# Distributed Transaction Manager
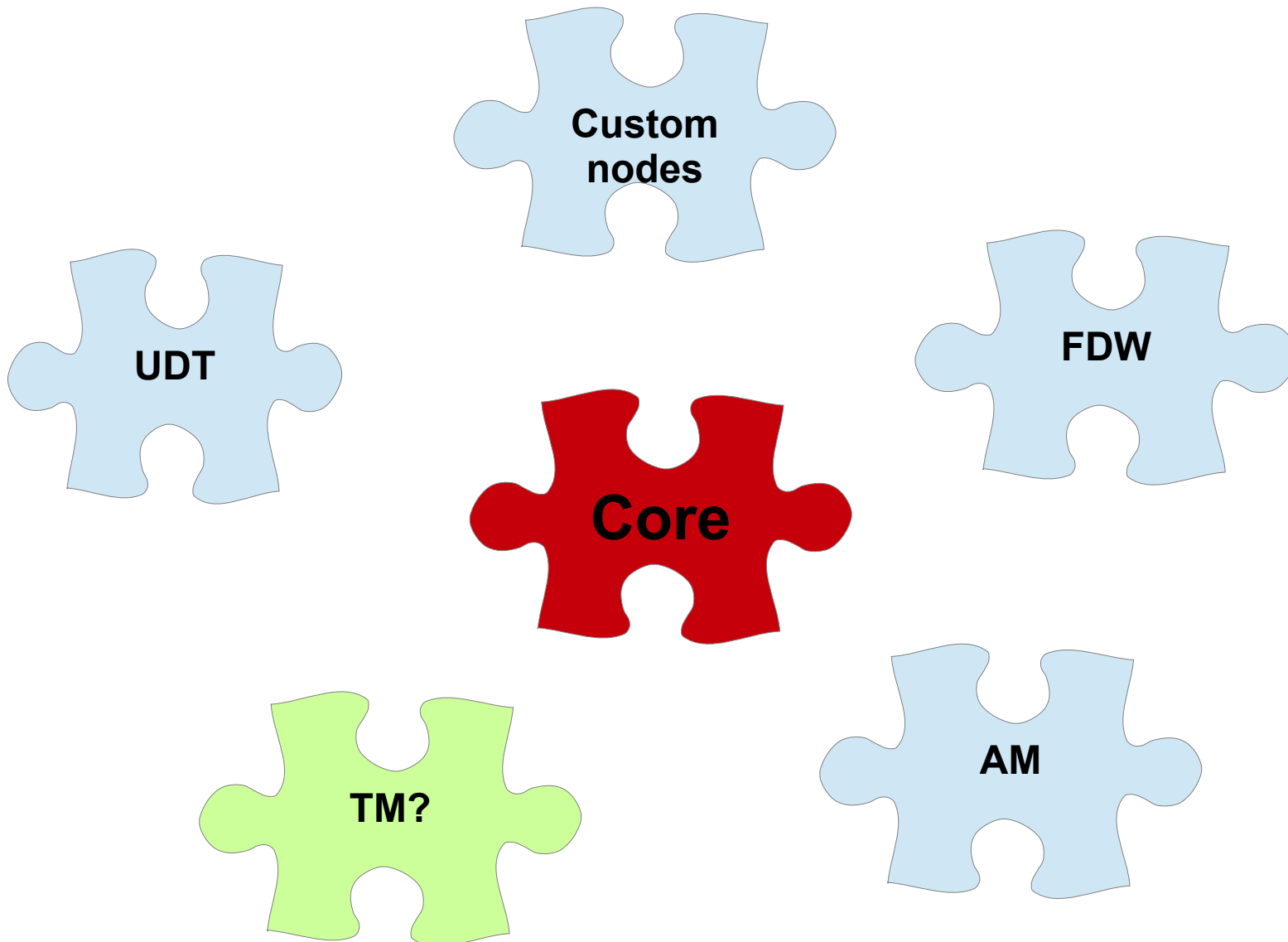
Stas Kelvich
Konstantin Knizhnik
Konstantin Pan

# Мы пойдём другим путём...

# Pluggable transaction API

# eXtensible Transaction API

- XidStatus (*GetTransactionStatus)(TransactionId xid, XLogRecPtr *lsn);

