

Guard Independence and Constraint-Preserving Snapshot Isolation

Stephen J. Hegner
Umeå University
Department of Computing Science
SE-901 87 Umeå, Sweden
hegner@cs.umu.se
<http://www.cs.umu.se/~hegner>

FoIKS 2014
Eighth International Symposium on
Foundations of Information and Knowledge Systems
Bordeaux, France
3 March 2014

The ACID Characterization

- The properties which a set of concurrent transactions should exhibit is often expressed via the acronym *ACID*:

Atomicity: For each transaction, either the complete result of its execution is recorded in the database, or else nothing about its results is recorded.

Consistency: The execution of any transaction in isolation preserves the integrity of the database.

Isolation: The execution of one running transaction must not affect the execution of another concurrently running transaction.

Durability: The results of the transactions are permanent in the database.

- *Atomicity*, *consistency*, and *durability* are properties of individual transactions and their execution.
- *Isolation* describes how transactions interact.
- The focus of this talk is *isolation* of atomic, consistent, and durable transactions.

Serial Execution of Transactions

Serial execution: A set of transactions runs *serially* if there is no temporal overlap in their operations.

- Serial execution is considered to define optimal isolation, even though the result may depend upon the order of execution.

T_1	T_2	x
Read $\langle x \rangle$		10000
Cpd $\langle x, 10\% \rangle$		10000
Write $\langle x \rangle$		11000
	Read $\langle x \rangle$	11000
	Wd $\langle x, 2000 \rangle$	11000
	Write $\langle x \rangle$	9000

T_1	T_2	x
	Read $\langle x \rangle$	10000
	Wd $\langle x, 2000 \rangle$	10000
	Write $\langle x \rangle$	8000
Read $\langle x \rangle$		8000
Cpd $\langle x, 10\% \rangle$		8000
Write $\langle x \rangle$		8800

- The operations *Cpd* = *compound* and *Wd* = *withdraw* operate internally and do not write the database.

Lost Updates

- If the steps of the transactions are interleaved in certain ways, isolation may be lost.
- One symptom of poor isolation is *lost updates*.

T_1	T_2	x
Read $\langle x \rangle$		10000
Cpd $\langle x, 10\% \rangle$		10000
	Read $\langle x \rangle$	10000
	Wd $\langle x, 2000 \rangle$	10000
	Write $\langle x \rangle$	8000
Write $\langle x \rangle$		11000

T_1	T_2	x
	Read $\langle x \rangle$	10000
	Wd $\langle x, 2000 \rangle$	10000
Read $\langle x \rangle$		10000
Cpd $\langle x, 10\% \rangle$		10000
Write $\langle x \rangle$		11000
	Write $\langle x \rangle$	8000

- In the schedule on the left, the result of T_2 is lost.
- In the schedule on the right, the result of T_1 is lost.

The Model of Operations, Transactions, and Schedules

- Model the database schema as a set of updateable objects.

Object-level model of operations: There are two basic operations:

Read: $r_T\langle x \rangle$ denotes that transaction T reads data object x .

Write: $w_T\langle x \rangle$ denotes that transaction T writes data object x .

- In particular, the changes which T makes to x during a write are **not** modelled.
- A *transaction* is then modelled as a sequence of such operations:

Examples: T_1 : $r_{T_1}\langle x_1 \rangle w_{T_1}\langle x_1 \rangle r_{T_1}\langle x_2 \rangle w_{T_1}\langle x_2 \rangle$ T_2 : $r_{T_2}\langle x_1 \rangle r_{T_2}\langle x_3 \rangle w_{T_2}\langle x_3 \rangle w_{T_2}\langle x_1 \rangle$

- A *schedule* for a set of transactions is an intertwining of their operation sequences which preserves the local order for each transaction.

