversion 2PL
 - How does system know a transaction is read only?

101

101

Snapshot Isolation

- Solution 2: Give snapshot of database state to **every** transaction
 - Updates alone use 2PL
 - Problem: anomalies such as lost update



- Partial solution: snapshot isolation level

102

102

Snapshot Isolation

- Transaction T
 - takes snapshot of committed data at start
 - reads/modifies data in its own snapshot
 - updates of *concurrent* transactions not visible to T
 - writes of T complete when it commits
- First-committer-wins rule:
 - Commit if no other concurrent transaction has already committed

103

103

Snapshot Read

- $X_0 = 100, Y_0 = 0$

T_1 : Deposit 50 in Y
rd $X_0 \rightarrow 100$
rd $Y_0 \rightarrow 0$

wr $Y_1 := 50$
rd $X_0 \rightarrow 100$
rd $Y_1 \rightarrow 50$

T_2 : Withdraw 50 from X

rd $Y_0 \rightarrow 0$
rd $X_0 \rightarrow 100$
wr $X_2 := 50$

- $X_2 = 50, Y_1 = 50$

104

104

Snapshot Write

- $X_0 = 100$

T_1 : Deposit 50 in X
rd $X_0 \rightarrow 100$

wr $X_1 := 150$
commit

T_2 : Withdraw 50 from X

rd $X_0 \rightarrow 100$

wr $X_2 := 50$

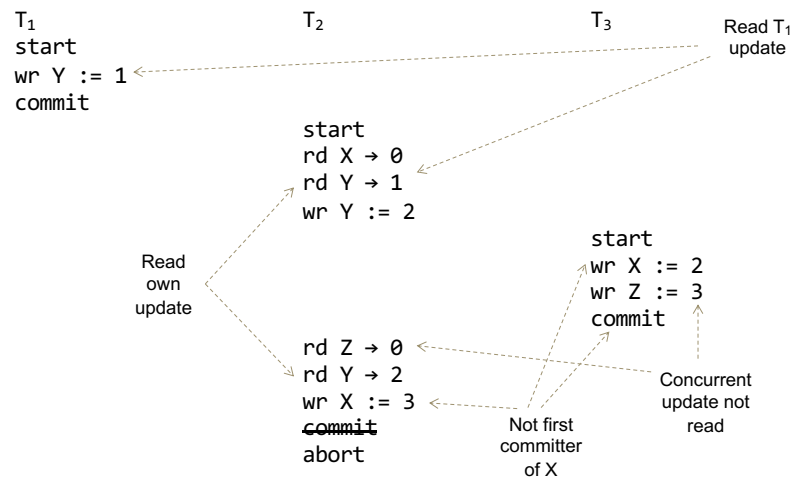
~~commit~~
abort

- $X_1 = 150$

105

105

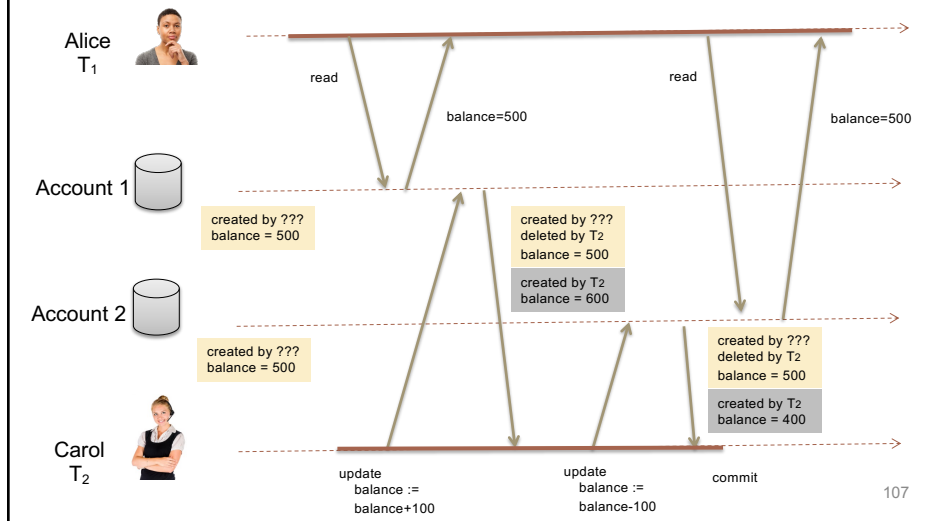
Snapshot Isolation



106

106

Implementing SI



107

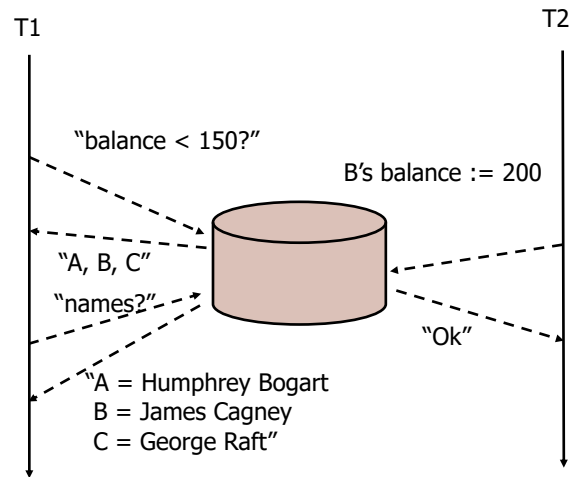
Benefits of SI

- Reading is never blocked
- Performance similar to Read Committed
- Avoids the usual anomalies
 - No dirty read
 - No lost update
 - No non-repeatable read
 - Predicate based selects are repeatable (no phantoms)

