

Precisely Serializable Snapshot Isolation (PSSI)

ICDE '11

Stephen Revilak, Patrick O'Neil, and Elizabeth O'Neil
University of Massachusetts Boston

Snapshot Isolation (SI)

- Introduced in 1995. Today, used in Oracle, Postgres, Microsoft SQL Server.
- Rather than locking, SI relies on multi-versioning. This gives SI two attractive properties:
 - Reads never wait for writes
 - Writes never wait for reads

How SI Works

- Data is versioned (labeled with the writer's transaction id).
- Each transaction has two timestamps: $\text{start}(T_i)$ and $\text{commit}(T_i)$.
 - T_i reads the most recent versions committed prior to $\text{start}(T_i)$ (the *snapshot*)
 - T_i 's writes become visible to new transactions that begin after $\text{commit}(T_i)$
- First committer wins to prevent lost updates.

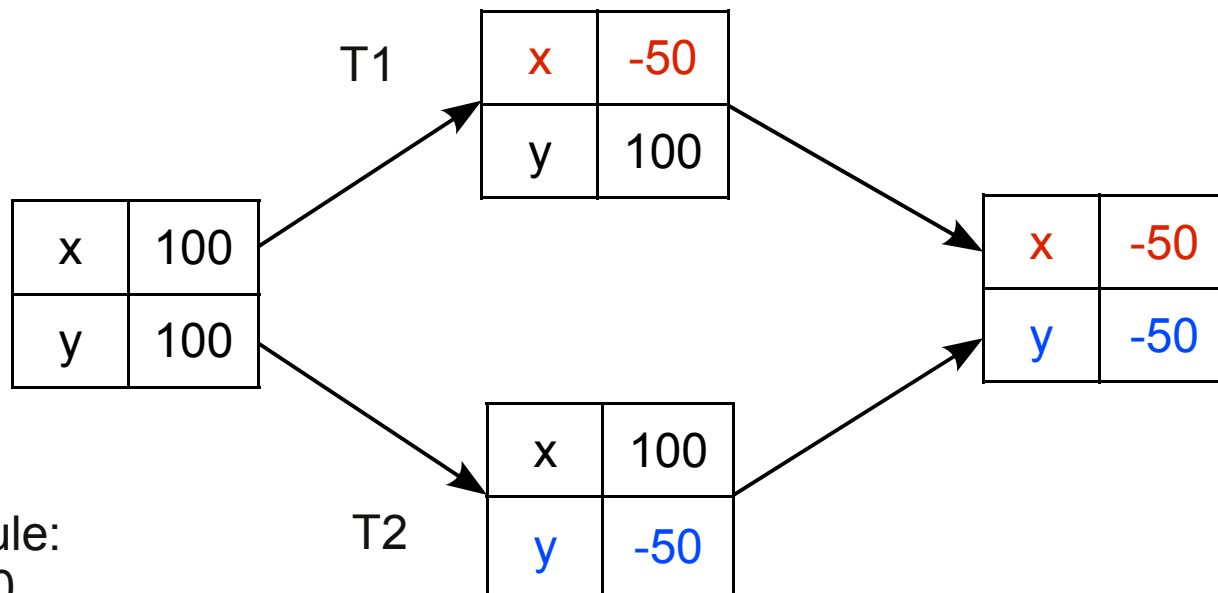
First Updater Wins (FUW)

- Two concurrent transactions: T_i , T_j .
- T_i updates x . Update takes $w_i(x)$ lock.
- T_j tries to update x . Two outcomes:
 - If T_i already committed, then T_j aborts immediately (no lock wait).
 - If T_i is still active, T_j goes into lock wait.
 - If T_i commits, then T_j aborts.
 - If T_i aborts, then T_j takes $w_j(x)$ lock and continues.

SI Anomalies - Write Skew

Example of an SI Anomaly:

H1: $r1(x0, 100), r1(y0, 100), r2(x0, 100), r2(y0, 100),$
 $w1(x1, -50), c1, w2(y2, -50), c2$



Business Rule:
 $(x + y) \geq 0$

Dependency Theory

- Write-write dependency: $T_i\text{-}ww \rightarrow T_j$
 - T_i writes x , T_j writes successor version of x .
- Write-read dependency: $T_i\text{-}wr \rightarrow T_j$
 - T_i writes x . T_j reads what T_i wrote.
 - Data item or predicate read, conflicting with data item write.
- Read-write anti-dependency: $T_i\text{-}rw \rightarrow T_j$
 - T_i reads x , T_j writes the successor version of x .
 - T_i , T_j may or may not be concurrent.