Examples: S_1 : $r_{T_1}\langle x_1 \rangle w_{T_1}\langle x_2 \rangle r_{T_1}\langle x_1 \rangle w_{T_1}\langle x_2 \rangle r_{T_2}\langle x_1 \rangle r_{T_2}\langle x_3 \rangle w_{T_2}\langle x_3 \rangle w_{T_2}\langle x_2 \rangle$
 S_2 : $r_{T_1}\langle x_1 \rangle w_{T_1}\langle x_1 \rangle r_{T_2}\langle x_1 \rangle r_{T_2}\langle x_3 \rangle w_{T_2}\langle x_3 \rangle r_{T_1}\langle x_2 \rangle w_{T_1}\langle x_2 \rangle w_{T_2}\langle x_2 \rangle$

- S_1 is a *serial* schedule for $\{T_1, T_2\}$, while S_2 is a non-serial schedule.

The Gold Standard for Isolation: View Serializability

Idea: A schedule is *view serializable* if it can be obtained by rearranging the operations of some serial schedule in such a way that:

- The read operations read from the same writer in each case (which might be the initial database state).
- The final writer of each data object is the same transaction in each case.
- Such a rearrangement does not change the final result of running the transactions.

Examples:

S_1 :	$r_{T_1}\langle x_1 \rangle w_{T_1}\langle x_2 \rangle r_{T_1}\langle x_1 \rangle w_{T_1}\langle x_2 \rangle r_{T_2}\langle x_1 \rangle r_{T_2}\langle x_3 \rangle w_{T_2}\langle x_3 \rangle w_{T_2}\langle x_2 \rangle$
S_2 :	$r_{T_1}\langle x_1 \rangle w_{T_1}\langle x_1 \rangle r_{T_2}\langle x_1 \rangle r_{T_2}\langle x_3 \rangle w_{T_2}\langle x_3 \rangle r_{T_1}\langle x_2 \rangle w_{T_1}\langle x_2 \rangle w_{T_2}\langle x_2 \rangle$
S_3 :	$r_{T_1}\langle x_1 \rangle r_{T_2}\langle x_1 \rangle w_{T_1}\langle x_1 \rangle r_{T_2}\langle x_3 \rangle w_{T_2}\langle x_3 \rangle r_{T_1}\langle x_2 \rangle w_{T_1}\langle x_2 \rangle w_{T_2}\langle x_2 \rangle$
S_4 :	$r_{T_1}\langle x_1 \rangle w_{T_1}\langle x_1 \rangle r_{T_2}\langle x_1 \rangle r_{T_2}\langle x_3 \rangle w_{T_2}\langle x_3 \rangle r_{T_1}\langle x_2 \rangle w_{T_2}\langle x_2 \rangle w_{T_1}\langle x_2 \rangle$

- S_1 and S_2 are view serializable.
- S_3 is not view serializable (changed read).
- S_4 is not view serializable (changed final write).

Guaranteeing View Serializability — SS2PL

System requirement: Need a scheduling algorithm which guarantees view-serializable schedules, not just a test for view serializability.

Strong strict two-phase locking (SS2PL): A lock-based solution.

- *Shared* (read) locks and *exclusive* (write) locks are required for all data access.
- Locks may be acquired at any time.
- All locks held until the transaction commits (ends).

Severe drawback; The locking requirements greatly limit concurrency.

- Querying on a non-indexed attribute would require locking the entire table until the end of the transaction!

Incorrect claim: Many DBMS textbooks incorrectly assert that SS2PL is widely used in practice to realize serializable isolation.

- Unknown to many users, the SQL `SERIALIZABLE` mode of isolation does **not** provide view serializability in many systems (e.g., Oracle).
- Even those systems which do implement SS2PL, it is not widely used due to poor performance.

Levels of Isolation of Transactions

Question: Isn't view serializability necessary to guarantee correct results?

Answer: That depends upon what is meant by "correct".

- Isolation is a matter of degree.

Real-world fact: Lower levels of isolation are used routinely.