108

108

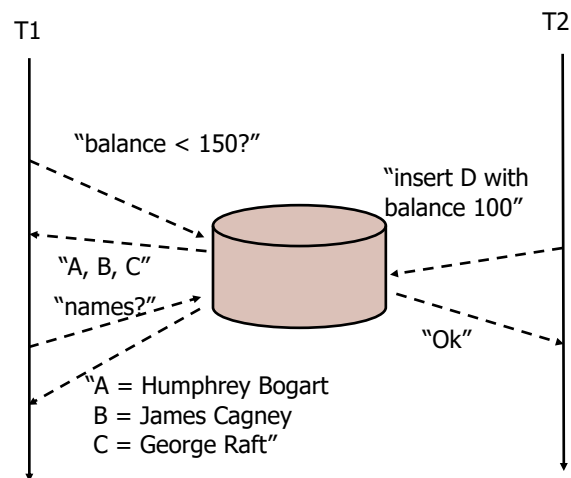
Unrepeatable Read



109

109

Phantom Read

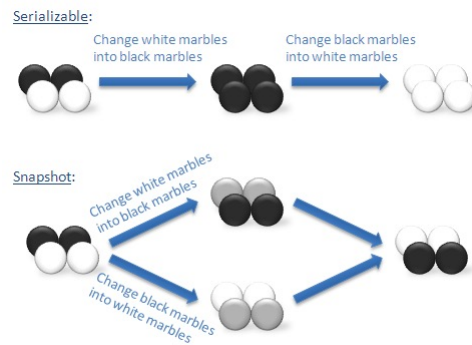


110

110

Problems with SI

- SI does not always give serializable executions



- Result: Integrity constraints can be violated

111

111

Write Skew

- Example
 - T1: $x := y$
 - T2: $y := x$
 - Initially $x = 3$ and $y = 17$
 - Serial execution: $x = ??, y = ??$
 - Snapshot isolation: $x = ??, y = ??$
- Example:
 - Find max order number among all purchase orders
 - Create a new purchase order with order number = previous max + 1

112

112

Write Skew and Phantoms

Alice



name	on_call
Alice	true
Bob	false
Carol	true

Carol



```
currently_on_call := (
  select count(*) from doctors
  where on_call = true
)
currently_on_call = 2
if (on_call > 1) {
  update doctors
  set on_call=false
  where name = "Alice"
}
commit
```

name	on_call
Alice	false
Bob	false
Carol	true

name	on_call
Alice	false
Bob	false
Carol	false

```
currently_on_call := (
  select count(*) from doctors
  where on_call = true
)
currently_on_call = 2
```

```
if (on_call > 1) {
  update doctors
  set on_call=false
  where name = "Carol"
}
commit
```

113

113

Write Skew and Phantoms

- Example: Meeting room booking system

```
BEGIN TRANSACTION;
-- Check for any existing bookings noon-1pm
SELECT COUNT(*) FROM bookings
  WHERE room_id = 123
     AND end_time > '2022-01-01 12:00'
     AND start_time < '2022-01-01 13:00';
-- If the previous query returned zero:
INSERT INTO bookings
  (room_id, start_time, end_time, user_id)
  VALUES (123, '2015-01-01 12:00', '2015-01-01 13:00', 666);
COMMIT;
```

114

114

Write Skew and Phantoms

- Pattern causing write skew:
 - Query based on some condition
 - INSERT, UPDATE, DELETE based on query
 - Phantom: Operation changes the result of query in another transaction
- What if query tests for absence in database, and operation does INSERT?
- Materializing conflicts
 - Locks just for query results

115

115

Snapshot Isolation and Correctness

- Execution may not be serializable but still correct!
- Example: reserving seats for a concert
 - Integrity constraint: a seat cannot be reserved by more than one person

116

116

Snapshot Isolation and Correctness

- Example: reserving seat X or Y for a concert

T_1 :	T_2 :
rd X \rightarrow free	
rd Y \rightarrow free	
	rd X \rightarrow free
	rd Y \rightarrow free
	wr Y := reserved
	commit
wr X := reserved	
commit	

117

117

Snapshot Isolation and Correctness

- Example: reserving seat X or Y for a concert

T_1 :	T_2 :
rd X \rightarrow free	
rd Y \rightarrow free	
	rd X \rightarrow free
	rd Y \rightarrow free
	wr X := reserved
	commit
wr X := reserved	
commit	
abort	

