

# IN2267 – Transaction Systems

## Week 4: Snapshot Isolation Concurrency Control

*Dr. Uwe Röhm  
The University of Sydney*

## Learning Objectives

### ■ Background

- ▶ Reprise: Concurrency Control Approaches
- ▶ Synchronization Problems & ANSI SQL Isolation Levels

### ■ Optimistic Concurrency Control

- ▶ Snapshot Isolation
- ▶ SI Implementation Details
- ▶ Serializable Snapshot Isolation

### ■ Outlook

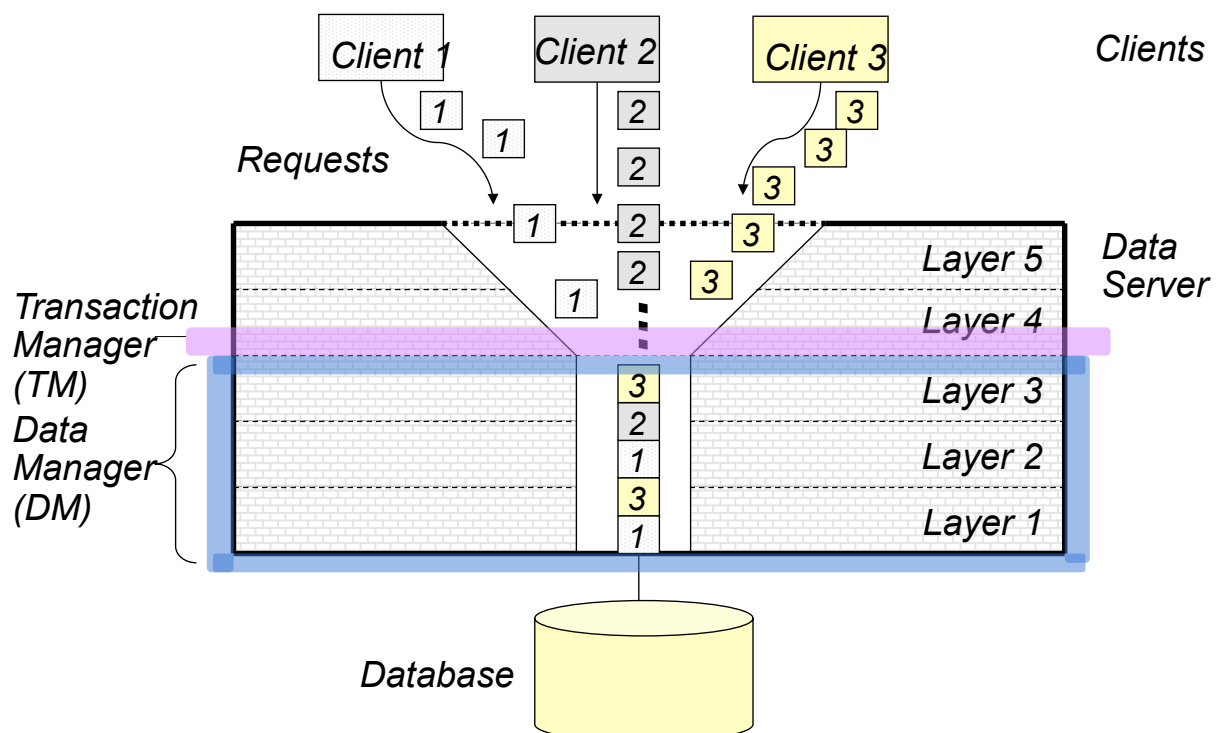
- ▶ SI on multi-core CPUs

*Based on USydney slides from U. Roehm and M. Cahill,  
and Weikum/Vossen (2002) "Transactional Information Systems"*

# Concurrency Control

- The **concurrency control** of a DBMS is responsible for enforcing serializability among concurrent transactions
  - ▶ Two important techniques: *Locking* and *Versioning*
- Note: In addition to serializable, DBMSs implement less stringent isolation levels
  - ▶ Serializable schedules correct for *all* applications
  - ▶ Less stringent levels *do not* guarantee correctness for all applications, but are correct for *some*
  - ▶ Application programmer is responsible for choosing appropriate level

## Transaction Scheduler



# Potential Anomalies

- **lost update** ('dirty write'):  
updating a value that was already updated by a concurrent, uncommitted transaction.
- **dirty read**:  
reading a value that was updated by a concurrent, uncommitted transaction
- **non-repeatable read** ('fuzzy read'):  
reading a value twice gives different results because of a concurrent update by a different transaction in between
- **phantom read**:  
using the same selection criteria on a table twice gives different result sets, because a concurrent updater deleted or inserted elements satisfying the selection criteria

# ANSI SQL Isolation Levels

- Defined in terms of anomalies
  - ▶ Anomaly prohibited at one level is also prohibited at all higher levels
  - ▶ READ UNCOMMITTED: all anomalies possible
  - ▶ READ COMMITTED: dirty read prohibited
  - ▶ REPEATABLE READ: reads of individual tuples are repeatable (but phantoms are possible)
  - ▶ SERIALIZABLE: phantoms prohibited; transaction execution is serializable
- Serializable is according to SQL standard the default...
  - ▶ In practice, most systems have weaker default level! (Oracle!)
- Lower degrees of consistency useful for gathering approximate information about the database, e.g., statistics for query optimizer.

# Comparison of SQL Isolation Levels

Isolation Level	Lost Update	Dirty Read	Unrepeatable Read	Phantom
<b>READ UNCOMMITTED</b>	not possible	<i>possible</i>	<i>possible</i>	<i>possible</i>
<b>READ COMMITTED</b>	not possible	not possible	<i>possible</i>	<i>possible</i>
<b>REPEATABLE READ</b>	not possible	not possible	not possible	<i>possible</i>
<b>SERIALIZABLE</b>	not possible	not possible	not possible	not possible

**Note:** ANSI SQL Isolation Level *SERIALIZABLE*  
!=  
Definition in serialisability theory  
(such as conflict serialisability)

## Locks in Relational Databases

- DBMS guarantees that each SQL statement is isolated
- Early (non-strict) lock release used to implement levels
  - ▶ Short-term locks - held for duration of single statement
  - ▶ Long-term locks - held until transaction completes (strict)
- At all levels, transactions obtain long-term write locks
- This means for isolation levels:
  - ▶ READ UNCOMMITTED - no read locks (dirty reads possible since transaction can read a write-locked item)
  - ▶ READ COMMITTED - short-term read locks on rows (non-repeatable reads possible since transaction releases read lock after reading)
  - ▶ REPEATABLE READ - long-term read locks on rows (phantoms possible)
  - ▶ SERIALIZABLE - combination of table, row, and index locks

# Agenda

## ■ Background

- ▶ Reprise: Concurrency Control Approaches
- ▶ Synchronization Problems & ANSI SQL Isolation Levels

## ■ Optimistic Concurrency Control

- ▶ Snapshot Isolation
- ▶ SI Implementation Details
- ▶ Serializable Snapshot Isolation

## ■ Outlook

- ▶ SI on multi-core CPUs

# Optimistic Concurrency Control

- Locking is a conservative approach in which conflicts are prevented. Disadvantages:
  - ▶ Lock management overhead.
  - ▶ Deadlock detection/resolution.
  - ▶ Lock contention for heavily used objects.
- If conflicts are rare, we might be able to gain concurrency by not locking, and instead checking for conflicts before transactions commit.
  - ▶ Optimistic, validating CC
  - ▶ Multiversion CC