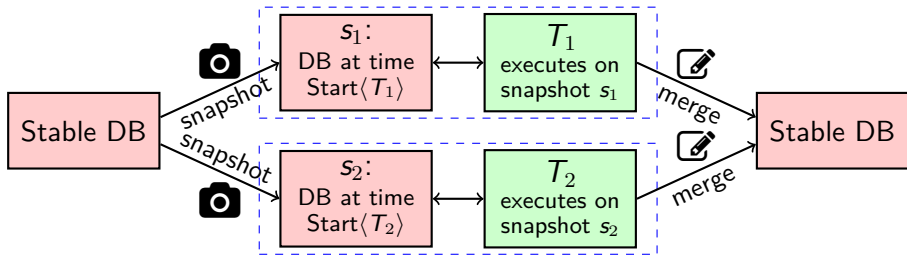
- The default level of isolation in many real systems is *read committed*, which guarantees that only committed data are read, but little more.
- Many transactions can tolerate such lower isolation levels without suffering serious consequences.
- The highest level, *view serializable*, is used only where absolutely essential, such as in financial transactions.

Multiversion Concurrency Control

MVCC: Most modern DBMSs employ *multiversion concurrency control*.

- There may be several *versions* of a given data object x .
- Rather than requiring locks, concurrency is achieved by allowing distinct transactions to operate on distinct versions of x .
- Differences must eventually be resolved, but typically not at the expense of long waits.
- In general, MVCC supports far more concurrency than single-version, lock-based approaches.
- One of the most common approaches within MVCC for achieving a high level of isolation with substantial concurrency is called *snapshot isolation*.
- Because the approach of this research is based upon it, it is worth a closer look.

Snapshot Isolation



- In *snapshot isolation (SI)*, each transaction operates on a *snapshot*:
 - a (private) copy of the database with values taken at the point in time at which the transactions begins.

First Committer Wins (FCW): T_i is allowed to commit its local writes to the stable DB only if no data object x which it writes has been committed, since its snapshot was created, to the stable DB by another transaction.

- Otherwise, it must abort and start over.

Advantages of Snapshot Isolation

- SI has some very attractive properties.

High Level of Isolation: Since each transaction operates on a private copy, isolation is achieved at what appears to be at a relatively high level.

Enhanced concurrency: No locks \Rightarrow writers do not block readers.

- Readers (almost) never have to wait for writers to finish.
- The attainable level of concurrency is far greater than that of SS2PL.
- For these reasons, SI is widely used in practice.

⚠ Real systems use *first updater wins (FUW)*, and there may be some blocking when foreign-key constraints are checked, but these are details which do not distort the main conclusions.

Question: Does SI provide serializable-level isolation?

Answer: That depends upon the definition of *serializable*.

Write Skew — Constraint Violation under SI

Fact: SI does not guarantee view-serializable isolation. \square

Example (write skew): x and y represent the balances of two accounts.

Integrity constraint: $x + y \geq 500\text{€}$ **Initial state:** $x = 300\text{€}$, $y = 300\text{€}$

T_1 : Withdraw 100€ from x

T_2 : Withdraw 100€ from y .

- Assume that these transactions run concurrently under SI.
- Each transaction run in isolation satisfies the integrity constraint.
- The final state is $(x, y) = (200\text{€}, 200\text{€})$, which violates the constraint.
- With serial execution, the second transaction will fail.
- Thus, SI does not guarantee view serializability.

The SQL Standard and Serializability

❖ SI satisfies the conditions set forth in the SQL standard for the SERIALIZABLE isolation level.

- The standard **defines** serializability as the absence of three types of transaction anomalies.

Apparent reason: The architects of the standard could not think of any nonserializable behavior which could arise in the absence of violations of those anomalies.

Consequence: Real systems are free to implement the SERIALIZABLE level of isolation as SI, and several do so.

- Unfortunately, many users mistakenly believe that SERIALIZABLE isolation in SQL must mean view serializable.

Opinion/Rant: The definition of SERIALIZABLE in the SQL standard is a poster child for why good theory is a necessary part of even the most practical endeavors.

Serializable Snapshot Isolation

Serializable SI (SSI): Augment SI to achieve true view serializability.