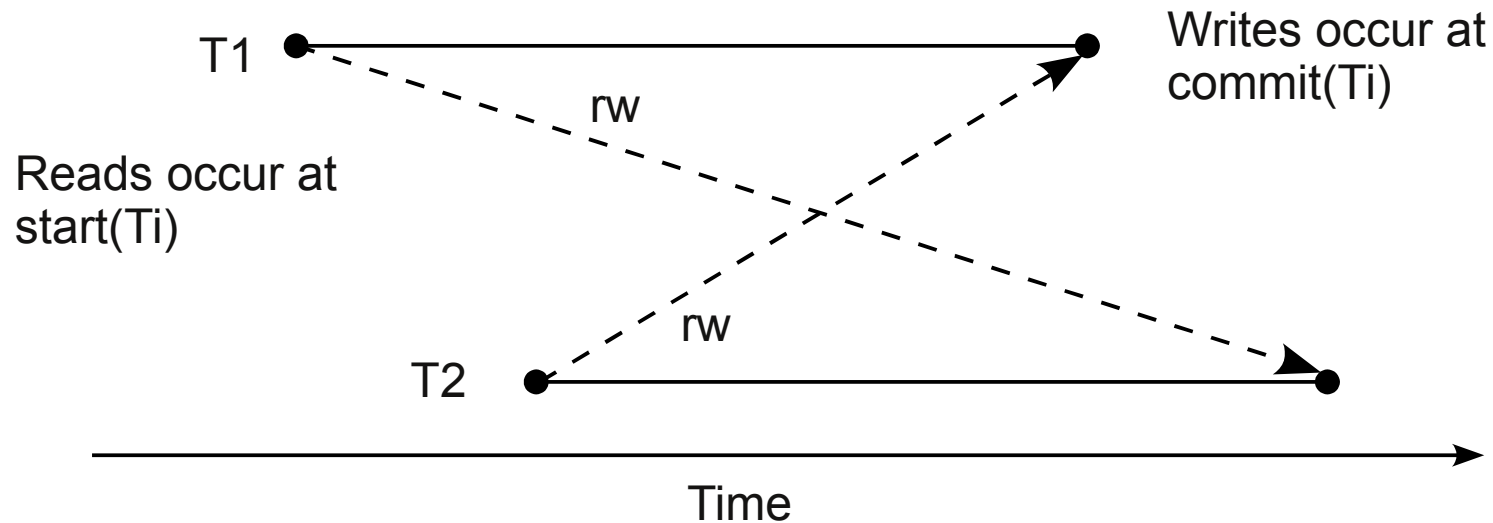
Dependency Serialization Graph

- We can model dependencies with a directed graph, called a Dependency Serialization Graph (DSG).
 - Each node is a committed transaction. Each edge is a dependency
- $DSG(H)$ is acyclic IFF history H is serializable.
- We can achieve serializability by preventing dependency cycles.

SI-RW Diagrams

H1: $r1(x0, 100), r1(y0, 100), r2(x0, 100), r2(y0, 100),$
 $w1(x1, -50), c1, w2(y2, -50), c2$

Below is a SI-RW diagram for H1.



Note: cycle $T1\text{-}rw \rightarrow T2\text{-}rw \rightarrow T1$ didn't happen until *after* $\text{commit}(T1)$.

PSSI's Strategy

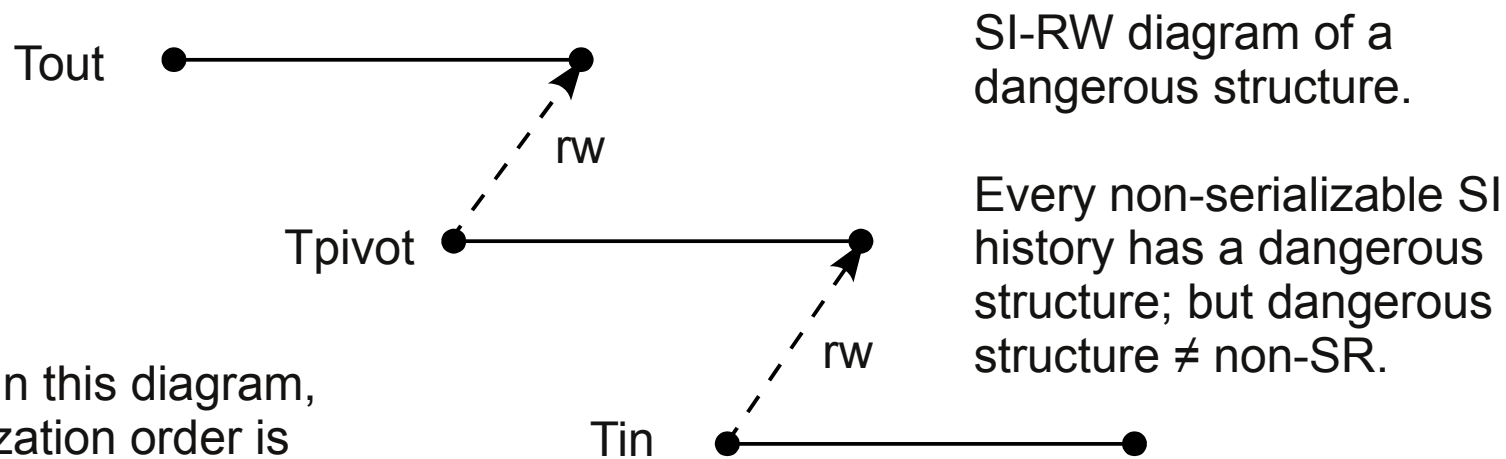
Ensure serializability by preventing cycles.

- Use a lock table to find dependencies.
- Use a cycle testing graph to find cycles.
- Use a variant of ARIES/IM to prevent phantom anomalies

We Implemented PSSI with MySQL 5.1.31's InnoDB Storage Engine.

ESSI - An Alternate Approach

- Rather than preventing cycles, ESSI prevents *essential dangerous structures*.
- ESSI guarantees Serializability, but is more conservative than PSSI.



SI-RW diagram of a dangerous structure.

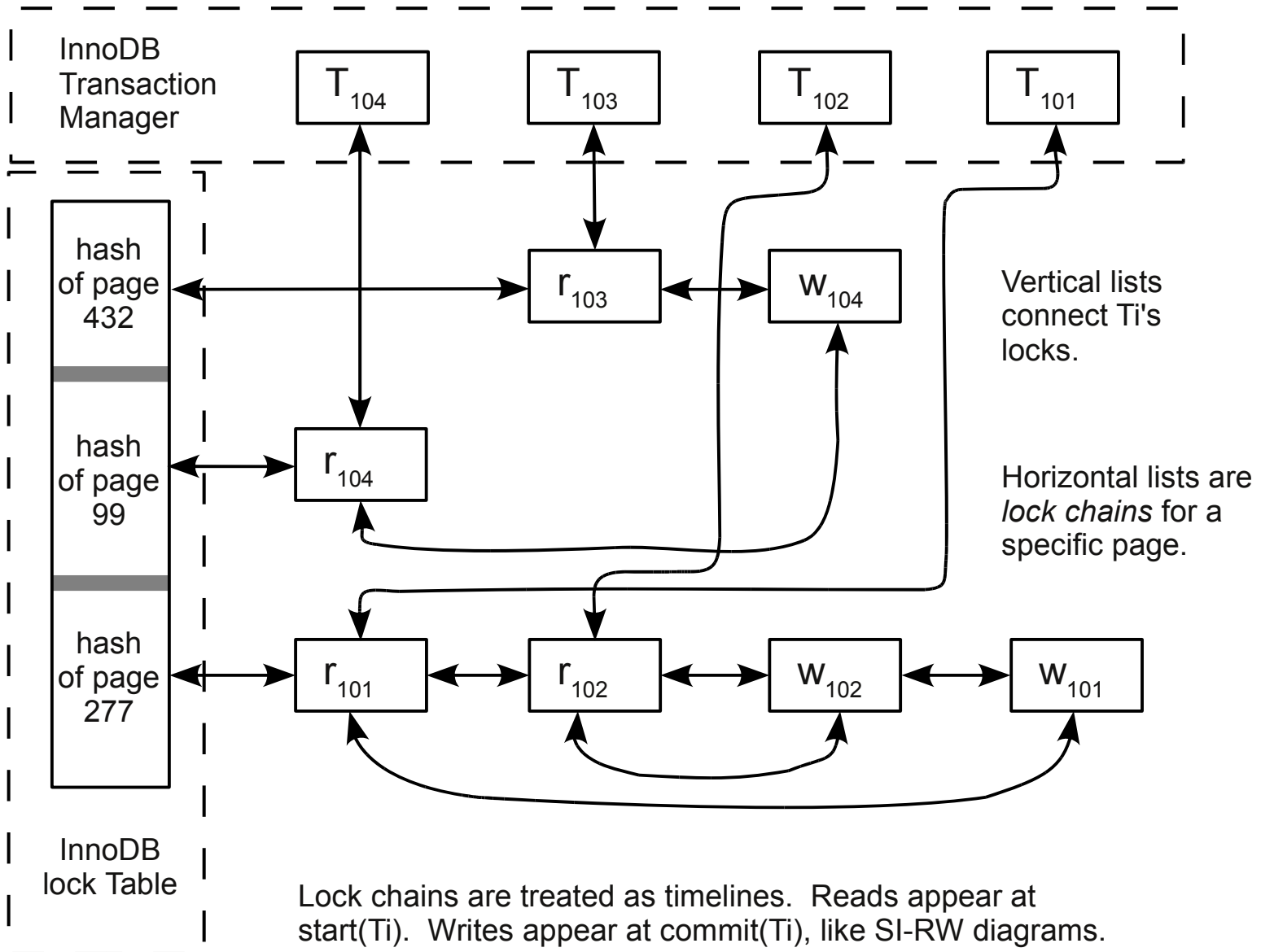
Every non-serializable SI history has a dangerous structure; but dangerous structure \neq non-SR.