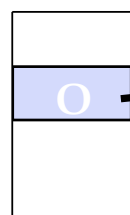
# Snapshot Isolation – Conceptual Idea

- Every transaction reads from a consistent **snapshot** (copy) of the database (the db state of when tx started)
- Writes are collected into a transaction's **writeset**
  - ▶ Writeset is not visible to concurrent transactions
- At commit time, the writeset is compared to the writesets of all concurrent transactions.
  - ▶ If they are disjoint (no overlap), then they are applied to the actual database => **commit**
  - ▶ If there's an overlap with the writeset of a concurrent, but already committed transaction, the later transaction must **abort**
    - => **"First Committer Wins" rule**

## In Practice: Snapshot Isolation (SI)

- A **multiversion concurrency control** mechanism which was described in SIGMOD '95 by H. Berenson, P. Bernstein, J. Gray, J. Melton, E. O'Neil, P. O'Neil
  - ▶ Incremental implementation of an optimistic concurrency control scheme
- **Core Idea:** Let writers create a "new" copy while readers use an appropriate "old" copy.

Current versions of DB objects



**VERSION POOL**  
(Older versions that may be useful for some active readers.)

- ❖ **Readers are always allowed to proceed.**
  - But may be blocked until writer commits.

# Reads with Snapshot Isolation

- **Multiversion database:** The old value of an item is *not* overwritten when it is updated (no ‘in-place updates’). Instead, a new version is created
- Read of an item does not necessarily give latest value
- Instead, use old versions (kept with timestamps) to find value that had been most recently committed *at the time the transaction started*
  - ▶ Exception: if the txn has modified the item, use the value it wrote itself
- The transaction sees a “snapshot” of the database, at an earlier time
  - ▶ Intuition: this should be consistent, if database was consistent before
  - ▶ No read locks necessary: a transaction reads all values from latest snapshot at time it started. Thus, **read/only transactions do not wait.**

# Writes with Snapshot Isolation

- A transaction  $T$  that has updated  $x$  can commit if no other transaction that concurrently updated  $x$  has committed
  - ▶ **“First-committer-wins” rule:**
    - ▶ Updater  $T$  will not be allowed to commit if any other transaction has committed and *installed a changed* value for that item, between  $T$ 's start (snapshot) and  $T$ 's commit
    - ▶ Similar to optimistic validation-based cc, but only write-sets are checked
- $T$  must hold X-lock on modified items at time of commit, to install them. In practice, commit-duration X-locks may be set when write executes. These help to allow conflicting modifications to be detected (and  $T$  aborted) when  $T$  tries to write the item, instead of waiting till  $T$  tries to commit.

# Benefits of SI

- Reading is *never* blocked, and also doesn't block other transactions' activities
  - ▶ Fast performance similar to Read Committed
- Avoids the usual anomalies
  - ▶ No dirty read
  - ▶ No lost update
  - ▶ No inconsistent read
  - ▶ Set-based selects are repeatable (no phantoms)
- ▶ Note: not Write-Skews – cf. later slides

# Who does this?

- Oracle: used for "Isolation Level Serializable"
  - ▶ But does *not* guarantee serializable execution as defined in standard transaction management theory!
- PostgreSQL: used for "Isolation Level Serializable"
  - ▶ As of version 9.1 guarantees serializable execution, but not earlier
- Available in Microsoft SQL Server 2005 and above as "Isolation Level Snapshot"
  - ▶ If mssql db is configured to provide snapshots
- Berkeley DB
- MySQL / InnoDB (sort of)



# Agenda

- Background
  - ▶ Reprise: Concurrency Control Approaches
  - ▶ Synchronization Problems & ANSI SQL Isolation Levels
  
- **Optimistic Concurrency Control**
  - ▶ Snapshot Isolation
  - ▶ **SI Implementation Details**
  - ▶ Serializable Snapshot Isolation
  
- Outlook
  - ▶ SI on multi-core CPUs

## SI Design Choices

- Tuple Versions
  - ▶ Store old versions or generate as required?
  
- Granularity
  - ▶ should individual records be versioned, or pages?
  - ▶ (or even tables?)
  
- How is a snapshot represented?  
("what is time?")

# SI Common Themes

- Almost every implementation takes locks for updates
  - ▶ This blocks other updates until commit / abort
  - ▶ Guarantees forward progress
  - ▶ Reduces conflict-abort-retry thrashing
- First-committer-wins implemented as “has a version been committed since I started?”

# PostgreSQL: Intro

- Full RDBMS, long history
- Provides SI when you ask for REPEATABLE READ or SERIALIZABLE
- Stores old versions of rows in the database
  - ▶ Needs regular VACUUMing

# pgsql: SnapshotData

```
33 typedef struct SnapshotData
34 {
35     SnapshotSatisfiesFunc satisfies; /* tuple test function */
36
37     /*
38      * The remaining fields are used only for MVCC snapshots, and are normally
39      * just zeroes in special snapshots. (But xmin and xmax are used
40      * specially by HeapTupleSatisfiesDirty.)
41      *
42      * An MVCC snapshot can never see the effects of XIDs >= xmax. It can see
43      * the effects of all older XIDs except those listed in the snapshot. Xmin
44      * is stored as an optimization to avoid needing to search the XID arrays
45      * for most tuples.
46      */
47     TransactionId xmin; /* all XID < xmin are visible to me */
48     TransactionId xmax; /* all XID >= xmax are invisible to me */
49     TransactionId *xip; /* array of xact IDs in progress */
50     uint32 xcnt; /* # of xact ids in xip[] */
51     /* note: all ids in xip[] satisfy xmin <= xip[i] < xmax */
52     int32 subxcnt; /* # of xact ids in subxip[], -1 if overflow */
53     TransactionId *subxip; /* array of subxact IDs in progress */
54
55     ...
56 } SnapshotData; [src/include/utils/snapshot.h]
```

IN2267 "Transaction Systems" - WS 2013/14 (Guest Lecture U. Röhm)

21

# Pgsql: Tuple Visibility

[src/backend/utils/time/tqual.c]

```
327 *      mao says 17 march 1993: the tests in this routine are correct;
328 *      if you think they're not, you're wrong, and you should think
329 *      about it again. i know, it happened to me. we don't need to
330 *      check commit time against the start time of this transaction
331 *      because 2ph locking protects us from doing the wrong thing.
332 *      if you mess around here, you'll break serializability. the only
333 *      problem with this code is that it does the wrong thing for system
334 *      catalog updates, because the catalogs aren't subject to 2ph, so
335 *      the serializability guarantees we provide don't extend to xacts
336 *      that do catalog accesses. this is unfortunate, but not critical
337 */
338 bool
339 HeapTupleSatisfiesNow(HeapTupleHeader tuple, Snapshot snapshot, Buffer buffer)
340 { ... }
```

