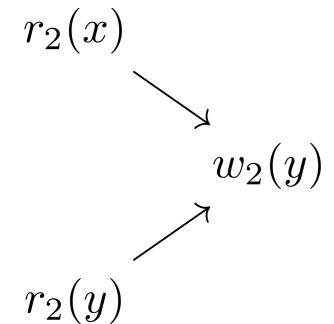
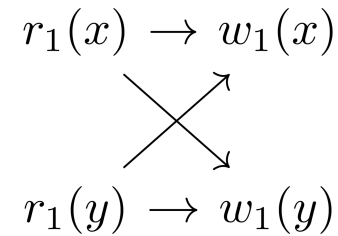


Data Management Algorithms
Consistency, Concurrency control, and Recovery
Consistency: Correctness Criteria

Boris Novikov

Transactions in Page Model

- Semantics of operations: $r(x)$ reads last preceding $w(x)$
- Transactions: partially ordered finite sets of operations
- Partial vs. total ordering



Histories and Schedules

A history H for finite set of transactions t_1, t_2, \dots, t_n :

- Operations $op(H) = \bigcup_i op(t_i) \cup (\{c_i | a_i\})$
- Partial ordering compatible with ordering in t_i : $p_i, q_i \in op(t_i) \quad p_i <_i q_i \Rightarrow p_i <_s q_i$
- For any x operations on x are ordered: $p_i(x) < q_j(x)$ or $q_j(x) < p_i(x)$
- Commit or abort follows all operations of the transaction: $p_i < c_i \vee p_i < a_i$

A schedule s is a prefix of H

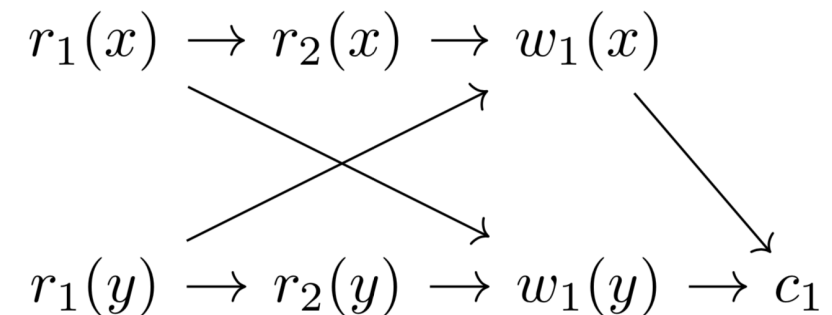
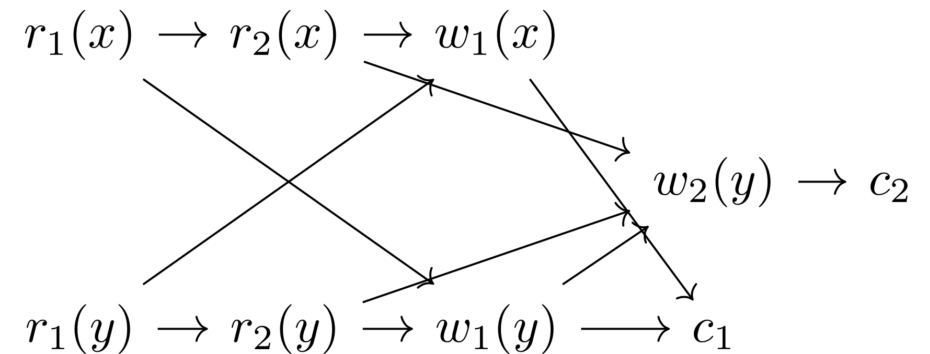
Histories and Schedules: Examples