Note: in this diagram, serialization order is opposite commit order.

PSSI Lock Table (1)

- Purpose: to detect conflicts, and find dependencies between transactions.
 - Keeps track of all reads and writes
- Looks like a traditional lock manager (but behaves differently).
- Requires separate notions of *lock conflict* and *lock wait*.
- Assumes that all dependencies can be found via conflicts on data items.

PSSI Lock Table (2)



Cycle Testing Graph: CTG

- Each node is a transaction T_i .
- Each edge is a $T_i \rightarrow T_j$ dependency.
- CTG is a *suffix* of a complete history, containing only:
 - Committed transactions with *potential* to become part of a cycle ("zombies").
 - The currently committing transaction T_k .
- T_k commits if $CTG \cup T_k$ is acyclic.

CTG Pruning

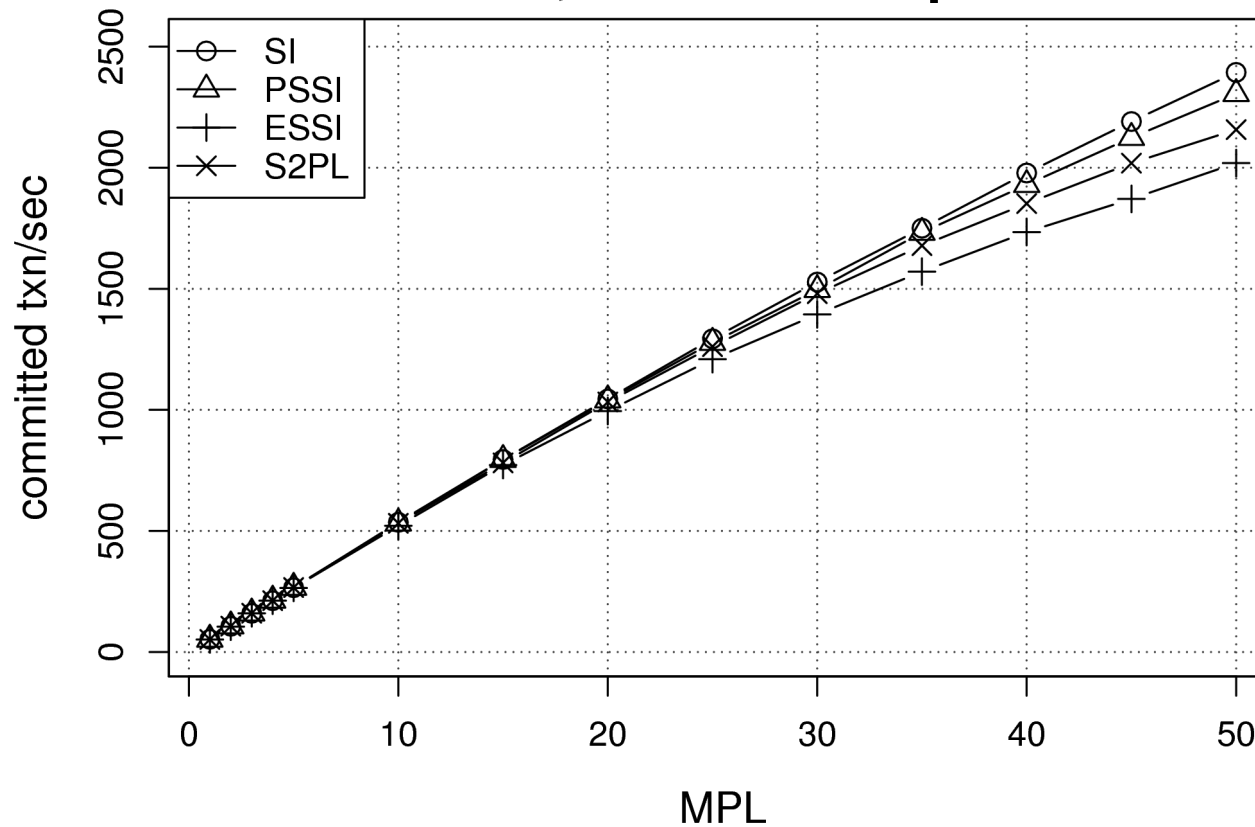
- If T_i can't become part of a (future) cycle, then we'd like to *prune* T_i from the CTG.
- We can prune T_i if:
 - T_i has no in edges, and
 - T_i committed before the oldest active transaction started.
- Pruning happens each time a transaction commits or aborts.