## ■ Tuple header defines a closed-open transaction-time interval

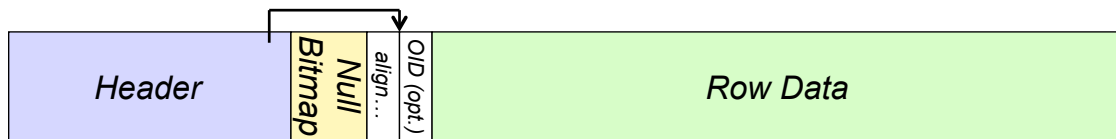
- ▶ Basic Idea: A Tuple is visible iff
  - ▶ xmin is a committed transaction ID < own transaction ID and not in-progress at transaction start.
  - ▶ xmax is either blank, or greater than the start transaction ID and in-progress at transaction start

IN2267 "Transaction Systems" - WS 2013/14 (Guest Lecture U. Röhm)

22

# Row Format in PostgreSQL

Row = RowHeader + NullBitmap + *alignment padding* + (OID) + RowData



## ■ RowHeader:

- ▶ 23 Bytes; cf. next slide

## ■ NullBitmap

- ▶ variable len: either 0 (if all NOT NULL) or  $((|Columns| + 7) / 8)$  Bytes
- ▶ one bit per attribute; 1 if NULL, 0 otherwise

## ■ OID (*PostgreSQL speciality as it supports objects*)

- ▶ *fix 4 Bytes (optional, depending whether table WITH OIDs or not)*

## ■ RowData = FixedColumns + VarColumns

- FixedColumns: directly stored & aligned!
- VarColumns = varattrib + userdata + aligned
  - varattrib = 4 bytes length words including 2 bits for compression/TOAST flags

# Row Format in PostgreSQL (cont' d)

## ■ Row Header structure

- ▶ 23 bytes (plus bitmap plus padding; cf. *t\_hoff* value as 'pointer')
- ▶ Cf. **src/include/access/htup.h**:

```
typedef struct HeapTupleHeaderData
```
- ▶ Some information on visibility of a tuple for current transaction snapshot or newer version (needed for snapshot isolation algorithm)
  - **t\_xmin** TransactionId 4 bytes insert XID stamp
  - **t\_xmax** TransactionId 4 bytes delete XID stamp
  - **t\_cid** CommandId 4 bytes insert CID stamp (*actual a UNION struct*)
  - **t\_ctid** ItemPointerData 6 bytes current TID of this or newer row version
- ▶ How long is this row? Is it variable length? Does it have NULLS?
  - **t\_natts** int16 2 bytes number of attributes
  - **t\_infomask** uint16 2 bytes various flag bits
    - e.g. HAS\_NULL | HASVARWIDTH | HASOID | locks(!)
  - **t\_hoff** uint8 1 byte /\* sizeof header incl. bitmap, padding \*/

# pgSQL SI: Tuple Visibility Example

Current Transaction ID: 100

Cre 30 Exp	Visible	In-Progress Transactions: 25
Cre 30 Exp 80	Skip	50
Cre 30 Exp 110	Visible	75
Cre 30 Exp 75	Visible	
Cre 50 Exp	Skip	
Cre 110 Exp	Skip	

## InnoDB: Intro

- Transactional backend for MySQL
  - ▶ MySQL supports different storage engines!
- Only needs to deal with read / writes of rows
  - ▶ MySQL looks after SQL and query processing
- Generates old values on demand
  - ▶ Uses "undo" records from the log

# InnoDB: Concurrency Control

## ■ MVCC but not SI

- ▶ Read-only Transactions (pure queries) read from a snapshot
- ▶ Locking reads (including updates) read most recently committed value
- ▶ No first-committer-wins rule

## InnoDB: read\_view

```
116 /* Read view lists the trx ids of those transactions for which a consistent
117 read should not see the modifications to the database. */
118
119 struct read_view_t{
...
124     trx_id_t low_limit_no; /* The view does not need to see the undo
125                            logs for transactions whose transaction number
126                            is strictly smaller (<) than this value: they
127                            can be removed in purge if not needed by other
128                            views */
129     trx_id_t low_limit_id; /* The read should not see any transaction
130                            with trx id >= this value; in other words, this is the "high water mark" */
131     trx_id_t up_limit_id; /* The read should see all trx ids which
132                            are strictly smaller (<) than this value; this is the "low water mark" */
133     ulint n_trx_ids; /* Number of cells in the trx_ids array */
134     trx_id_t* trx_ids; /* Additional trx ids which the read should
135                            not see: typically, these are the active
136                            transactions at the time when the read is
137                            serialized, except the reading transaction
138                            itself; the trx ids in this array are in a
139                            descending order */
140     trx_id_t creator_trx_id; /* trx id of creating transaction, or
141                            (0, 0) used in purge */
142     UT_LIST_NODE_T(read_view_t) view_list;
143     /* List of read views in trx_sys */
144 };
```

[ storage/innobase/include/read0read.h ]

# InnoDB: read old versions

```
438 /* Constructs the version of a clustered index record which a consistent
439 read should see. We assume that the trx id stored in rec is such that
440 the consistent read should not see rec in its present version. */
442 ulint
443 row_vers_build_for_consistent_read(
444 /*=====*/
445     /* out: DB_SUCCESS or DB_MISSING_HISTORY */
446     const rec_t* rec, /* in: record in a clustered index; */
447     read_view_t* view, /* in: the consistent read view */
448     rec_t** old_vers) /* out, own: old version, or NULL if the
449                        record does not exist in the view, that is,
450                        it was freshly inserted afterwards */
451 {
452     ...
453     for (;;) {
454         ...
455         err = trx_undo_prev_version_build(rec, mtr, version, index, *offsets, heap, &prev_version);
456         ...
457         if (prev_version == NULL) {
458             /* It was a freshly inserted version */
459             *old_vers = NULL;
460             err = DB_SUCCESS;
461             break;
462         }
463         ...
464         trx_id = row_get_rec_trx_id(prev_version, index, *offsets);
465         if (read_view_sees_trx_id(view, trx_id)) {
466             ...
467             err = DB_SUCCESS;
468             break;
469         }
470         version = prev_version;
471     } /* for (;;) */
472 }
```

IN2267 "Transaction Systems" - WS 2013/14 (Guest Lecture U. Röhm)

29

## pgsql and InnoDB: Lessons

- MVCC is used by default
  - ▶ Has to work well under most conditions
  
- Old versions need space somewhere
  - ▶ postgresql: in the database file
  - ▶ InnoDB: undo records in the log buffer

# Summary: SI Design Space

	<b>BDB</b>	<b>pgsql</b>	<b>InnoDB</b>
<b>old versions</b>	store in cache	store on disk	generate on demand
<b>granularity</b>	page	record	record
<b>transaction time</b>	LSNs	snapshot of active txnIDs	snapshot of active txnIDs

## Agenda

### ■ Background

- ▶ Reprise: Concurrency Control Approaches
- ▶ Synchronization Problems & ANSI SQL Isolation Levels

### ■ **Optimistic Concurrency Control**

- ▶ Snapshot Isolation
- ▶ SI Implementation Details
- ▶ **Serialisable Snapshot Isolation**

