# No False Negatives: Accepting All Useful Schedules in a Fast Serializable Many-Core System

Dominik Durner, Thomas Neumann

April 10, 2019

Technische Universität München

▶ Concurrency control schemes only approximate the class of serializable schedules, such as 2PL, OCC, TicToc

▶ Therefore, unexpected behavior and also unnecessary aborts are introduced

▶ Spurious aborts due to implementation artifacts that are hard to understand

- ▶ Concurrency control schemes only approximate the class of serializable schedules, such as 2PL, OCC, TicToc
- ▶ Therefore, unexpected behavior and also unnecessary aborts are introduced
- ▶ Spurious aborts due to implementation artifacts that are hard to understand
- ▶ For example, 2PL cannot accept:

▶ Concurrency control schemes only approximate the class of serializable schedules, such as 2PL, OCC, TicToc

▶ Therefore, unexpected behavior and also unnecessary aborts are introduced

▶ Spurious aborts due to implementation artifacts that are hard to understand

▶ For example, 2PL cannot accept:

- Concurrency control schemes only approximate the class of serializable schedules, such as 2PL, OCC, TicToc
- Therefore, unexpected behavior and also unnecessary aborts are introduced
- Spurious aborts due to implementation artifacts that are hard to understand
- For example, 2PL cannot accept:



- Only Serialization Graph Testing (SGT) accepts all valid schedules
- SGT seems to be too expensive and not scalable

▶ Conflict graphs allow to accept all conflict serializable schedules

| all schedules |
|---|
| CSR |

▶ Conflict graphs allow to accept all conflict serializable schedules

▶ Recoverability is independent of serializability

| all schedules | RC | |
|---|---|---|
| CSR | | |

- Conflict graphs allow to accept all conflict serializable schedules
- Recoverability is independent of serializability
- DBMS users expect to see committed changes

| all schedules | RC | |
|---|---|---|
| CSR | | |
| OCSR | | |
| COCSR | | |

- Conflict graphs allow to accept all conflict serializable schedules
- Recoverability is independent of serializability
- DBMS users expect to see committed changes

| all schedules | RC | |
|---|---|---|
| CSR | | |
| OCSR | | |
| COCSR | | |

Note that $S2PL \subsetneq COCSR \cap RC$

Our approach leverages the conflict graph and

1. accepts all useful *COCSR ∩ RC* schedules
2. meets users' expectations
3. has low overhead for maintaining the graph
4. scales to many-core systems

- ▶ Theorem: $s \in CSR \Leftrightarrow CG(s)$ is acylic
- ▶ Update $CG(s)$ at operation arrival and allow if $CG(s)$ is acyclic
- ▶ Remove all outgoing edges of a node at its deletion

▶ Theorem: $s \in CSR \Leftrightarrow CG(s)$ is acylic
▶ Update $CG(s)$ at operation arrival and allow if $CG(s)$ is acyclic
▶ Remove all outgoing edges of a node at its deletion

Example: $s = r_0[x] \, w_0[x]$

► Theorem: $s \in CSR \Leftrightarrow CG(s)$ is acylic
► Update $CG(s)$ at operation arrival and allow if $CG(s)$ is acyclic
► Remove all outgoing edges of a node at its deletion

Example: $s = r_0[x] \, w_0[x] \, r_1[x]$

- ▶ Theorem: $s \in CSR \Leftrightarrow CG(s)$ is acylic
- ▶ Update $CG(s)$ at operation arrival and allow if $CG(s)$ is acyclic
- ▶ Remove all outgoing edges of a node at its deletion

Example: $s = r_0[x]\, w_0[x]\, r_1[x]\, r_2[x]\, w_2[x]$

- ▶ Theorem: $s \in CSR \Leftrightarrow CG(s)$ is acylic
- ▶ Update $CG(s)$ at operation arrival and allow if $CG(s)$ is acyclic
- ▶ Remove all outgoing edges of a node at its deletion

Example: $s = r_0[x] \, w_0[x] \, r_1[x] \, r_2[x] \, w_2[x] \, w_2[y] \, c_2 \, c_0 \, c_1$

$$\Rightarrow s \in CSR$$

► SGT has the best theoretical properties of accepting all valid schedules
► However, previous work fails to implement SGT efficiently in practice

▶ SGT has the best theoretical properties of accepting all valid schedules

▶ However, previous work fails to implement SGT efficiently in practice

> We developed the first practical and scalable algorithm that leverages
> the theoretical superior concept of graph-based serialization testing

Pitfall: Deletion of a committed node $t_c$

Pitfall: Deletion of a committed node $t_c$

Example: $s = r_0[x]\, w_0[x]\, r_1[x]\, r_2[x]\, w_2[x]\, w_2[y]\, c_2$

Pitfall: Deletion of a committed node $t_c$

Example: $s = r_0[x]\, w_0[x]\, r_1[x]\, r_2[x]\, w_2[x]\, w_2[y]\, c_2$

Pitfall: Deletion of a committed node $t_c$

Example: $s = r_0[x]\, w_0[x]\, r_1[x]\, r_2[x]\, w_2[x]\, w_2[y]\, c_2\, \mathbf{r_0[y]}\, c_0\, c_1$