SI Cycles: a Parametrized Workload

- 1MM rows, 100-bytes per row.
 - kseq, krandseq, kval
- sKuN. Each transaction runs:
 - K select statements (by krandseq),
 - N update statements (by krandseq).
 - 3ms +/- 50% delay between statements.
- Vary K, N, MPL, size of hotspot.

s5u1 CTPS

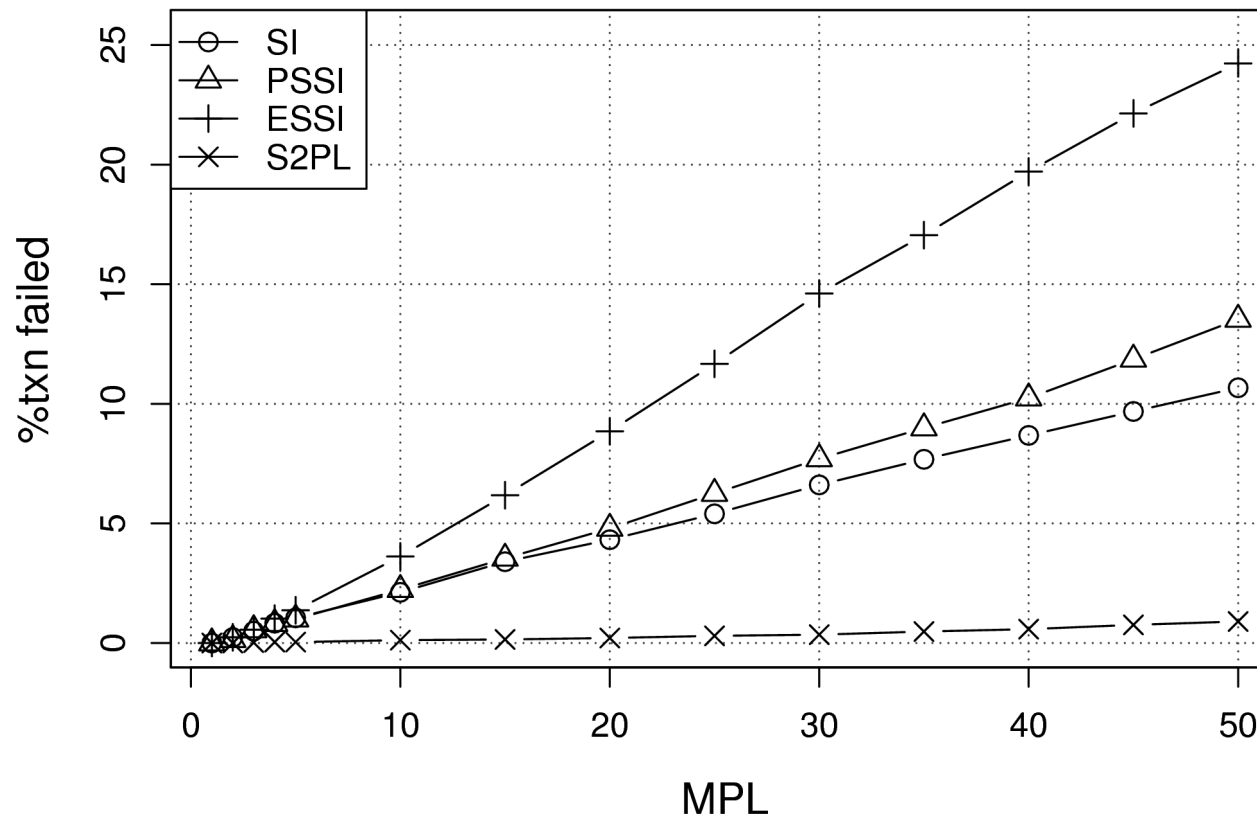
Committed Transactions Per Second s5u1, 400-row hotspot



SI performs best, but it's making mistakes.
The other three isolation levels are serializable.

s5u1 Abort Rates

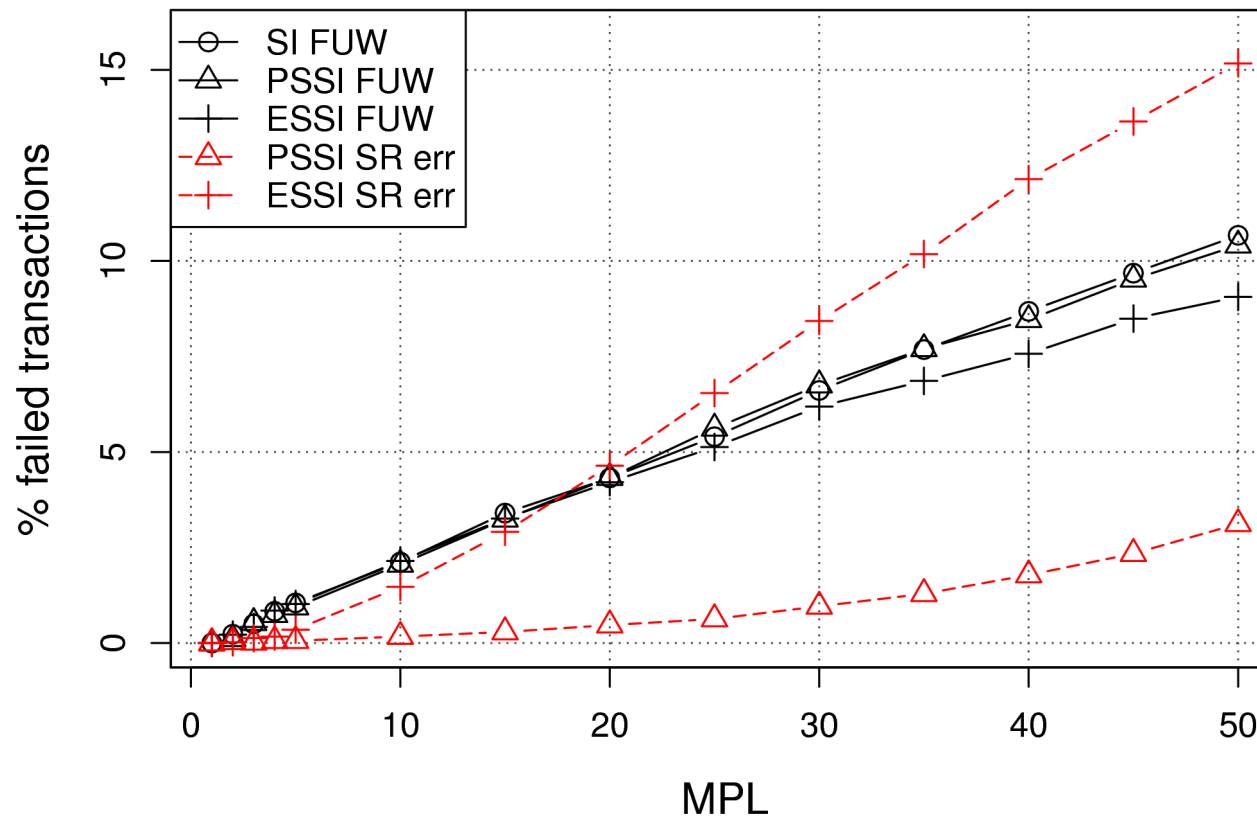
**%Failed Transactions
s5u1, 400-row hotspot**