### ■ Outlook

- ▶ SI on multi-core CPUs

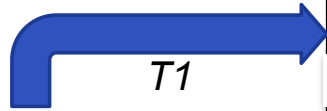


# Problem with SI: Snapshot Isolation $\neq$ Serializable

Inherent constraint: for every date, there is at least 1 doctor on duty

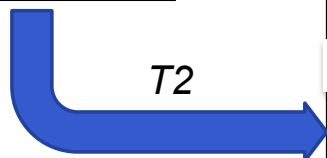
**Invariant violated!**

Doctor	Shift	Status
House	12 June	on duty
Grey	12 June	on duty



Doctor	Shift	Status
House	12 June	<b>reserve</b>

**Oops...**



Doctor	Shift	Status
Grey	12 June	<b>reserve</b>

## Write Skew

- Formally this is known as **Write Skew Problem**
  - ▶ SI breaks serializability when transactions modify *different* items, each based on a previous state of the item the other modified
  - ▶ This is fairly rare in practice
    - Eg the TPC-C benchmark runs correctly under SI: when transactions conflict due to modifying different data, there is also a shared item they both modify too (like a total quantity) so SI will abort one of them

# Vendor Advice

- Oracle: *“Database inconsistencies can result unless such application-level consistency checks are coded with this in mind, even when using serializable transactions.”*
- *“PostgreSQL's Serializable mode does not guarantee serializable execution...”*
  - ▶ **FIXED since PostgreSQL 9.1!!!**
- SQL Server: *only gives performance advices, but keeps quiet on the correctness issue...*

# Serializable SI

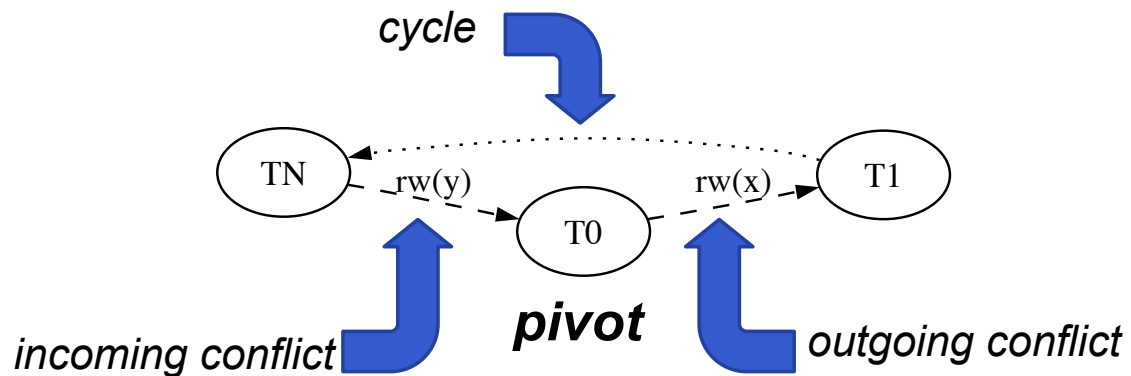
- Theory exists about how write-skews can be detected
  - ▶ A. Fekete, D. Liarakapis, E. O'Neil, P. O'Neil, D. Shasha in TODS2005: *“Making Snapshot Isolation Serializable”*
    - Analyze the graph of transaction conflicts
    - Conditions on the graph for application to be serializable at SI; def. of *dangerous structure*
- Solution: Two Approaches
  - ▶ Introduce artificial ww-conflicts to application to trigger first-committer-wins rule
    - Requires semantic program analysis before -> NP Hard
  - ▶ Modify SI CC to identify ‘dangerous patterns’ in concurrent snapshot transactions and abort one of them
    - false positives are possible
    - PhD thesis of Michael Cahill at U Sydney [SIGMOD2008]
    - Fortunately, we now have a system that implements this → PostgreSQL-9.1

# Dangerous Structures => SI Anomalies

[Fekete et al., TODS2005]

## SI Serialisability Test:

Build the static dependency graph,  
then check for any “**dangerous structures**”



## Limitations of This Approach

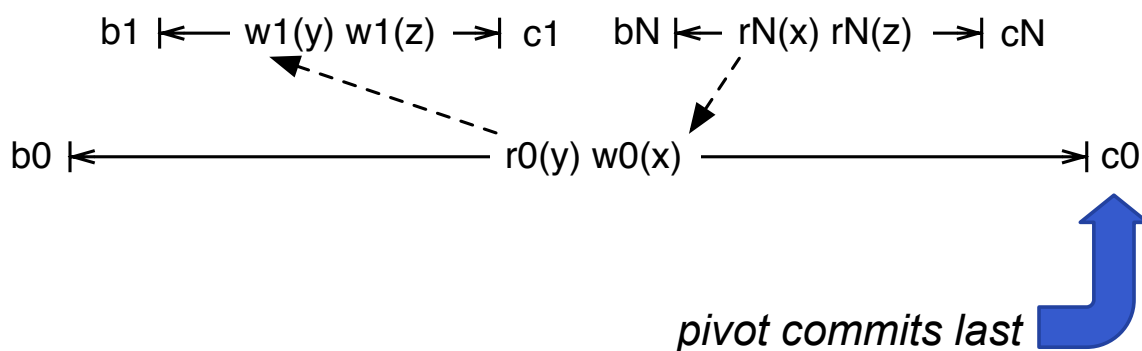
- Determining the conflict graph is non-trivial
- Repeat for every change to the application
- Ad hoc queries not supported
- Difficult to automate: reasoning required to avoid false positives
- What to do with the outcome?
  - ▶ Standard approach as of 2005 was that the applications needed to get changed, e.g. by introducing artificial writes to 'promote' a rw-dependency to a ww-dependency

# “Online” SSI Approach

[Cahill, SIGMOD2008]

- New algorithm for serializable isolation
  - ▶ Online, dynamic
  - ▶ Modifications to standard Snapshot Isolation
- Core Idea:
  - ▶ Detect read-write conflicts at runtime
  - ▶ Abort transactions with consecutive rw-edges
  - ▶ Don't do full cycle detection