- void (*SetTransactionStatus)(TransactionId xid, int nsubxids, TransactionId *subxids, XidStatus status, XLogRecPtr lsn);

- Snapshot (*GetSnapshot)(Snapshot snapshot);

- TransactionId (*GetNewTransactionId)(bool isSubXact);

- TransactionId (*GetOldestXmin)(Relation rel, bool ignoreVacuum);

- bool (*IsInProgress)(TransactionId xid);

- TransactionId (*GetGlobalTransactionId)(void);

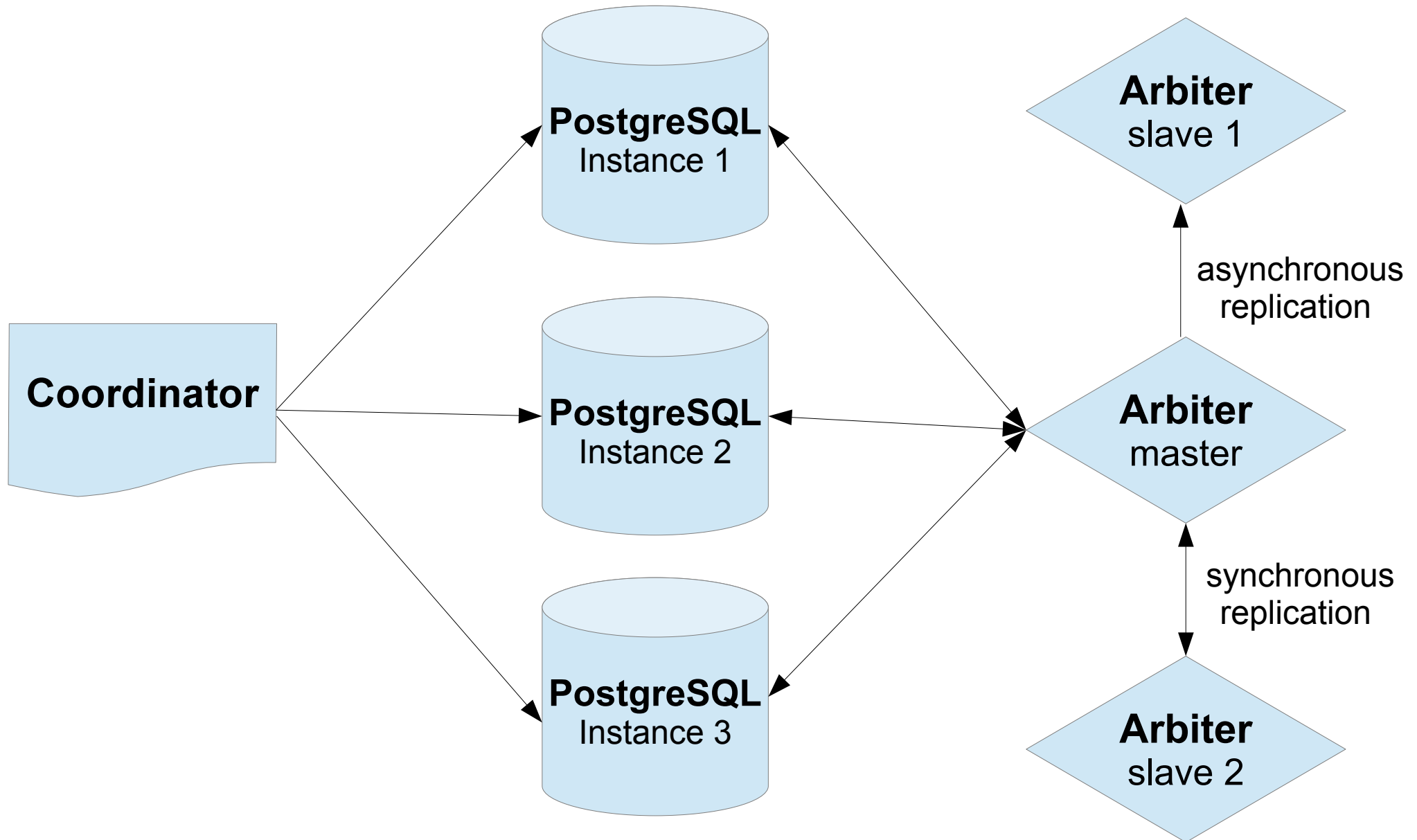- bool (*IsInSnapshot)(TransactionId xid, Snapshot snapshot);