SI, PSSI, ESSI throughput is primarily determined by error rate.
S2PL proceeds with caution.

s5u1 Abort Detail

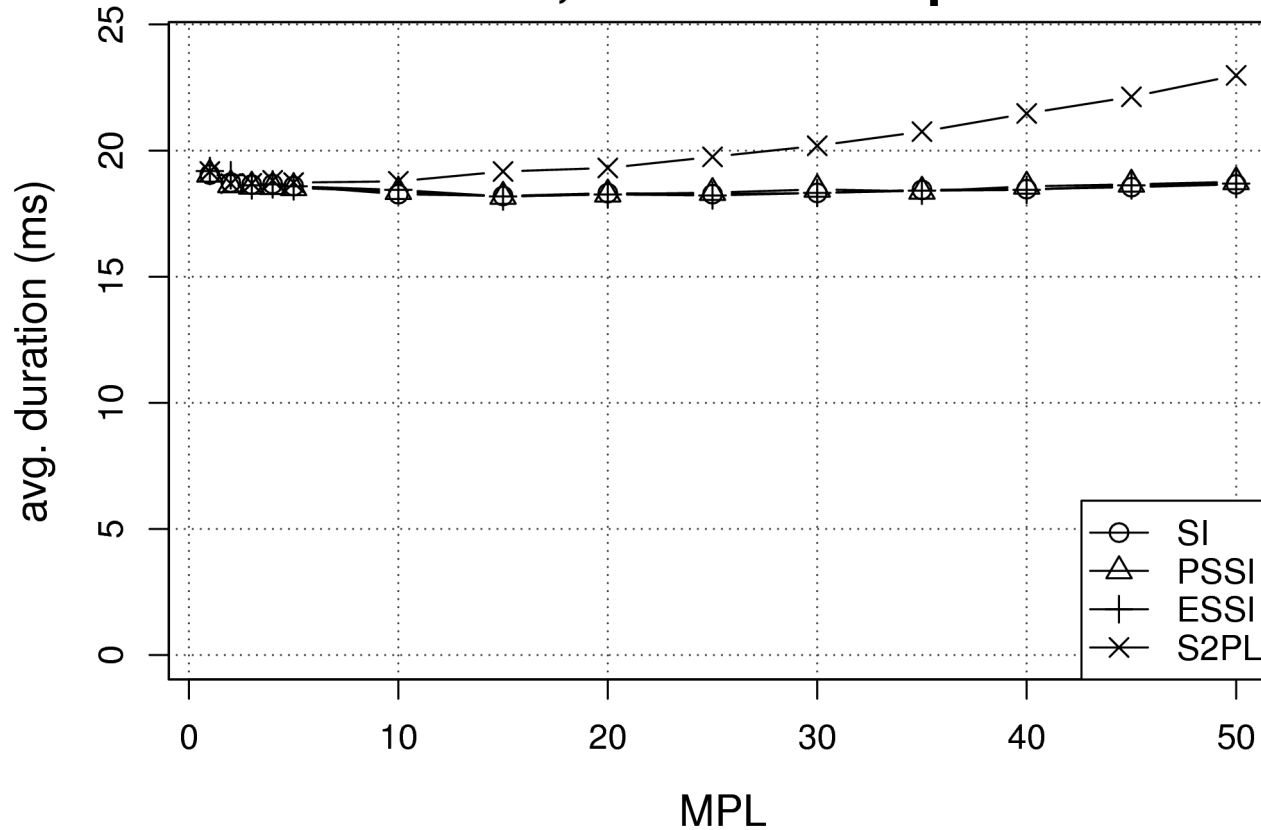
%Failures by Type
s5u1, 400-row hotspot



Black=FUW errors. Red=Serialization Errors.

s5u1 Durations

**Avg. Duration of Committed Txns
s5u1, 400-row hotspot**



S2PL duration increases due to lock waits.

SI, PSSI, ESSI duration increases slightly (only ww conflicts wait)

Summary

- Snapshot Isolation: good throughput, but doesn't provide serializability.
- PSSI: makes SI Serializable, with a minimum of false-positive aborts.
 - Lock Table to find dependencies.
 - Cycle Testing Graph to find cycles.
 - Pruning to keep the CTG small.

Thank You

Extra slides

Challenges in Making SI Serializable

- T_i might not become part of a cycle until after $\text{commit}(T_i)$
 - DBMS must retain information about transactions after they've committed
- In general, commit order does not match serialization order
(We'll see examples of this shortly)

InnoDB and Durability

- Controlled via `innodb_flush_log_at_trx_commit`.
- *flush=2*. Write logs at commit, flush logs asynchronously 1x/second.
 - High performance. Not very durable.
- *flush=1*. Write and flush logs at commit (group commit)
 - Lower performance. *Almost* durable.