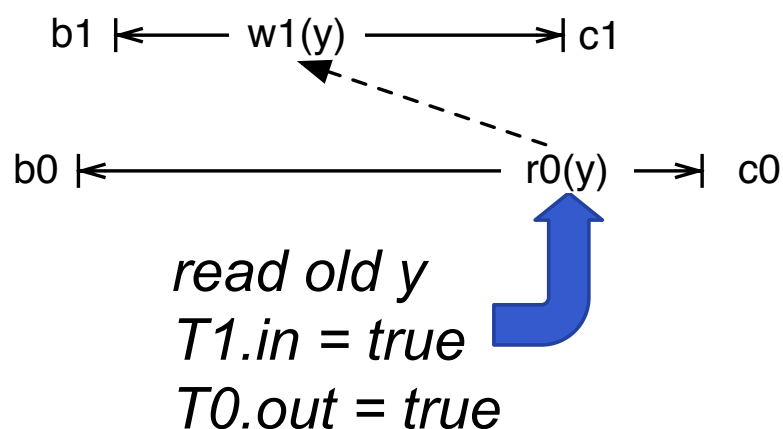
## SI Anomalies: a Simple Case



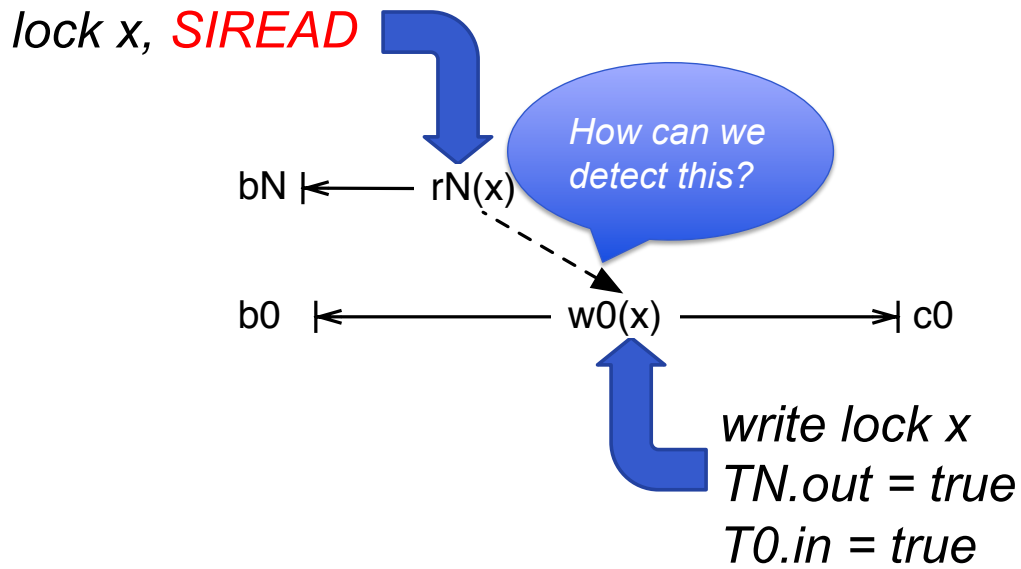
# The Algorithm in a Nutshell

- Add two flags to each transaction (*in* & *out*)
- Set  $T_0.out$  if *rw*-conflict  $T_0 \rightarrow T_1$
- Set  $T_0.in$  if *rw*-conflict  $T_N \rightarrow T_0$
- Abort  $T_0$  (the pivot) if both  $T_0.in$  and  $T_0.out$  are set
  - ▶ If  $T_0$  has already committed, abort the conflicting transaction

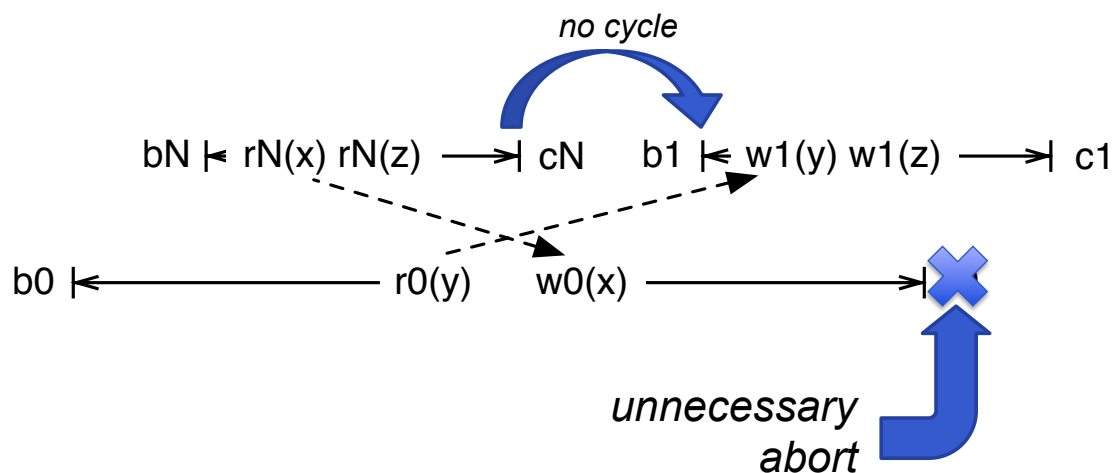
## Detection: Write before Read



# Detection: Read before Write



# Main Disadvantage: False Positives



# SSI Variants

- SSI:  
Original SIGMOD08 paper (and more detailed in TODS2009)
- Precise Serialisable Snapshot Isolation (PSSI)
  - ▶ Revilak *et al* in ICDE 2011
  - ▶ In essence a full serialization graph test on top of Cahill's SSI
- Revilak's ICDE2011 paper also did their own implementation of TODS2009 algorithm in InnoDB -> referred to as ESSI
- PostgreSQL implementation of SSI by D. Ports (VLDB2012)
  - ▶ Including some optimisations for read-only transactions which do not need to take SIREAD locks on a *safe snapshot*

# Design Decisions for SSI

- Local versus Global Dependency Tracking
  - ▶ anti-dependencies tracked per transaction or in a separate global data structure?
- Approximate versus accurate serialisability check
  - ▶ Check only for a *dangerous structure*, or perform a full cycle test?
- Ongoing checks versus commit-time check
  - ▶ check for potential abort with each update operation or at commit?

# Previous Work on SSI

	<b>SSI</b> [Cahill, SIGMOD08]	<b>ESSI</b> [Revilak, ICDE11; Cahill, TODS09]	<b>PSSI</b> [Revilak, ICDE11]	<b>pgSSI</b> [Ports, VLDB12]
<b>Tracking</b>	local	local	global	local
<b>Data Structure</b>	two <b>Bits</b> per transact.	two <b>Pointers</b> per transact.	cycle testing <b>Graph (CTG)</b>	two <b>Lists</b> per transact.
<b>Check</b>	<i>dangerous structure</i>	<i>dangerous structure</i> (using CTG)	<i>cycle test</i>	<i>dangerous structure</i>
<b>When</b>	each update	at commit	at commit	each update

## SSI in PostgreSQL 9.1

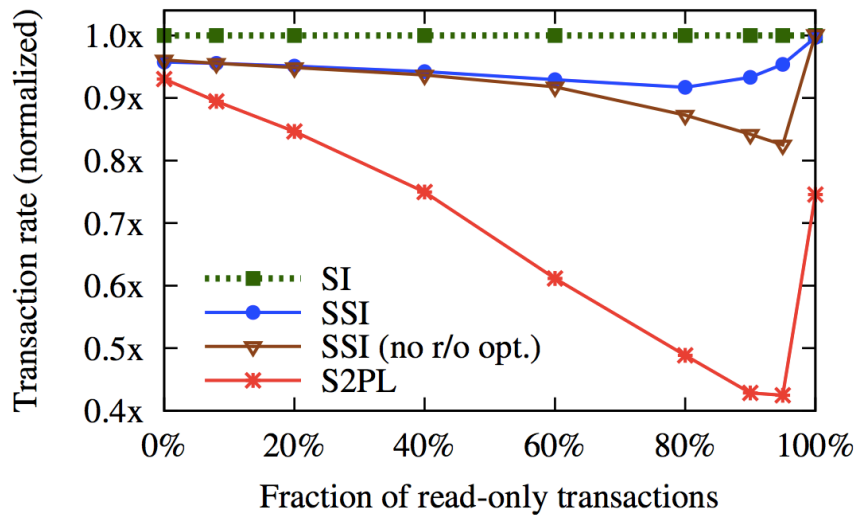
[Ports, VLDB2012]

- Basically follows Cahill's approach
- Some additional optimisations:
  - ▶ Identify 'safe snapshot' situations for read-only transactions
  - ▶ If a read-only tx is found to run on a safe snapshot, will never abort and does not need to take and SIREAD locks (less overhead)
  - ▶ For long running transactions (e.g. backup):
    - Allow DBA to delay them until they are guaranteed to run on a safe snapshot
- Fully integrated into DBMS
  - ▶ Needed a bit more complex code than Michael's prototype ;)
- ISOLATION LEVEL SERIALISABLE is now SSI
- ISOLATION LEVEL REPEATABLE READ is the former SI
- Default isolation level is still read committed...