- Partially ordered history

- Totally ordered history

$r_1(x) r_2(x) r_1(y) w_1(x) r_2(y) w_1(y) c_1 w_2(y) c_2$

- Schedule



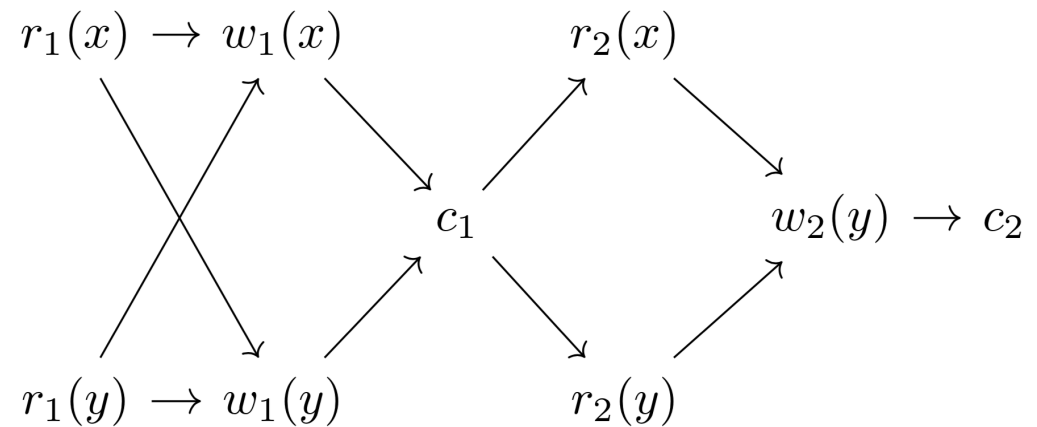
Anomalies of Concurrent Execution

- Lost update: $r_1(x) r_2(x) w_1(x) w_2(x) c_1 c_2$
- Inconsistent read (read skew): $r_1(x) r_2(x) w_1(x) w_1(y) r_2(y) c_1 c_2$
- Dirty read: $r_1(x) w_1(x) r_2(x) a_1 c_2$
- 20+ more
- Some of anomalies cannot be expressed in the page model

Serial Histories and Schedules

History is serial if:

- Transactions (but not necessarily operations) can be totally ordered:
 $t_1 < t_2 < \dots < t_n$
- If $t_i < t_j$ then $\forall p_i \forall q_j \quad p_i <_s q_j$
- A serial schedule is a prefix of a serial history
- Serial histories are correct. Why?



Transaction Sets for a Schedule

- All transactions: *trans(s)*
- Committed: *commit(s)*
- Aborted *abort(s)*
- Not yet completed: *active(s)*
-

Requirements for Correctness Criteria

- The class of correct schedules must be wide
- The correctness criteria should be computationally efficient
- Aborted transactions are not essential for correctness
-

The Approach

- Define an equivalence of schedules
- A schedule is correct iff it is equivalent to some serial schedule on the same set of transactions
- For any equivalence of schedules EQ, the correct schedules are also called **EQ-serializable**
- Several different equivalences will be considered

Herbrand Semantics: Informally

- Any read $r_i(x)$ returns the value written by the last preceding write $w_j(x)$ (executed by any transaction t_j)
- A value written by $W_i(x)$ depends on
 - *The functions of an application that issued the transaction (including any user input)
 - *Any read operations $r_i(y)\dots$ of the same transaction before $w_i(x)$

Herbrand Semantics for Operations

- Transaction t_0 writes the initial state of all pages x in the database

$$H_s(r_i(x)) = H_s(w_j(x)) \quad w_j(x) <_s r_i(x), \quad \forall w_k(x) \quad w_k(x) < r_i(x) \Rightarrow w_k(x) < w_j(x)$$

$$H_s(w_i(x)) = f_{ix}(r_i(y^1), \dots, r_i(y^m)) \quad r_i(y^k) < sw_i(x)$$

Where f_{ix} represent the application logic.