# New commit callback events

- XACT_EVENT_START,

- XACT_EVENT_COMMIT,

- XACT_EVENT_PARALLEL_COMMIT,

- XACT_EVENT_ABORT,

- XACT_EVENT_PARALLEL_ABORT,

- XACT_EVENT_PREPARE,

- XACT_EVENT_PRE_COMMIT,

- XACT_EVENT_PARALLEL_PRE_COMMIT,

- XACT_EVENT_PRE_PREPARE,

- XACT_EVENT_COMMIT_PREPARED,

- XACT_EVENT_ABORT_PREPARED

# Different DTM implementations

| | Local transactions | 2PC | Arbiter | Examples |
|---|---|---|---|---|
| Snapshot sharing | | | ⌄ | XL, DTM |
| Timestamp | ⌄ | ⌄ | | Spanner, Cockroach, tsDTM |
| Incremental | ⌄ | | ⌄ | SAP HANA |

# DTM architecture

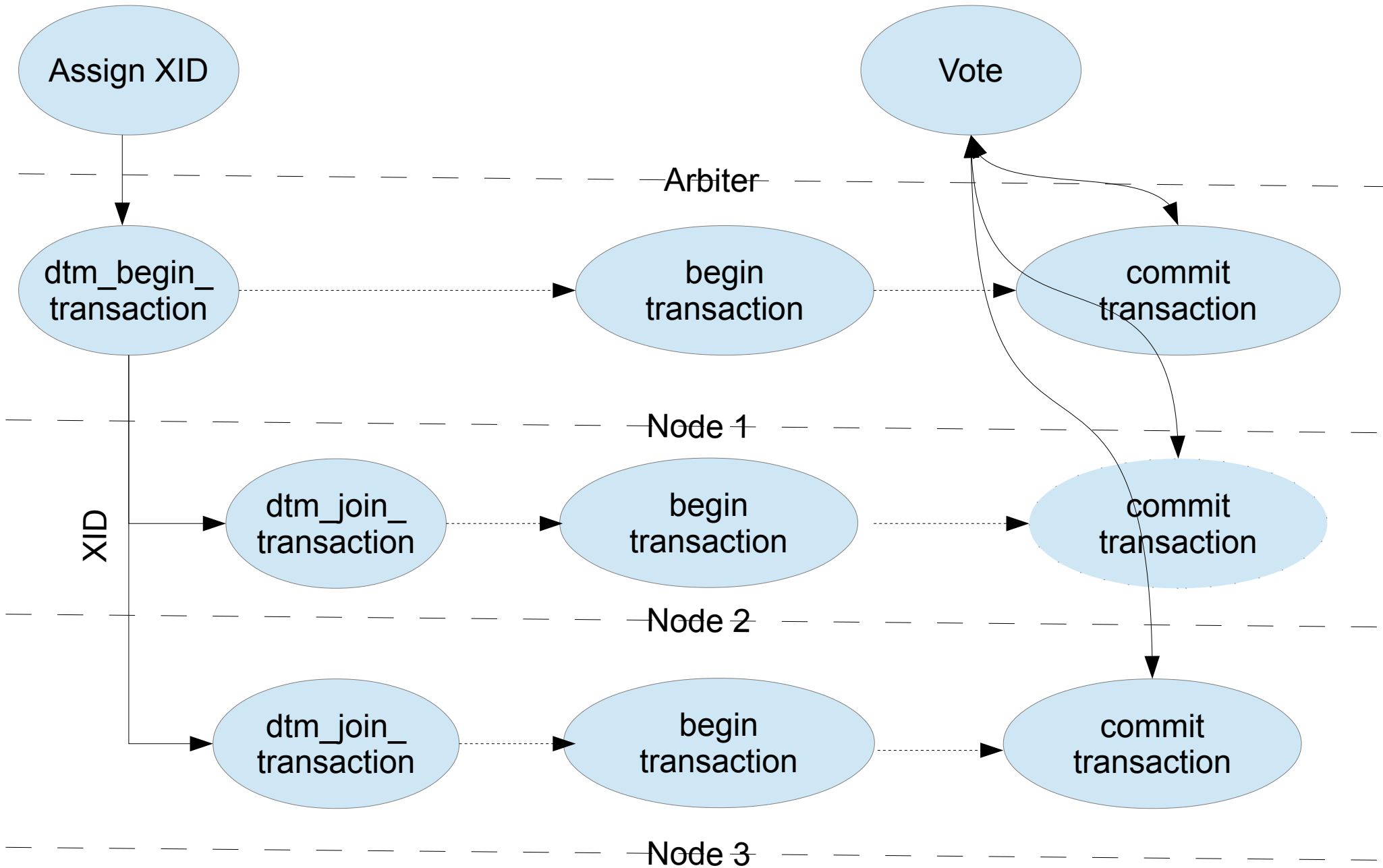# DTM from client's point of view