InnoDB and Group Commit

- Our paper states: "InnoDB does not offer group commit". MySQL 5.1.31's InnoDB offers group commit, but it's broken in some configurations.
- cf. MySQL bug #13669: *Group commit is broken in 5.0.*

(Fixed in the current InnoDB, but not the version we used for development)

Testing Environment

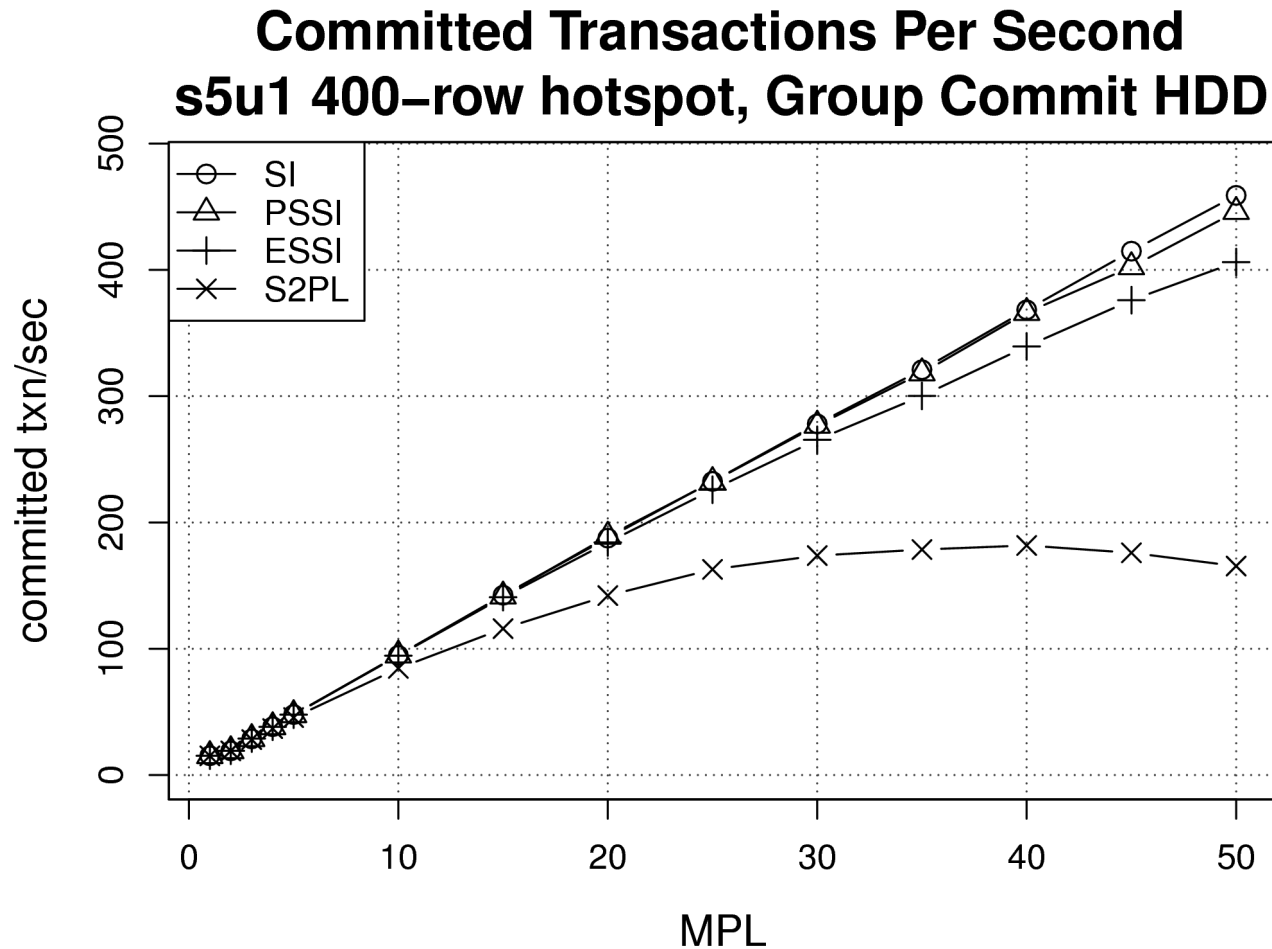
Two Linux machines

- Quad Core server. 2.67 GHz, 4GB ram, 64-bit OS.
- Dual Core client. 2.53 GHz, 2GB ram, 32-bit OS.
- Java benchmarking program. Plain JDBC + MySQL Connector/J

PSSI Commit Sequence

- Find dependencies that $\text{commit}(T_i)$ would cause.
 - Add T_i + dependencies to CTG.
- Perform a cycle test.
 - If cycle, abort. Otherwise, commit continues.
- Rest of commit (write logs, release write locks, handle FUW, etc).
- Prune the CTG.