Herbrand Semantics for Schedules

- Herbrand universe is a set of expressions $HU = \{f_{ix}(e_1, \dots, e_m), e_k \in HU\}$ (this definition is recursive)
- Herbrand semantics for a schedule is a mapping from the database to Herbrand universe: $H(s) : D \rightarrow HU$

Final State Serializability - FSR

- Schedules s, s' are final-state equivalent iff $op(s) = op(s') \wedge H(s) = H(s')$.
- A schedule is final state serializable, $s \in \text{FSR}$, if it is FS-equivalent to a serial schedule

FSR Examples

$$r_1(x) \ r_2(y) \ w_1(y) \ r_3(z) \ w_3(z) \ r_2(x) \ w_2(z) \ w_1(x)$$

$$r_3(z) \ w_3(z) \ r_2(y)r_2(x) \ w_2(z) \ r_1(x) \ w_1(y) \ w_1(x)$$

are FS-equivalent

$$r_1(x) \ r_2(y) \ w_1(y) \ w_2(y) \ c_1 \ c_2 \ f_{2y}(f_{0y}())$$

$$r_1(x) \ w_1(y) \ r_2(y) \ w_2(y) \ c_1 \ c_2 \ f_{2y}(f_{1y}(f_{0y}()))$$

Are not equivalent

Reads_From

- Relation $\text{reads_From}(t_j, x, r_i)$ is true iff $w_j(x)$ is the last write operation preceding $r_i(x)$
- Relations *Reads_From* coincide for two schedules if Herbrand semantics for all operations coincide

View Equivalence

- FSR prevents lost updates, but does not prevent inconsistent reads
- All transactions in equivalent schedules should see the same state of the database
- Schedules are view equivalent if $H_s(r_i(x)) = H_{s'}(r_i(x))$ for all read operations of all transactions
- Actually Herbrand semantics coincide also for write operations

View Serializability - VSR

- Two schedules are VS-equivalent iff the relations *Reads_From* coincide,
- FSR \neq VSR.
- Proof: $w_1(x) r_2(x) r_2(y) w_1(y) c_1 c_2 \in \text{FSR} \setminus \text{VSR}$

VSR Properties

- Perfect from application semantics viewpoint
- Cannot be efficiently verified

Conflicts

Two operations of page model are inconflict if:

- Both process same page (implies the ordering)
- Belong to different transactions
- At least one of two is the write operation.

The first operation in the conflict always precedes the second in the schedule.

The set of all conflicts in a schedule s is $conf(s)$.

Conflict Serializability - CSR

- Schedules s, s' are conflict-equivalent iff $op(s) = op(s') \wedge conf(s) = conf(s')$
- A schedule is conflict-serializable ($s \in CSR$) if it is conflict-equivalent to a serial schedule.

CSR \subset VSR

CSR \subseteq VSR:

- If $r_i(x)$ reads from $w_j(x)$ in s and from $w_k(x)$ in s' then $(w_k(x), w_j(x)) \in \text{conf}(s)$
but $(w_j(x), w_k(x)) \in \text{conf}(s)$

CSR \neq VSR:

- $s = w_1(x) w_2(x) w_2(y) c_2 w_1(y) c_1 w_3(x) w_3(y) c_3$

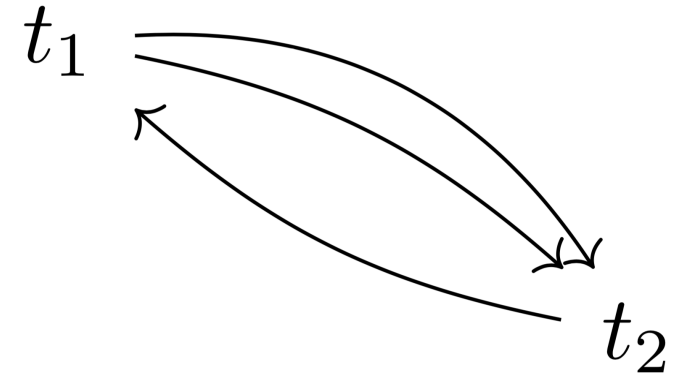
Conflict Graph

- Graph $G(s)$:

*nodes: $trans(s)$,

*edges: $conf(s)$

$r_1(x) \ r_1(y) \ w_2(x) \ w_2(y) \ w_1(x) \ c_1 \ c_2$



CSR Criterion

$s \in \text{CSR}$ iff $G(s)$ is acyclic

- Same conflicts produce same graph for original and equivalent serial. The graph for serial schedule is acyclic.
- If the graph is acyclic, the topological sort produces equivalent serial graph.