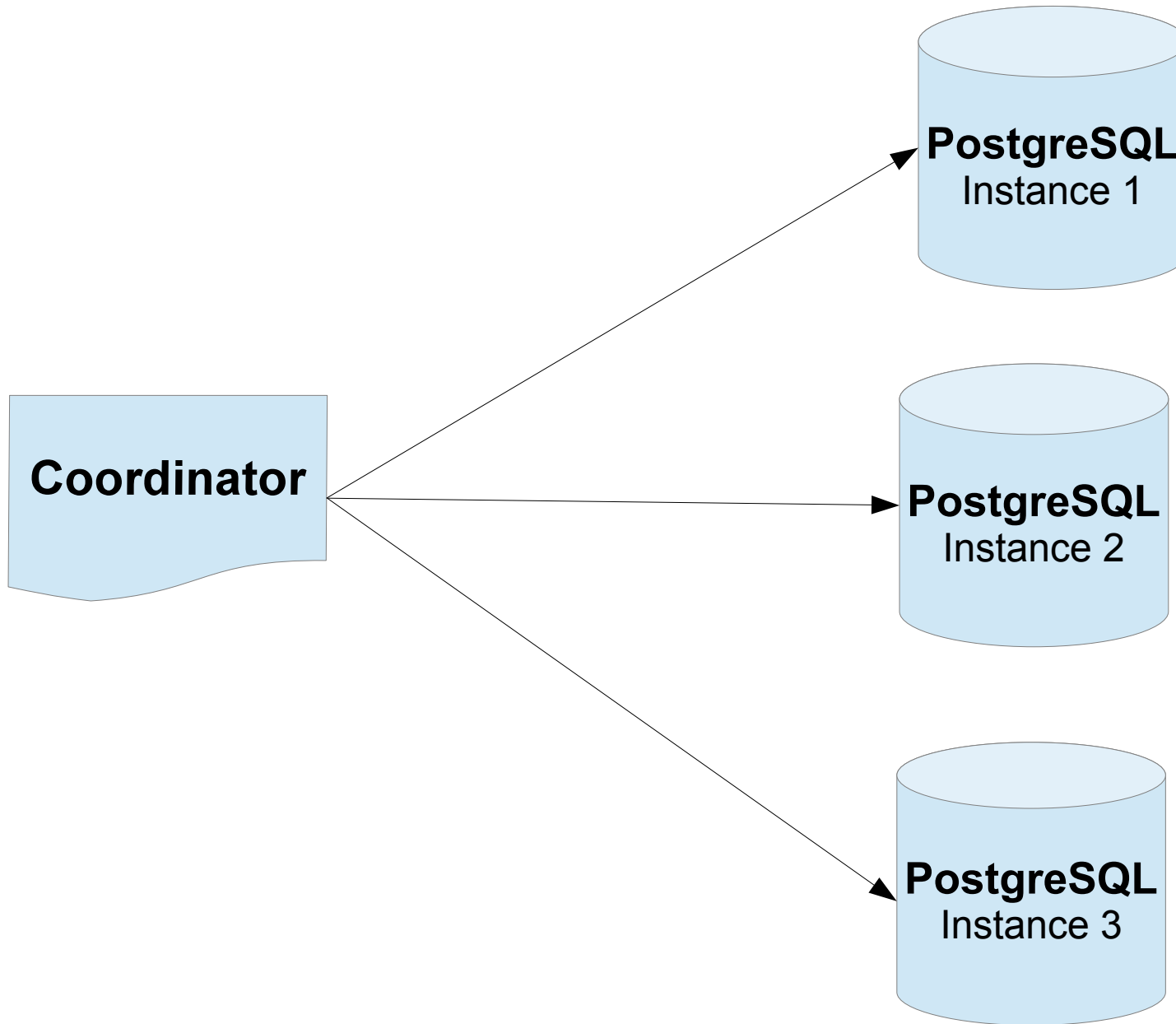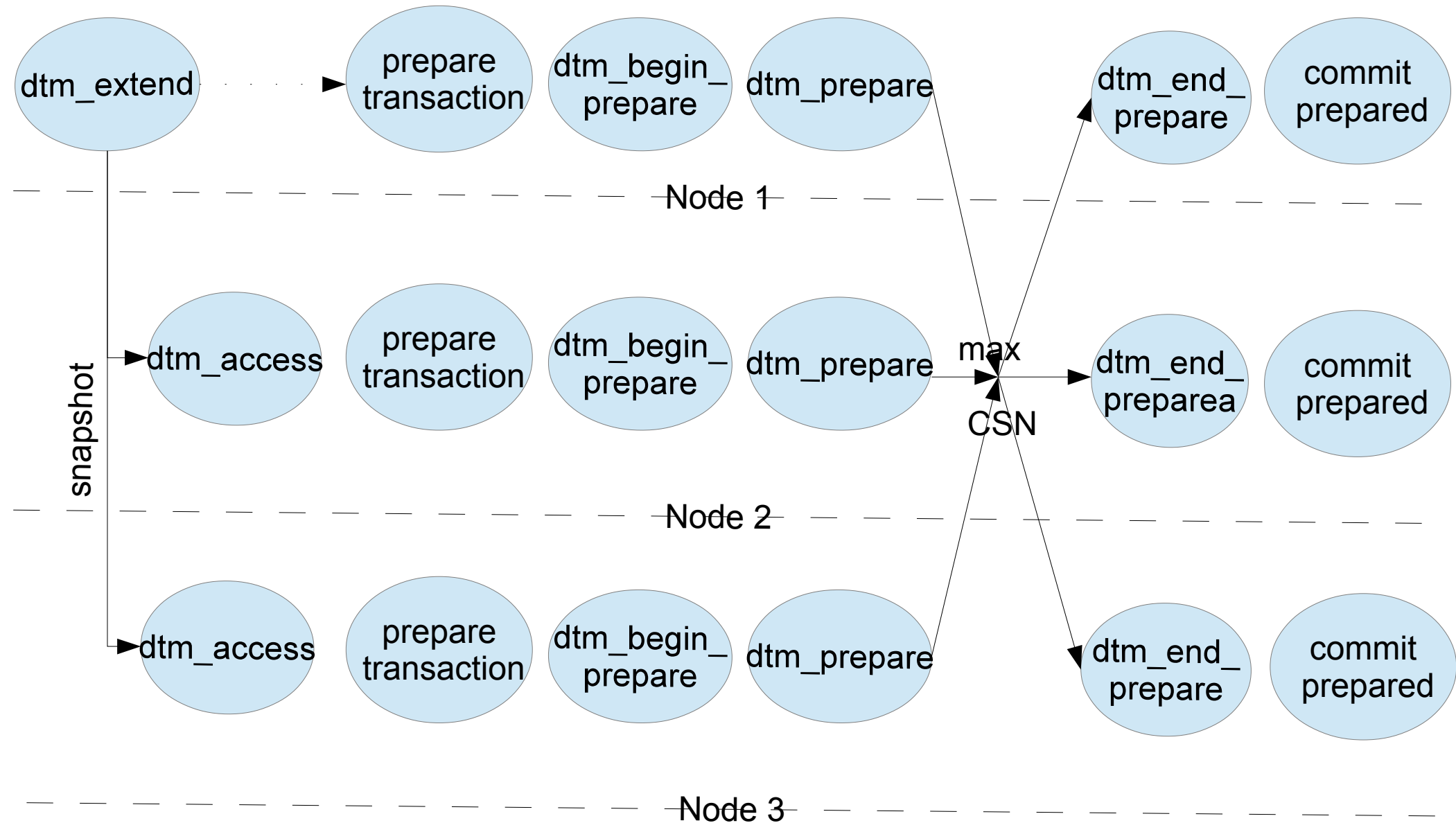| Primary server | Secondary server |
|---|---|
| create extension pg_dtm; | create extension pg_dtm; |
| select dtm_begin_transaction(); <br><br> begin transaction; <br><br> update...; <br><br> commit; | select dtm_join_transaction(xid); <br><br> begin transaction; <br><br> ppdate...; <br><br> commit; |