DSG: The approach relies on properties of the *direct serialization graph*, whose vertices are committed transactions. The edges are as follows:

rw-edge: $T_i \xrightarrow{rw} T_j$ T_i reads a data object and T_j is the next writer of that object.

ww-edge and wr-edge: Similar; details not central here.

Dangerous structure in DSG: $T_i \xrightarrow{rw} T_j \xrightarrow{rw} T_k$ ($T_i = T_k$ possible) *occurring in a cycle* with $\{T_i, T_j\}$ and $\{T_j, T_k\}$ concurrent.

Theorem [Fekete et al 2005]: If a schedule for SI is not view serializable, the DSG must contain a dangerous structure. \square

Optimistic strategy: Serializable SI (SSI):

- It is too expensive to maintain the entire DSG.
- Look for *potential* dangerous structures (need not be part of a cycle) and require one transaction to terminate to preserve serializability.
- This requires testing only three transactions at a time.
- But there will be false positives.

Serializable Snapshot Isolation — Practice and Limitations

Use in PostgreSQL: As of version 9.1, SSI is used to implement SERIALIZABLE isolation in PostgreSQL.

- Thus, SERIALIZABLE isolation is finally truly view serializability.
- Ordinary SI is still available as REPEATABLE READ isolation.
- Before version 9.1, both isolation levels were implemented as SI.

Question: Why is there a need for anything more?

Answers:

- SSI results in more false positives (with consequent aborts and reruns) than does ordinary SI.
- For some transaction mixes (particularly interactive and long-running), this may be a significant or severe drawback.

Question: Is there something in between SI and SSI?

Answer: Yes, *constraint-preserving SI (CPSI)*, the topic of this research.

- Ensures that constraints will be satisfied (no write skew).
- Much simpler algorithm with no false positives.

Permutation – Nonserializability without Constraint Violation

Example (SI permutation): $n \in \mathbb{N}$;

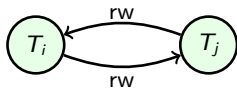
- d_0, d_1, \dots, d_{n-1} data objects.
 - $\tau_0, \tau_1, \dots, \tau_{n-1}$ transactions with $\tau_i: d_i \leftarrow d_{(i+1) \bmod n}$.
 - The n transactions, run concurrently under SI, effect a permutation of the values of the d_i 's (shift to the left).
-
- $\tau_i \xrightarrow{rw\langle d_i \rangle} \tau_{(i+1) \bmod n}$ denotes that τ_i reads d_i and $\tau_{(i+1) \bmod n}$ writes it.
- This behavior cannot be view serializable since if τ_i is run first, the old value of d_i is lost.
 - However, if any transaction (say τ_i) is removed, the result of running all transactions concurrently under SI is serializable.
 - Run them in this order: $\tau_{i+1} \dots \tau_{n-1} \tau_0 \dots \tau_{i-1}$.

Observation: For any $n \in \mathbb{N}$, there is a set of n transaction which, when run concurrently under SI, results in nonserializable behavior, yet any proper subset produces serializable behavior under SI.

Constraint-Preserving Snapshot Isolation

Question: Can such large cycles also occur when the condition to be preserved is constraint integrity and not full view serializability?

Theorem: Let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions running under SI according to some schedule \mathbf{S} . If the DSG is free of two-vertex cycles of the form



then the result is guaranteed to satisfy all integrity constraints. \square

Comment: These include all integrity constraints, even those implemented via triggers, for example.

Comment: The DSG might have larger cycles containing dangerous structures, but only those of the form identified above can result in constraint violation.

Two Types of Reads under SI

- A finer-grained version of the main theorem may be obtained by distinguishing two types of reads under SI:

Example: Let the database schema have three data objects w , x , and y with the constraint $x + y \geq 500$.

- Transaction T defined by $x \leftarrow x - w$.
- y is the *guard* of the transaction; it must be read in order to verify that the update will satisfy the integrity constraint.
- w must be read only to determine the update; it is not used in the checking the integrity constraint.