Commutativity

- Commutativity rules
 - C1: $r_i(x)r_j(x) \sim r_j(x)r_i(x) \quad i \neq j$
 - C2: $r_i(x)w_j(y) \sim w_j(y)r_i(x) \quad x \neq y, i \neq j$
 - C3: $w_i(x)w_j(y) \sim w_j(y)w_i(x) \quad x \neq y, i \neq j$
 - C4: $p_i(x), q_j(x)$ unordered $\Rightarrow p_i(x)q_j(y) \quad x \neq y \vee (p = r \wedge q = r)$

Commutativity Equivalence

- Equivalence: finite sequence of transformations C1 - C4 transforms $s \rightarrow s'$.
- A schedule s is commutativity-serializable iff $s \in \text{CSR}$.
- Operations commute iff they are not in conflict.
- The advantage is that semantics is hidden.

OCSR - Transaction Order Preserving CSR

$w_1(x) \ r_2(x) \ c_2 \ w_3(y) \ c_3 \ w_1(y) \ c_1$

- Serialization order: $t_3 \rightarrow t_1 \rightarrow t_2$,
- actual order $t_2 \ t_3$
- OCSR order of transactions is preserved
- $OCSR \subset CSR$
 - Proof?

COCSR: Commit Order Preserving CSR

- Definition; $s \in \text{COCSR}$ iff for any conflict $p_i < q_j \text{ and } q_j < p_i \Rightarrow c_i < c_j$
 \subset
- $\text{COCSR} \subset \text{CSR}$
- $\text{COCSR} \subset \text{OCSR}$

Relationships between CSR Subclasses

- $\text{COCSR} \subset \text{OCSR} \subset \text{CSR}$
- $s_1 = w_1(x) r_2(x) c_2 w_3(y) c_3 w_1(y) c_1 \in \text{CSR} \setminus \text{OCSR}$
- $s_2 = w_3(y) c_3 w_1(x) r_2(x) c_2 w_1(y) c_1 \in \text{OCSR} \setminus \text{COCSR}$
- $s_3 = w_3(y) c_3 w_1(x) r_2(x) w_1(y) c_1 c_2 \in \text{COCSR} \setminus \{\text{serial}\}$

Aborts and System Failures

- Transactions in *active(s)* can be aborted
- Correctness should not depend on aborts
- Any prefix of a correct schedule should be correct in the same sense

Prefix and Commit Closeness

- FSR is not prefix-closed
- VSR, CSR are prefix-closed
- $CP(s)$ is a projection of a schedule s consisting of transactions in $commot(s)$
- A class of schedules X is commit-closed if $CP(s') \in X$ for any prefix s' of $s \in X$.
- $CMCSR \subset CMVSR \subset CMFSR$
- $CMFSR \subset FSR$, $CMVSR \subset VSR$, $CVCSR = CSR$

Serializability Summary

$$s_1 = w_1(x) w_2(x) w_2(y) c_2 w_1(y) c_1 \notin \text{FSR}$$

$$s_2 = w_1(x) r_2(x) w_2(y) c_2 r_1(y) w_1(y) c_1 w_3(x) w_3(y) c_3 \\ \in \text{FSR} \setminus \text{CMFSR}$$

$$s_3 = w_1(x) r_2(x) w_2(y) w_1(y) c_1 c_2 \\ \in \text{CMFSR} \cup \text{VSR} \setminus \text{CVVSR}$$

$$s_4 = w_1(x) w_2(x) w_2(y) c_2 w_1(y) c_1 w_3(x) w_3(y) c_3 \\ \in \text{VSR} \setminus \text{CMFSR}$$