Pitfall: Deletion of a committed node $t_c$

Example: $s = r_0[x]\, w_0[x]\, r_1[x]\, r_2[x]\, w_2[x]\, w_2[y]\, c_2\, \mathbf{r_0[y]}\, c_0\, c_1$



$\Rightarrow s \notin CSR$, but not detectable if $t_2$ was deleted

Pitfall: Deletion of a committed node $t_c$

Example: $s = r_0[x]\, w_0[x]\, r_1[x]\, r_2[x]\, w_2[x]\, w_2[y]\, c_2\, \mathbf{r_0[y]}\, c_0\, c_1$



$\Rightarrow s \notin CSR$, but not detectable if $t_2$ was deleted

Deletion of committed node is only allowed if all incoming edges are removed

Every transaction commit needs to wait until it is not dependent on in-flight results

Example: $s = r_0[x]\, w_0[x]\, r_1[x]\, r_2[x]\, w_2[x]\, w_2[y]\, c_2\, c_0\, c_1$

Every transaction commit needs to wait until it is not dependent on in-flight results

Example: $s = r_0[x]\, w_0[x]\, r_1[x]\, r_2[x]\, w_2[x]\, w_2[y]\, c_2\, \cancel{c_0}\, \cancel{c_1}\, \mathbf{a_0}\, \mathbf{a_1}$

Every transaction commit needs to wait until it is not dependent on in-flight results

Example: $s = r_0[x]\, w_0[x]\, r_1[x]\, r_2[x]\, w_2[x]\, w_2[y]\, c_2\, \cancel{c_0\, c_1}\, \mathbf{a_0}\, \mathbf{a_1}$



No incoming write-read, write-write edge from an uncommitted node allowed

No (uncommitted) incoming edge at commit time to preserve the commit order

No (uncommitted) incoming edge at commit time to preserve the commit order

Example: $s = r_0[x]\ w_1[x]\ c_1$

No (uncommitted) incoming edge at commit time to preserve the commit order

Example: $s = r_0[x]\ w_1[x]\ \bowtie\!\!\!\!/\ d_1$

No (uncommitted) incoming edge at commit time to preserve the commit order

Example: $s = r_0[x]\ w_1[x]\ \bowtie\!\!\!\!/\ d_1\ r_2[y]\ c_2$

No (uncommitted) incoming edge at commit time to preserve the commit order

Example: $s = r_0[x] \ w_1[x] \ \bowtie \ d_1 \ r_2[y] \ c_2 \ w_0[y]$

No (uncommitted) incoming edge at commit time to preserve the commit order

Example: $s = r_0[x]\ w_1[x]\ \bowtie\!\!\!\!/\ d_1\ r_2[y]\ c_2\ w_0[y]\ c_0$

No (uncommitted) incoming edge at commit time to preserve the commit order

Example: $s = r_0[x]\ w_1[x]\ \bowtie\!\!\!\!/\ d_1\ r_2[y]\ c_2\ w_0[y]\ c_0\ c_1$

No (uncommitted) incoming edge at commit time to preserve the commit order

Example: $s = r_0[x] \ w_1[x] \ \bcancel{\bowtie} \ d_1 \ r_2[y] \ c_2 \ w_0[y] \ c_0 \ c_1$

$s_{orig} = r_0[x] \ w_1[x] \ c_1 \ r_2[y] \ c_2 \ w_0[y] \ c_0$

with $s' = t_2 \ t_0 \ t_1$, but $s_{orig} \notin COCSR$

No (uncommitted) incoming edge at commit time to preserve the commit order

Example: $s = r_0[x]\ w_1[x]\ \cancel{w_1}\ d_1\ r_2[y]\ c_2\ w_0[y]\ c_0\ c_1$

$s_{orig} = r_0[x]\ w_1[x]\ c_1\ r_2[y]\ c_2\ w_0[y]\ c_0$

with $s' = t_2\ t_0\ t_1$, but $s_{orig} \notin COCSR$

All useful $COCSR \cap RC$ schedules accepted due to commit delays

Committed nodes are deleted directly including all outgoing edges

- ▶ No incoming edges to commit simplifies cycle check
- ▶ Conflict graph is accessed concurrently by multiple threads
- ▶ No other transaction is allowed to modify a node during its final check

- ▶ No incoming edges to commit simplifies cycle check
- ▶ Conflict graph is accessed concurrently by multiple threads
- ▶ No other transaction is allowed to modify a node during its final check

Transaction local shared/exclusive locks help to scale the graph

Setup:
- ▶ 4-socket Intel Xeon server (60 cores) with 1TB DRAM
- ▶ Every transaction is scheduled on one worker thread
- ▶ Aborts require undos and restarts of the aborted transactions

Algorithms:
- ▶ Our SGT-based approach
- ▶ TicToc
- ▶ 2PL with row based atomic read-write locks and deadlock prevention

# SmallBank Medium Contention (1000 Customers)

Our SGT has competitive throughput while reducing aborts significantly!

accepts all useful
*COCSR ∩ RC* schedules

accepts all useful
*COCSR ∩ RC* schedules



reduces aborted schedules
and meets users' expectations

accepts all useful
*COCSR ∩ RC* schedules



has low protocol overhead
and scales to many-core systems



reduces aborted schedules
and meets users' expectations