# Performance Penalty for Correctness?

- Obviously, tracking dependencies and aborting transactions (some of them false aborts) doesn't come for free
- What are the costs for being correct?



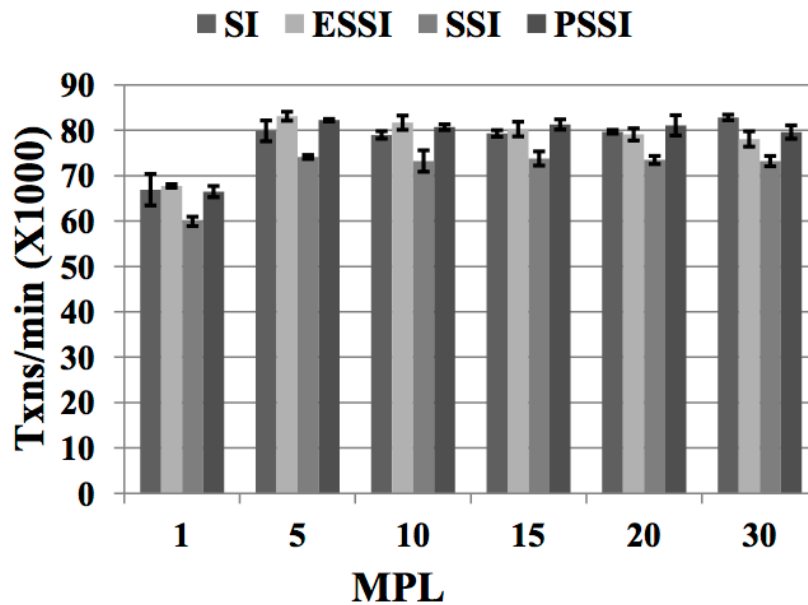
psql 9.1 with TPC-C, in-memory configuration (25 warehouses) [VLDB2012]

## Research Question

Is SSI still as fast on a multicore server?

# The Price of Serialisability

- MySQL 5.1.3 with InnoDB on 24 core Xeon Server
  - ▶ Implemented SSI, ESSI and PSSI in same engine (Linux)
  - ▶ quantified overhead at ca. 10%



IN2267 "Transaction Systems" - WS 2013/14 (Guest Lecture U. Röhm)

51

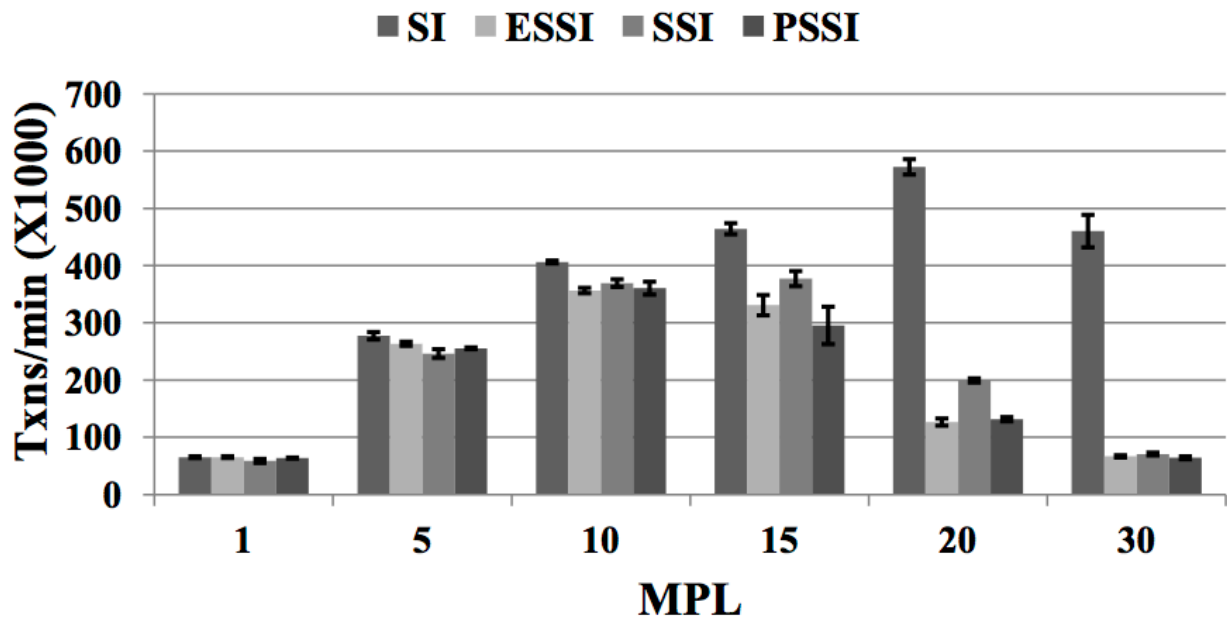
## Scalability with Number of Cores?

- The previous figure was for just a single core server
  - ▶ with a (75%-RO-25%-RU) workload
- What happens if we enable all 24 cores of the server?

IN2267 "Transaction Systems" - WS 2013/14 (Guest Lecture U. Röhm)

52

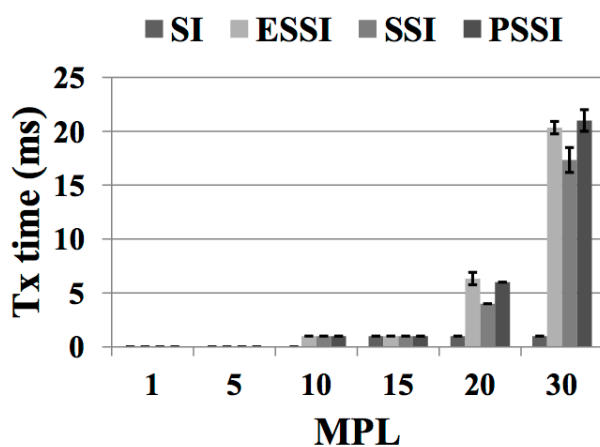
# The Big Picture



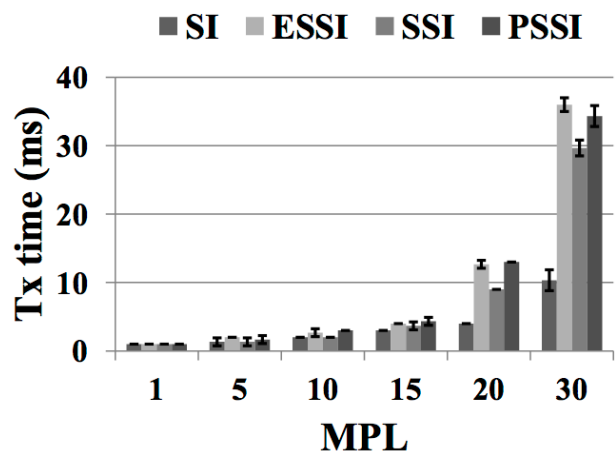
On 24 cores with a (75%-RO-25%-RU) workload