118


118

Preventing Lost Updates

- Atomic Writes

```
UPDATE counters SET value = value + 1 WHERE key = 'foo';
```

- Explicit Locking (SFU)

```
BEGIN TRANSACTION;  
SELECT * FROM figures  
WHERE name = 'robot' AND game_id = 222  
FOR UPDATE;   
-- Check whether move is valid, then update the position of a  
piece  
UPDATE figures SET position = 'c4' WHERE id = 1234;  
COMMIT;
```

119

119

Preventing Lost Updates

- Atomic Writes
- Explicit Locking
- Automatically detect lost updates
 - Necessary for snapshot isolation?
- Compare-and-set

```
-- May or may not be safe, depending on DB implementation  
-- (what if WHERE clause reads from snapshot)  
UPDATE wiki_pages SET content = 'new content'  
WHERE id = 1234 AND content = 'old content';
```

120

120

Preventing Lost Updates

- Conflict resolution and replication
 - Perform operation independently and resolve conflicts
 - Requires commutative operations
 - Otherwise Last Writer Wins (LWW) loses updates

121

121

Preventing Write Skew

- Atomic single-object operations
 - Multiple objects involved
- Automatic detection of lost updates
 - Not in practice
- Constraints
 - Requires multi-object constraints
- Explicitly lock rows
- Serializable isolation

122

122

Serializable Snapshot Isolation (SSI)

- Optimistic concurrency control
- Detect read of stale MVCC object
 - Have uncommitted writes committed by transaction end?
- Detect writes that affect prior reads
 - Use index locks to detect prior reads, but don't block

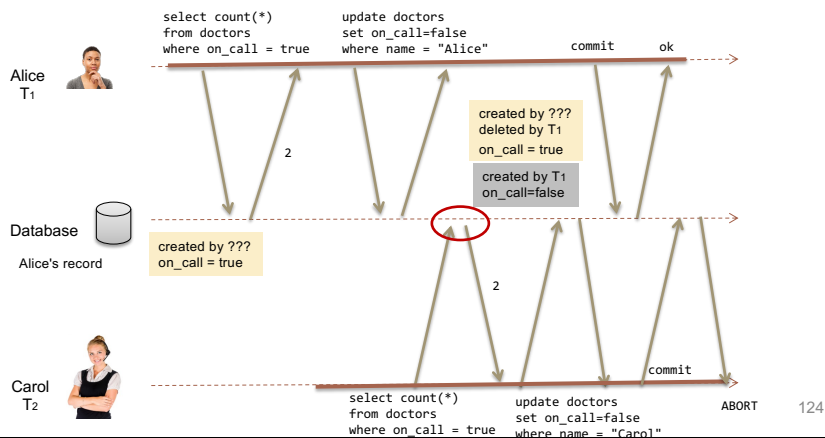
123

123

Serializability

name	on_call
Alice	true
Bob	false
Carol	true

- Serializable Snapshot Isolation: Detecting stale MVCC reads



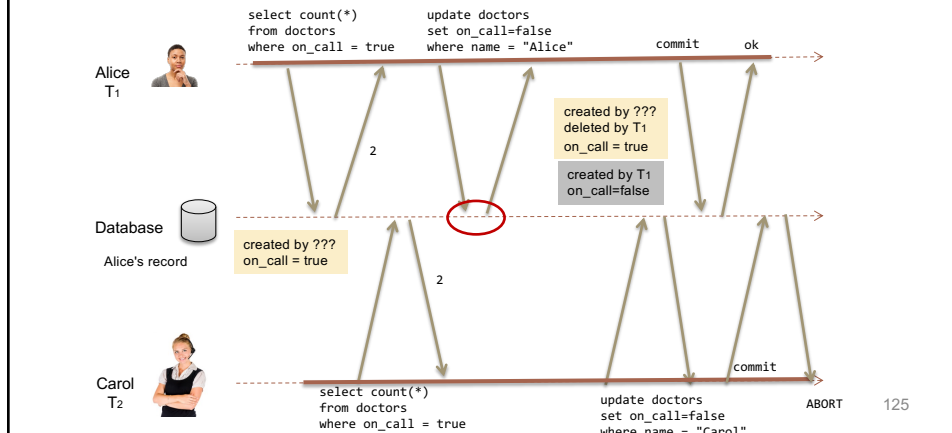
124

124

Serializability

name	on_call
Alice	true
Bob	false
Carol	true

- Serializable Snapshot Isolation: Detecting writes that affect prior reads



125

SI In Oracle and PostgreSQL

- Oracle
 - "Serializable" implemented by Snapshot Isolation!
 - "First updater wins"
- PostgreSQL:
 - <9.1: "Serializable" implemented by Snapshot Isolation (!)
 - >=9.1: "Serializable Snapshot Isolation" (SSI)

126

126