# DTM from client's point of view

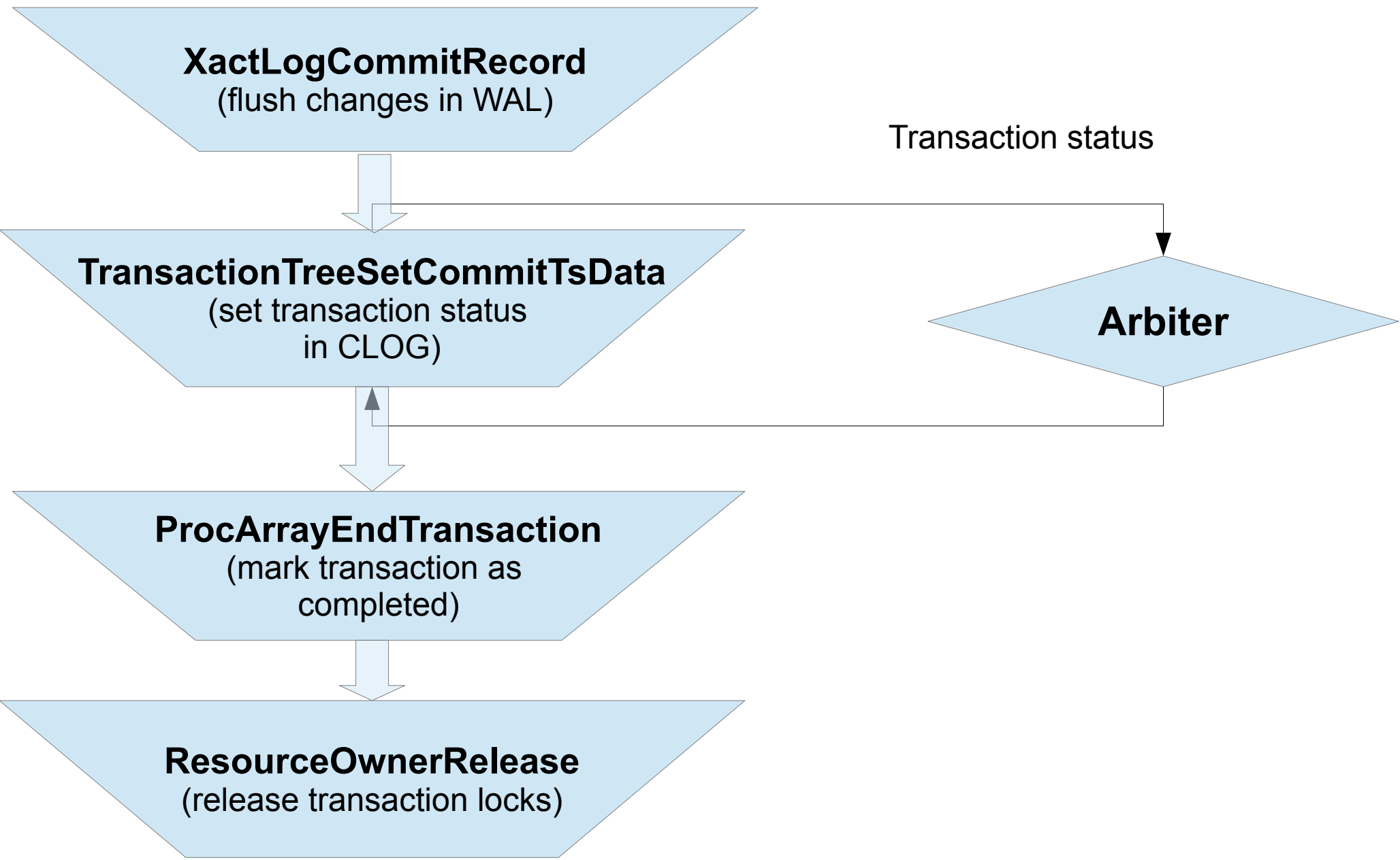| Primary server | Secondary server |
|---|---|
| create extension pg_dtm; | create extension pg_dtm; |
| select dtm_begin_transaction();<br><br>begin transaction;<br><br>update...;<br><br>commit; | select dtm_join_transaction(xid);<br><br>begin transaction;<br><br>ppdate...;<br><br>commit; |