The value of y when T commits is critical: If the value of the guard y of T is changed by another concurrent transaction, there is a risk that the constraint will be violated.

Only the snapshot value of w is important for constraint satisfaction: A change to the value of w by another concurrent transaction will not affect whether or not the constraint is satisfied.

The Guard of a Data Object

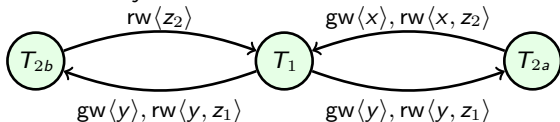
Guard of a transaction: The *guard* of a transaction T is the set of all data objects which must be read by T in order to verify the integrity constraints, but which are not written by T .

Example: Data objects: $\{x, y, z_1, z_2\}$; Constraint: $x + y \geq 500$.

Transaction	Write Set	Read Set	Guard Set
$T_1 : x \leftarrow x - z_1; z_2 \leftarrow z_2 - 10$	$\{x, z_2\}$	$\{y, z_1\}$	$\{y\}$
$T_{2a} : y \leftarrow y + z_2; z_1 \leftarrow z_1/2$	$\{y, z_1\}$	$\{x, z_2\}$	$\{x\}$
$T_{2b} : y \leftarrow y + z_2 ; z_1 \leftarrow z_1/2$	$\{y, z_1\}$	$\{z_2\}$	\emptyset

gw-edge $T_i \xrightarrow{gw} T_j$ in the DSG: T_j writes the guard of T_i .

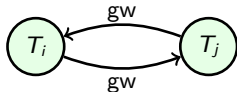
- $T_i \xrightarrow{gw} T_j \Rightarrow T_i \xrightarrow{rw} T_j$ but not conversely.



Note: T_{2a} and T_{2b} are alternatives; they cannot run concurrently.

Guard Independence

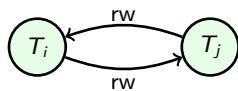
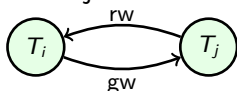
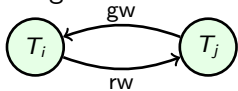
Guard independence of two transactions T_1 and T_2 is the formalization of the condition that a cycle of the form



does not exist.

Theorem: Let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions running under SI according to some schedule \mathbf{S} . If every pair of *concurrent* transactions is guard independent, then the result is guaranteed to satisfy all integrity constraints. Furthermore, the test is essentially free of false positives. \square

Remark: Cycles of the following three forms are allowed, as long as the rw-edges do not involve guard objects:

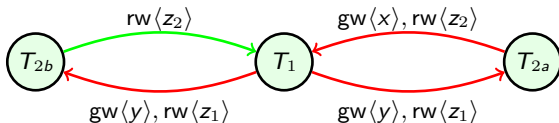


- These would identify a dangerous structure in SSI and result in the termination of one of the transactions.

Example of Guard Independence

Example: Data objects: $\{x, y, z_1, z_2\}$; Constraint: $x + y \geq 500$.

Transaction	Write Set	Read Set	Guard Set
$T_1 : x \leftarrow x - z_1; z_2 \leftarrow z_2 - 10$	$\{x, z_2\}$	$\{y, z_1\}$	$\{y\}$
$T_{2a} : y \leftarrow y + z_2; z_1 \leftarrow z_1/2$	$\{y, z_1\}$	$\{x, z_2\}$	$\{x\}$
$T_{2b} : y \leftarrow y + z_2 ; z_1 \leftarrow z_1/2$	$\{y, z_1\}$	$\{z_2\}$	\emptyset



Note: $rw\langle \alpha \rangle$ not shown if $gw\langle \alpha \rangle$ also holds for data object α on an edge.

- T_1 and T_{2b} are guard independent, while T_1 and T_{2a} are not.

Note: T_{2a} and T_{2b} are alternatives; they cannot run concurrently.