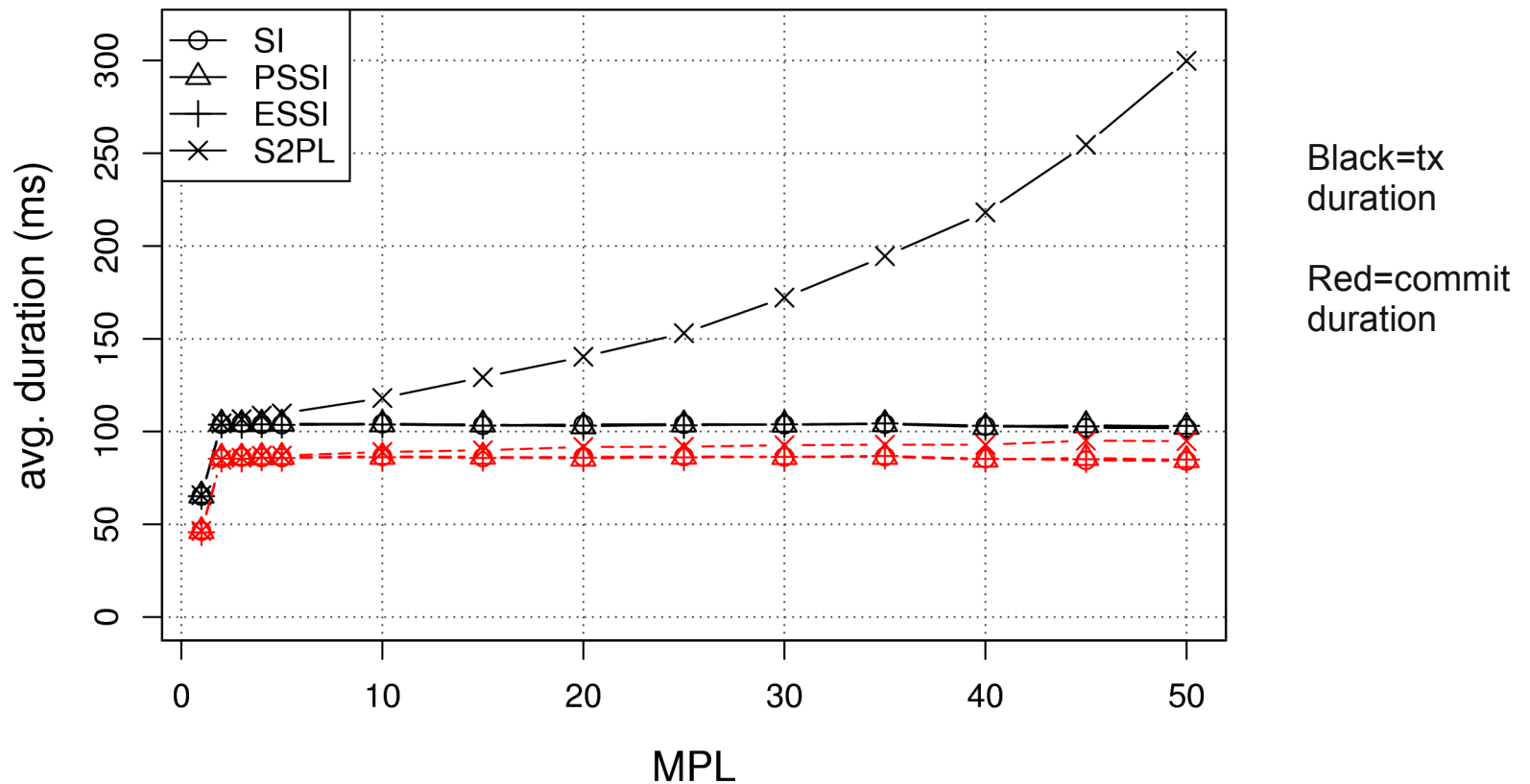
Group Commit HDD CTPS



Nearly 80% lower throughput (!) than flushing 1x/sec.
Note: InnoDB modified so S2PL releases locks *after* flush.

Group Commit HDD Durations

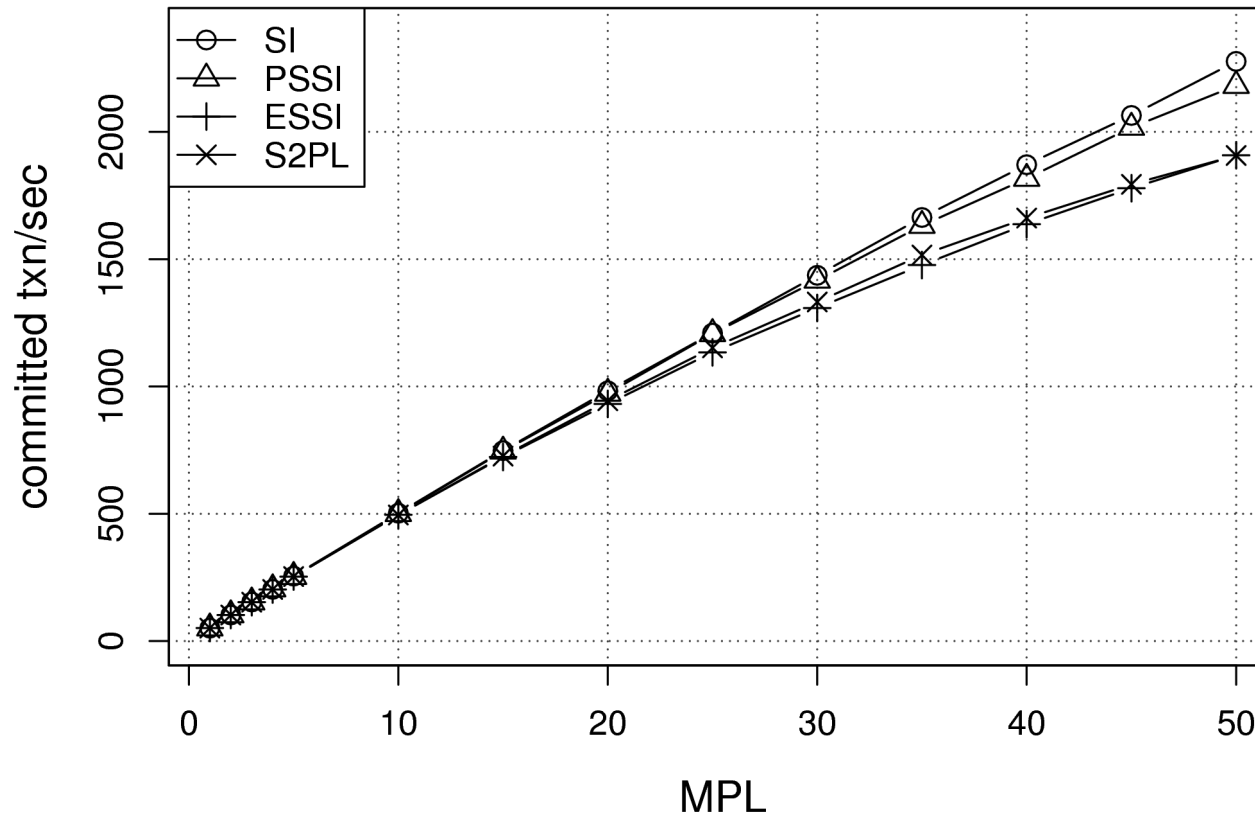
Transaction and Commit Duration
s5u1 400-row hotspot, Group Commit HDD