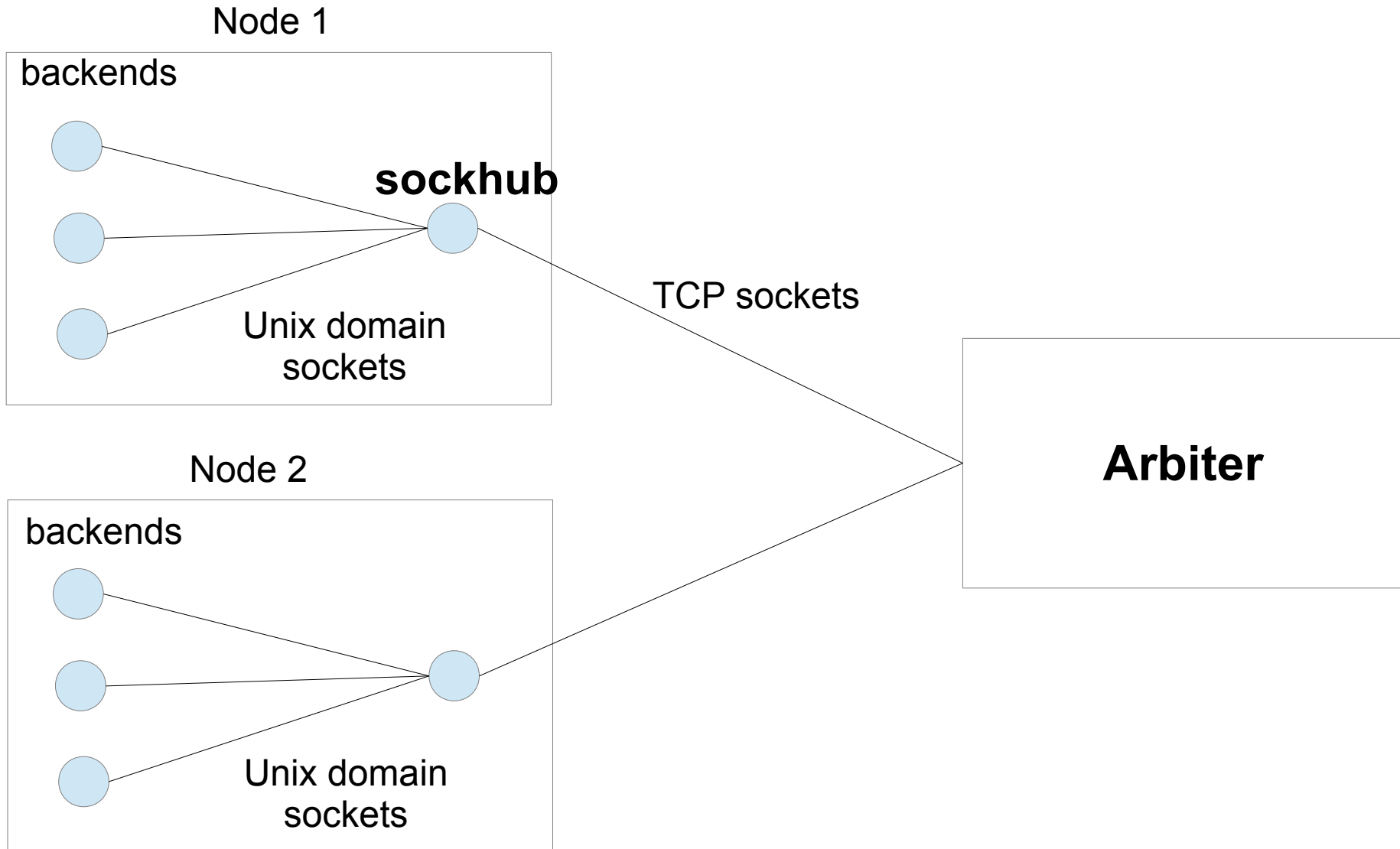# DTM transaction control flow