# Transaction Runtimes

## Read-Only Transactions

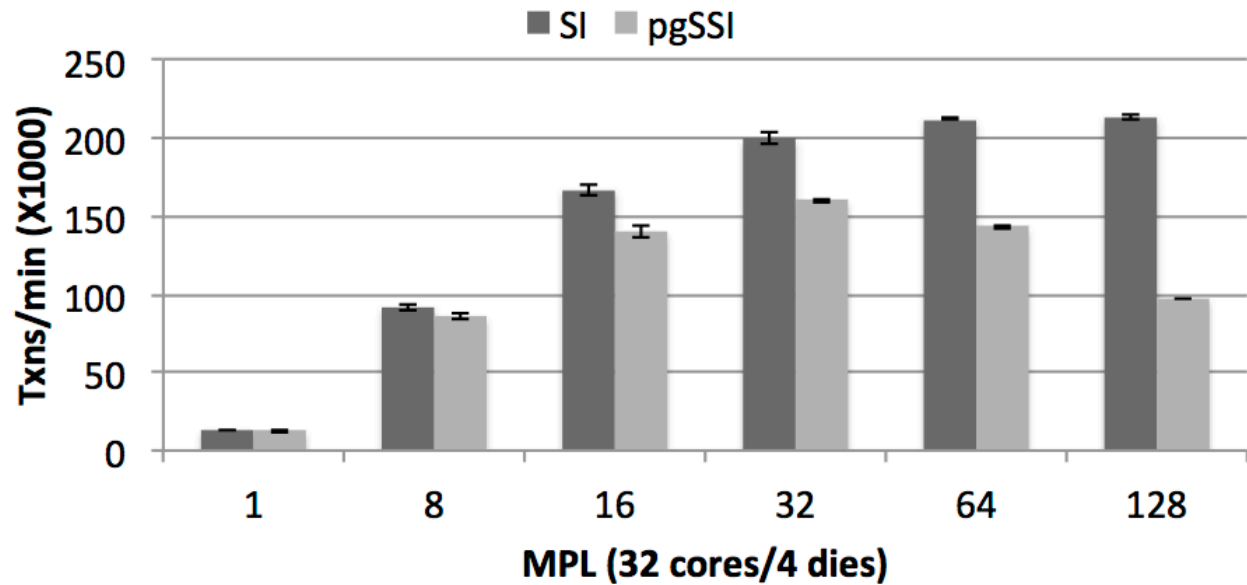


## Update Transactions



Transaction runtimes increase massively with MPL on 24 cores.

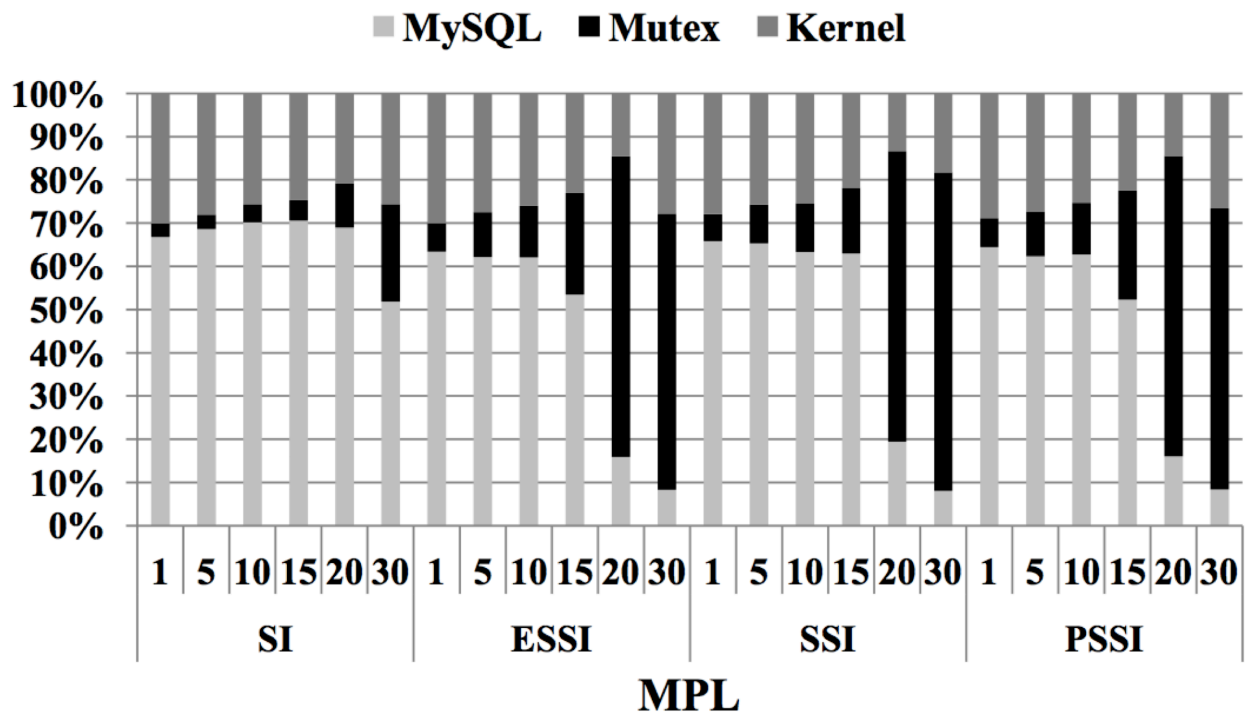
# Same Picture with Latest Postgres 9.2



On 32 cores with a (75%-RO-25%-RU) workload

## Why?

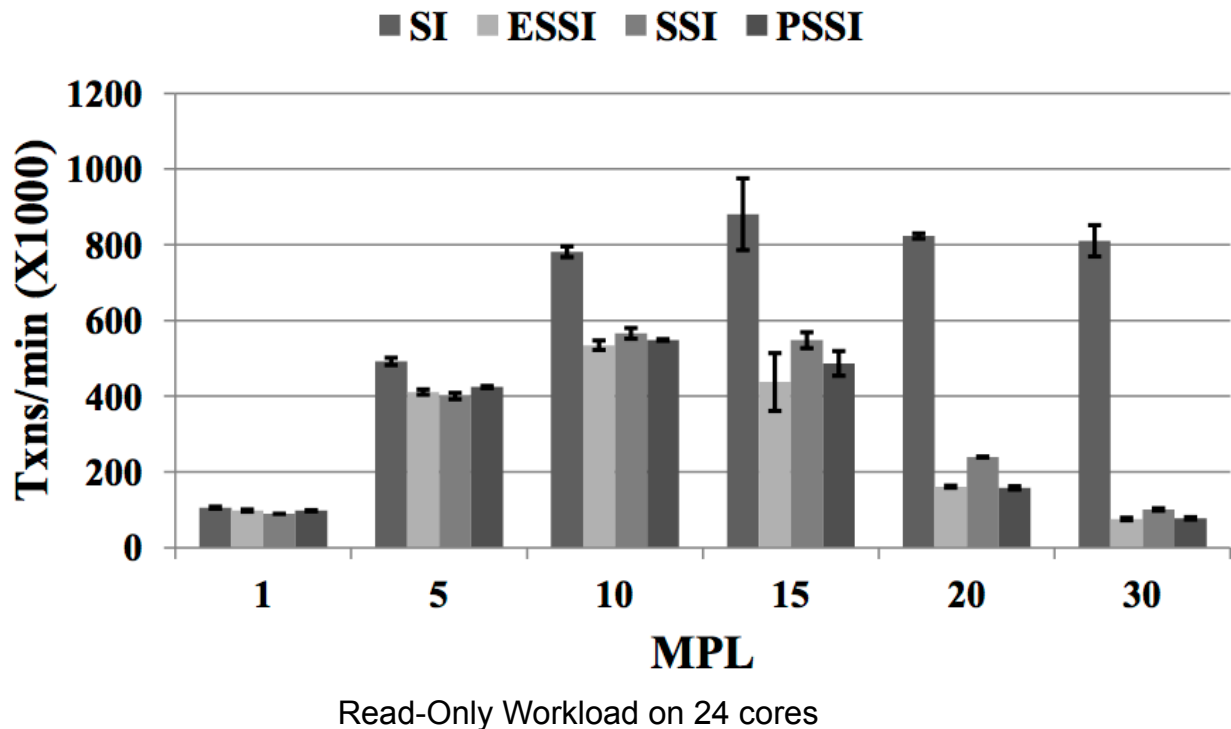
# Profiled Execution Times



## Reason: Mutex Contention

- More and more time is spend just waiting
- In order to avoid race conditions...
- ...but which 'race' ?
  