SI, PSSI, ESSI duration now > 100ms; much longer than flush=2's 18ms

Group Commit SSD CTPS

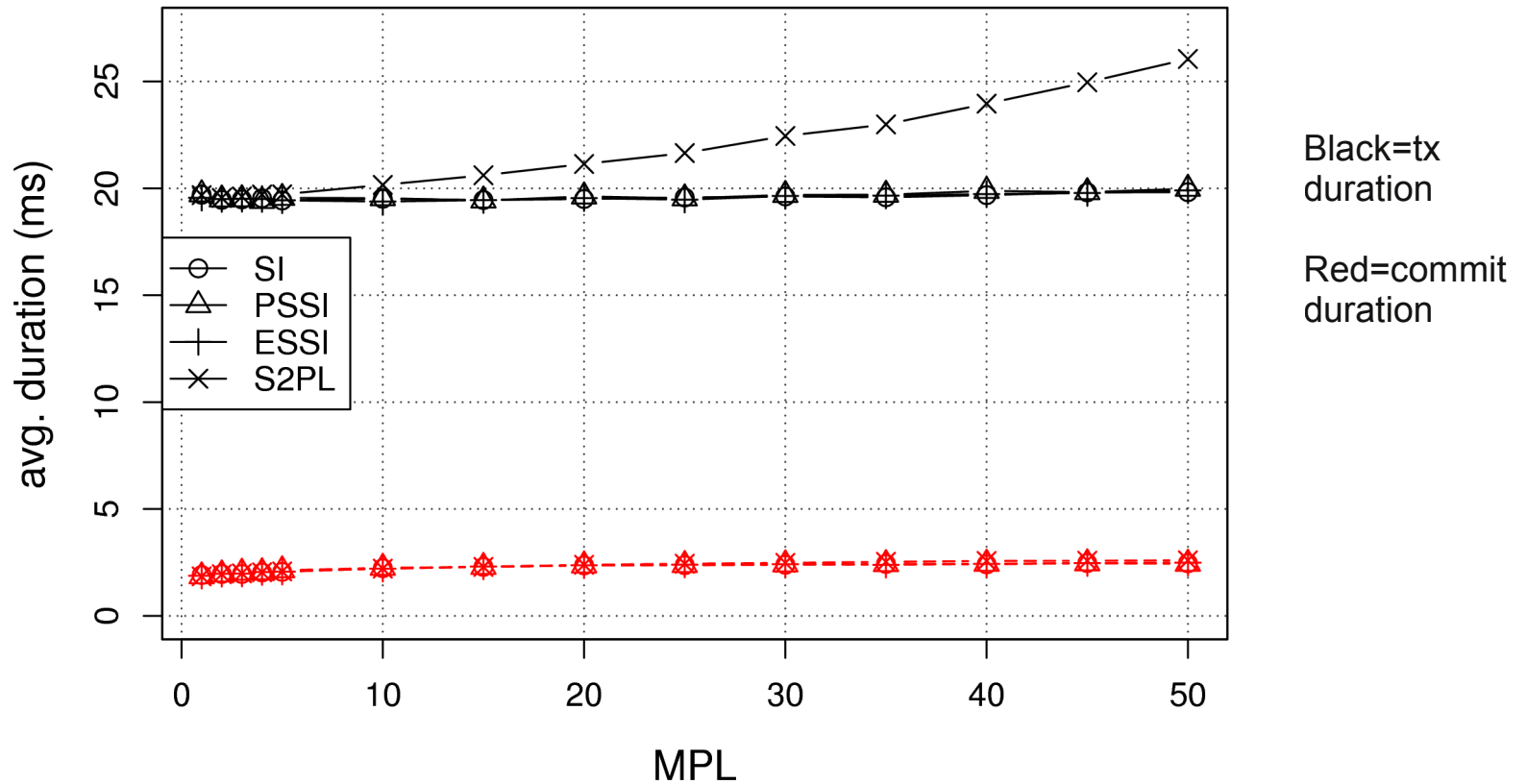
**Committed Transactions Per Second
s5u1 400-row hotspot, Group Commit SSD**



With SSD, throughput almost as good as flushing 1x/sec

Group Commit SSD Durations

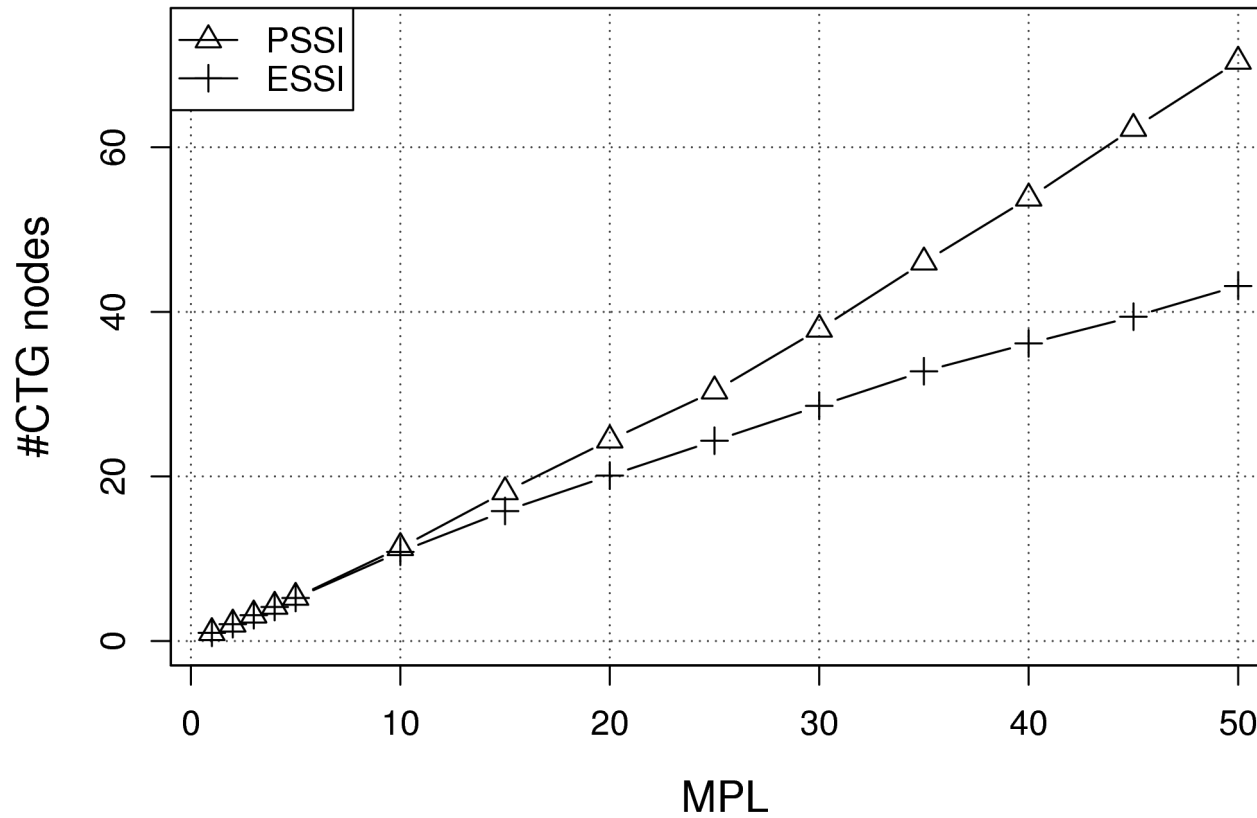
Transaction and Commit Duration
s5u1 400-row hotspot, Group Commit SSD



Transaction duration ~19ms, close to original ~18ms

Avg. CTG Size

Average CTG Size
s5u1, 400-row hotspot



Cycle Length Distribution