Theory — Write-Commuting Transactions

- The formal proof of the main result rests upon the notion of *write-commuting transactions*.

Definition: A transaction T is *consistent* if, after reading its initial snapshot data, its updates result in a database state which satisfies the integrity constraints.

Write Commute: Two transactions T_1 and T_2 *write commute* if they run concurrently and, whenever each is consistent (possibly reading different initial snapshots), they may commit in either order without violating any integrity constraints.

Observation: If all concurrent transactions write commute, there can be no constraint violations provide each transaction is consistent in isolation. \square

Theorem: Guard independence implies write commutativity. \square

Example of Write Commutativity

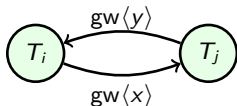
Example: Data objects: $\{x, y, z_1, z_2\}$; Constraint: $x + y \geq 500$.

Transaction	Write Set	Read Set	Guard Set
$T_1 : x \leftarrow x - z_1; z_2 \leftarrow z_2 - 10$	$\{x, z_2\}$	$\{y, z_1\}$	$\{y\}$
$T_{2a} : y \leftarrow y + z_2; z_1 \leftarrow z_1/2$	$\{y, z_1\}$	$\{x, z_2\}$	$\{x\}$
$T_{2b} : y \leftarrow y + z_2 ; z_1 \leftarrow z_1/2$	$\{y, z_1\}$	$\{z_2\}$	\emptyset

- T_1 and T_{2a} do not write commute:
 - Consider: $x = y = 300$, $z_1 = 100$, $z_2 = -100$ initially, and the two read the same snapshot.
- T_1 and T_{2b} write commute:
 - T_{2b} always succeeds, and it can only “improve” the likelihood of the constraint being satisfied after it commits.
 - Thus, if T_1 succeeds in isolation, there can be no constraint violation if it commits after T_{2b} .
 - Since T_1 and T_{2b} have disjoint write sets, the result of T_1 committing before T_{2b} is the same as committing afterwards.

Essentially Free of False Positives

Question: What does *essentially free of false positives* mean?



Answer: *Some* changes of the data objects x and y would result in a constraint violation.

- Since only an object-level model of operations is employed, information about which changes are made is not used to determine correctness.
- However, in contrast to the situation for ordinary SSI, the presence of other concurrent transactions is not necessary for a violation of the isolation condition to be possible.

Compound Data Objects as Views

Data objects as views: In the paper, data objects are modelled as *views* on the main schema.

- This provides a powerful mathematical framework for modelling the necessary ideas.
- Traditional *row objects* (a single row of a relation, identified by primary key) are modelled as the selection operation on that row.

Example: If R is a relation with single key attribute K , then $\sigma_{K=k_0}(R)$ is the view which models the data object whose key value is a .

- The advantage of the view model is that it also allows modelling of data objects other than view.

Example: If $R[ABC]$ is constrained by the FD $K \rightarrow AB$, then it is governed by the join dependency $\bowtie [KA, KB]$.

- In this case, the data objects defined by the projections $\pi_{KA}(\sigma_{K=k_0}(R))$ and $\pi_{KB}(\sigma_{K=k_0}(R))$ may be updated by separate transactions if only A and B , and not K , are modified.

Conclusions and Further Directions

Conclusions:

New Isolation Level: A new isolation level, *constraint-preserving snapshot isolation (CPSI)*, has been investigated.

SI < CPSI < Ser: It is at a strictly higher level than snapshot isolation, and a strictly lower level than view serializability.

- The test for adherence is much simpler than that for serializable snapshot isolation, with far less risk of false positives.

Further Directions:

Implementation and performance studies: It would be very useful to see how this approach fares in various situations.

Extension to a value-level model: Work is underway to extend the approach to a *value-level model*, in which the transaction manager has simple information about the nature of the updates which the transactions perform.

- This type of extension is critical for *interactive transactions*, in which abort and rerun is not an acceptable strategy for resolving conflicts.