- And it can get even worse:

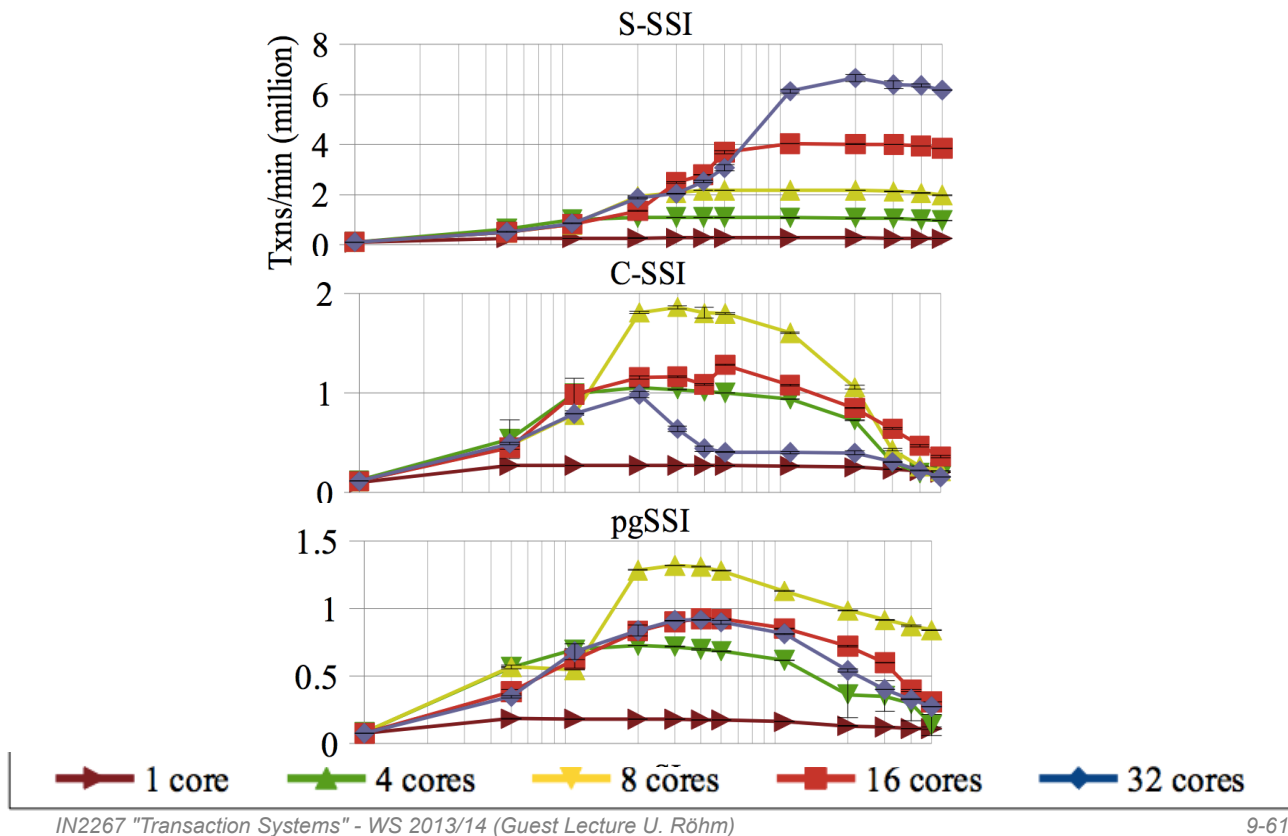
# “Readers Never Block”?



## Solution: Using Latch-free Data Structures for internal state of CC

- Latches are short-term locks (e.g. mutexes) that protect critical code sections from race conditions
  - ▶ E.g. only one thread is allowed to change the global transaction list
- Latch-Free Data Structures
  - ▶ Allow concurrent r/w access to in-memory data structures using atomic CPU operations such as Compare-And-Swap
  - ▶ Pro: Non-blocking, no latches needed anymore
  - ▶ Con: More complex; deletions require tombstones and some form of later garbage collection
- Proof of Concept: SSI with MySQL
  - ▶ Latch-free implementations of
    - read-write conflict checks, and
    - Consistent reads (read\_view)

# Performance Improvements



## Summary

- Transaction Management is the backbone of DBMSs
- Pessimistic Concurrency Control
  - ▶ lock-based concurrency control schemes detect conflicts between concurrent transactions by incompatible locks on data items
  - ▶ Strict-2PL: Avoids cascading aborts, but deadlocks possible
  - ▶ Deadlocks can either be prevented or detected.
- Optimistic Concurrency Control
  - ▶ aims to minimize CC overheads in an "optimistic" environment where reads are common and writes are rare.
  - ▶ Multiversion Timestamp CC is a variant which ensures that read-only transactions are never restarted; they can always read a suitable older version. Additional overhead of version maintenance.
  - ▶ Snapshot Isolation as popular CC nowadays
    - But does not guarantee serializable executions!

# References

- H. Berenson, P. Bernstein, J. Gray, J. Melton, E. O'Neil, P. O'Neil in SIGMOD1995: *"A Critique of ANSI SQL Isolation Levels"*
- A. Bernstein, P. Lewis and S. Lu in ICDE2000: *"Semantic Conditions for Correctness at Different Isolation Levels"*
- A. Fekete, D. Liarokapis, E. O'Neil, P. O'Neil, D. Shasha in TODS2005: *"Making Snapshot Isolation Serializable"*
- M. Alomari, M. Cahill, A. Fekete, U. Röhm in ICDE2008: *"The Cost of Serializability on Platforms That Use Snapshot Isolation"*
- **M. Cahill, U. Röhm and A. Fekete in SIGMOD2008:** *"Serialisable Isolation for Snapshot Databases"*
- **D. Ports and K. Grittner in VLDB2012:** *"Serialisable Snapshot Isolation in PostgreSQL"*
- **H. Jung, H. Han, A. Fekete, U. Röhm and H. Y. Yeom, DASFAA 2013:** *"Performance of Serializable Snapshot Isolation on Multicore Servers"*
- **H. Han, S. Park, H. Jung, A. Fekete, U. Röhm and H. Y. Yeom, ICDE2014:** *"Scalable Serializable Snapshot Isolation for Multicore Systems"*