$$s_5 = w_1(x) r_2(x) w_2(y) w_1(y) c_1 c_2 w_3(x) w_3(y) c_3 \\ \in \text{CMFSR} \setminus \text{VSR}$$

$$s_6 = w_1(x) w_2(x) w_2(y) c_2 w_1(y) w_3(x) w_3(y) c_3 w_1(z) c_1 \\ \in \text{CMVSR} \setminus \text{CSR}$$

$$s_7 = w_1(x) w_2(x) w_2(y) c_2 w_1(z) c_1 \\ \in \text{CSR} \setminus \text{OCSR}$$

$$s_8 = w_3(y) c_3 w_1(x) r_2(x) c_2 w_1(y) c_1 \\ \in \text{OCSR} \setminus \text{COCSR}$$

$$s_9 = w_3(y) c_3 w_1(x) r_2(x) w_1(y) c_1 c_2 \\ \in \text{COCSR} \setminus \{\text{serial}\}$$

$$s_{10} = w_1(x) w_1(y) c_1 w_2(x) w_2(y) c_2 \in \{\text{serial}\}$$

Snapshot Isolation - SI

- STS(t), CTS(t) - transaction start and commit timestamps
- $E(t) = (STS(t), CTS(t))$ - transaction execution time interval
- RS(t), WS(t) - transaction read set and write sets

SI rules:

- ◆ SI1 Any read operation on t returns values committed before STS(t)
- ◆ SI2 $E(t_1) \cap E(t_2) = \emptyset \vee WS(t_1) \cap WS(t_2) = \emptyset$

SI and Serializability

- $r_1(x) w_1(x) r_2(x) w_2(x) c_2 r_1(y) w_1(y) c_1 \in \text{CSR} \setminus \text{SI}$
- $r_1(x) r_2(y) w_1(y) w_2(x) c_1 c_2 \in \text{SI} \setminus \text{FSR}$

SI Anomalies

- Write skew $r_1(x) r_2(y) w_1(y) w_2(x) c_1 c_2$
- Read-only transaction anomaly:
 $r_2(x) r_2(y) r_1(y) w_1(y) c_1 r_3(x) r_3(y) c_3 w_2(x) c_2$

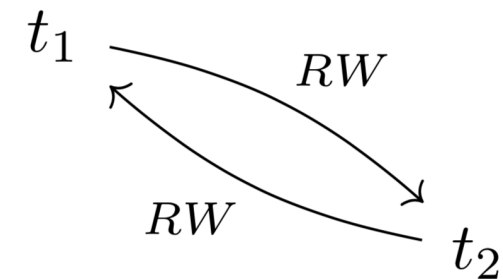
Dependencies

- WR $t_1 \rightarrow t_2 : w_1(x) \rightarrow r_2(x)$
- WW $t_1 \rightarrow t_2 : w_1(x) \rightarrow w_2(x)$ - excluded by SI2
- RW $t_1 \rightarrow t_2 : r_1(x) \rightarrow w_2(x)$

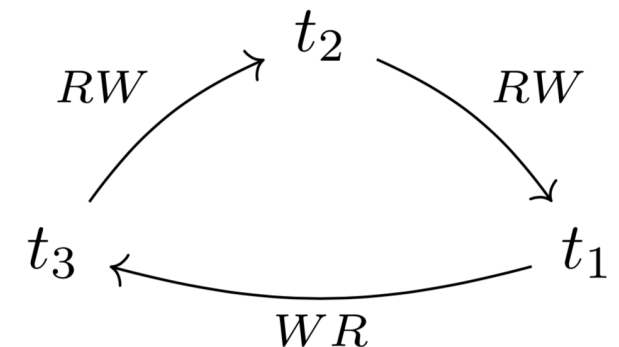
Serializable SI - SSI

- A schedule in SSI iff dependency graph is acyclic
- Only two kinds of cycles may exist under SI
- Checking for 2 edges: fast but false positives

Write skew:



Read-only:



Relaxed Isolation

- READ UNCOMMITTED
- READ COMMITTED
- REPEATABLE READ
- SERIALIZABLE