

Addendum to “High-Performance Concurrency Control Mechanisms for Main-Memory Databases”

Per-Ake Larson¹, Spyros Blanas², Cristian Diaconu¹, Craig Freedman¹, Jignesh M. Patel², Mike Zwilling¹

¹Microsoft, ²University of Wisconsin – Madison

1. Single-version locking scheduler

Proving the single-version locking scheme correct is trivial, as the scheduler is a 2PL scheduler.

2. Multi-version pessimistic (locking) scheduler

The multi-version pessimistic (locking) scheme is in fact a MV2PL scheduler. Holding a certify (commit) lock on a data item in MV2PL is exactly like having the NoMoreReadLocks bit set in the latest version of the data item in our implementation (see Section 4.2.1). Section 5.5.2 of [WV02] describes MV2PL in detail and proves it only admits 1SR multi-version histories.

3. Multi-version optimistic scheduler

Let us now prove that the multi-version optimistic scheduler only admits 1SR multi-version histories. We use the notation and theorems from Section 5.2 of [BHG87]. The multi-version optimistic scheduler behaves like a MVTO scheduler, with the changes described below.

Let transaction Tx be a committed transaction with a Begin timestamp of TxBegin and an End timestamp of TxEnd.

Property 1: Timestamps are assigned in a monotonically increasing order, and each transaction has a unique begin and end timestamp, such that $TxBegin < TxEnd$.

Property 2: A given version is valid for the interval specified by the begin and end timestamps. There is a total order \ll of versions for a given datum, as determined by the timestamp order of the non-overlapping version validity intervals.

Property 3: The transaction Tx reads the latest committed version as of TxRead (where $TxBegin \leq TxRead < TxEnd$) and validates (that is, repeats) the read of the latest committed version as of TxEnd. The transaction fails if the two reads return different versions.

Property 4: Updates or deletes to a version V first check the visibility of V. Checking the visibility of V is equivalent to reading V. Therefore, a write is always preceded by a read: if transaction Tx writes Vnew, then transaction Tx has first read Vold, where $Vold \ll Vnew$. Moreover, there exists no version V such that $Vold \ll V \ll Vnew$, otherwise Tx would have never committed: it would have failed during the Active phase when changing the end timestamp of Vold (see Section 3.1, paragraph “Update version”)¹.

¹ Notice that all our concurrency control algorithms enforce a stronger property: they use the first-writer-wins rule to abort transactions that participate in a write-write conflict before it is determined whether the first writer will commit. The more relaxed property described here is sufficient to prove correctness.

Property 5: The transaction T_x logically writes at $T_x\text{End}$, because the version is invisible to other transactions until $T_x\text{End}$ (see Section 2.5, tables 1, 2 on visibility of versions where TB or TE contains the ID of transactions in the Active phase).

The multi-version serialization graph $MVSG(H, \ll)$ is a graph defined on a multi-version history H and version total order \ll . The MVSG has nodes for the committed transactions in H , and by definition, an edge $T_i \rightarrow T_j$ exists in the MVSG (where i, j, k are distinct), if and only if:

- A) T_i writes V_i and T_j reads V_i
or
- B) T_i writes V_i and T_k reads V_j , where $V_i \ll V_j$
or
- C) T_i reads V_k and T_j writes V_j , where $V_k \ll V_j$

Let us prove that every edge in the MVSG is ordered with respect to the end timestamp order of the transactions involved. That is, we will prove that any directed edge $T_i \rightarrow T_j$ will always point from a transaction T_i to a transaction T_j such that $T_i\text{End} < T_j\text{End}$.

For (A), let T_j read V_i at $T_j\text{Read} < T_j\text{End}$. If $T_j\text{Read} < T_i\text{End}$, from Properties 3 and 5, it would have been impossible to read V_i , as it's uncommitted. Therefore $T_i\text{End} < T_j\text{Read} < T_j\text{End}$, therefore the $T_i \rightarrow T_j$ edge is ordered with respect to the end timestamp, as $T_i\text{End} < T_j\text{End}$.

For (B), T_k reads V_j , therefore the read is preceded by T_j writing V_j and committing. Furthermore from Property 3, $T_j\text{End} < T_k\text{Read}$, or V_j wouldn't be visible at $T_k\text{Read}$. From Property 4, since T_j wrote V_j , that means that T_j read V_i (or any later version) at $T_j\text{Read}$, where $T_j\text{Read} < T_j\text{End}$. Therefore, from (A), $T_i\text{End} < T_j\text{Read} < T_j\text{End}$, so the $T_i \rightarrow T_j$ edge is ordered with respect to the end timestamp, as $T_i\text{End} < T_j\text{End}$.

For (C), let T_i read V_k at $T_i\text{Read}$, where $T_i\text{Read} < T_i\text{End}$. T_j writes V_j at $T_j\text{End}$. There are three cases:

- 1) $T_j\text{End} < T_i\text{Read} < T_i\text{End}$. This violates property 3, as V_j was a committed version at $T_i\text{Read}$, therefore T_i would have read V_j , not V_k . (See Section 2.5, "Version Visibility".)
- 2) $T_i\text{Read} < T_j\text{End} < T_i\text{End}$. This violates property 3, as T_i would repeat the read at $T_i\text{End}$ during validation, and would read V_j . However T_i read V_k at $T_i\text{Read}$, therefore T_i would fail validation and would never participate in the MVSG. (See Section 3.2, "Preparation Phase").
- 3) $T_i\text{Read} < T_i\text{End} < T_j\text{End}$. T_i would validate the read at $T_i\text{End}$, read V_k again, and commit. Therefore the edge $T_i \rightarrow T_j$ is ordered with respect to the end timestamp, as $T_i\text{End} < T_j\text{End}$.

Therefore, from (A) and (B), all edges in the MVSG are ordered with respect to the transaction end timestamp order, and thus (from Property 1) they cannot be involved in a cycle. Hence, the MV histories that are accepted by our multi-version optimistic concurrency control scheduler are 1SR.

Bibliography

[BHG87] Philip A. Bernstein, Vassos Hadzilacos and Nathan Goodman: "Concurrency Control and Recovery in Database Systems", Addison-Wesley, 1987, ISBN 0-201-10715-5.

[WV02] Gerhard Weikum, Gottfried Vossen: "Transactional Information Systems: Theory, Algorithms, and the Practice of Concurrency Control and Recovery", Morgan Kaufmann, 2002, ISBN 1-55860-508-8.