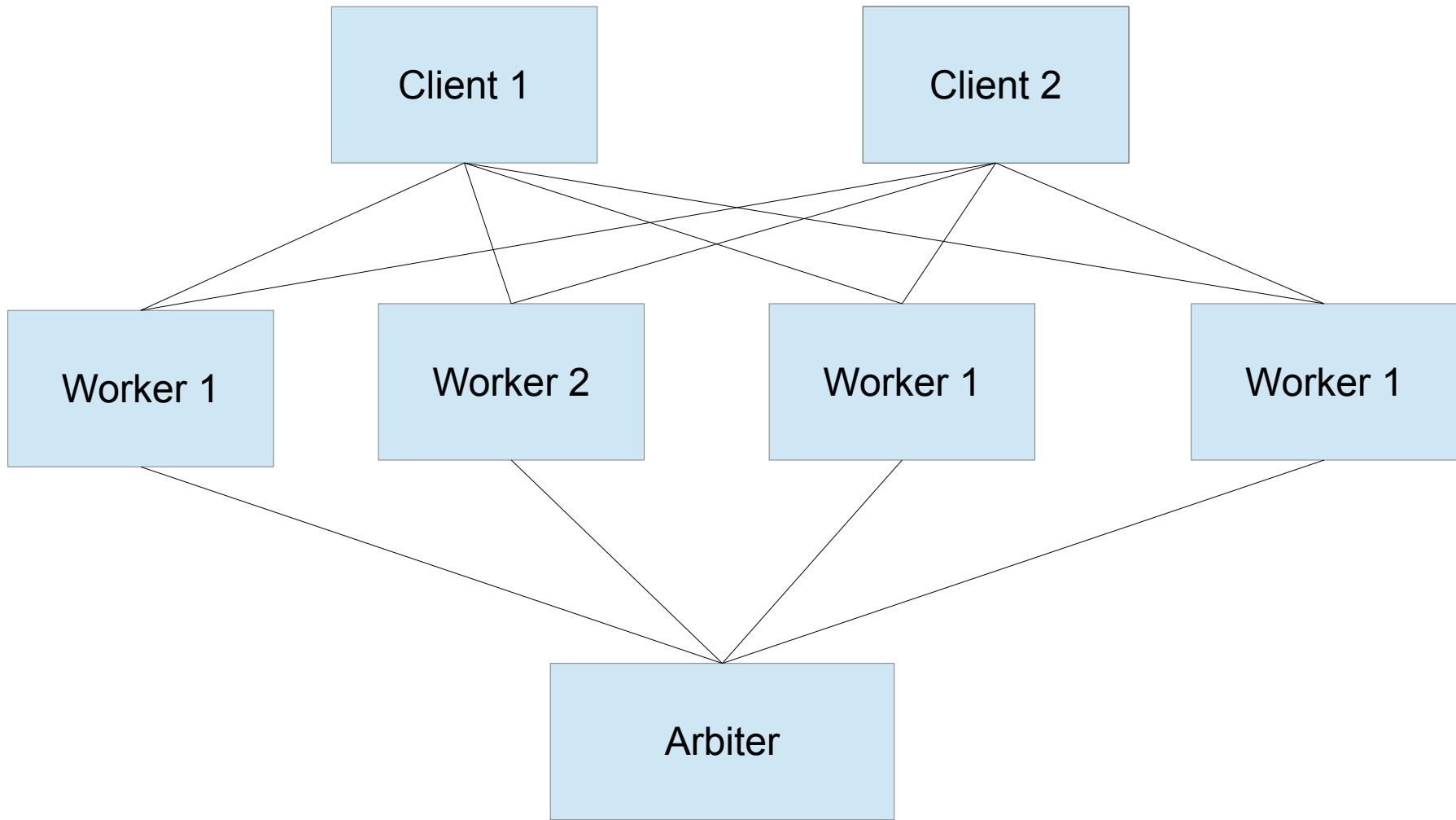
# tsDTM architecture

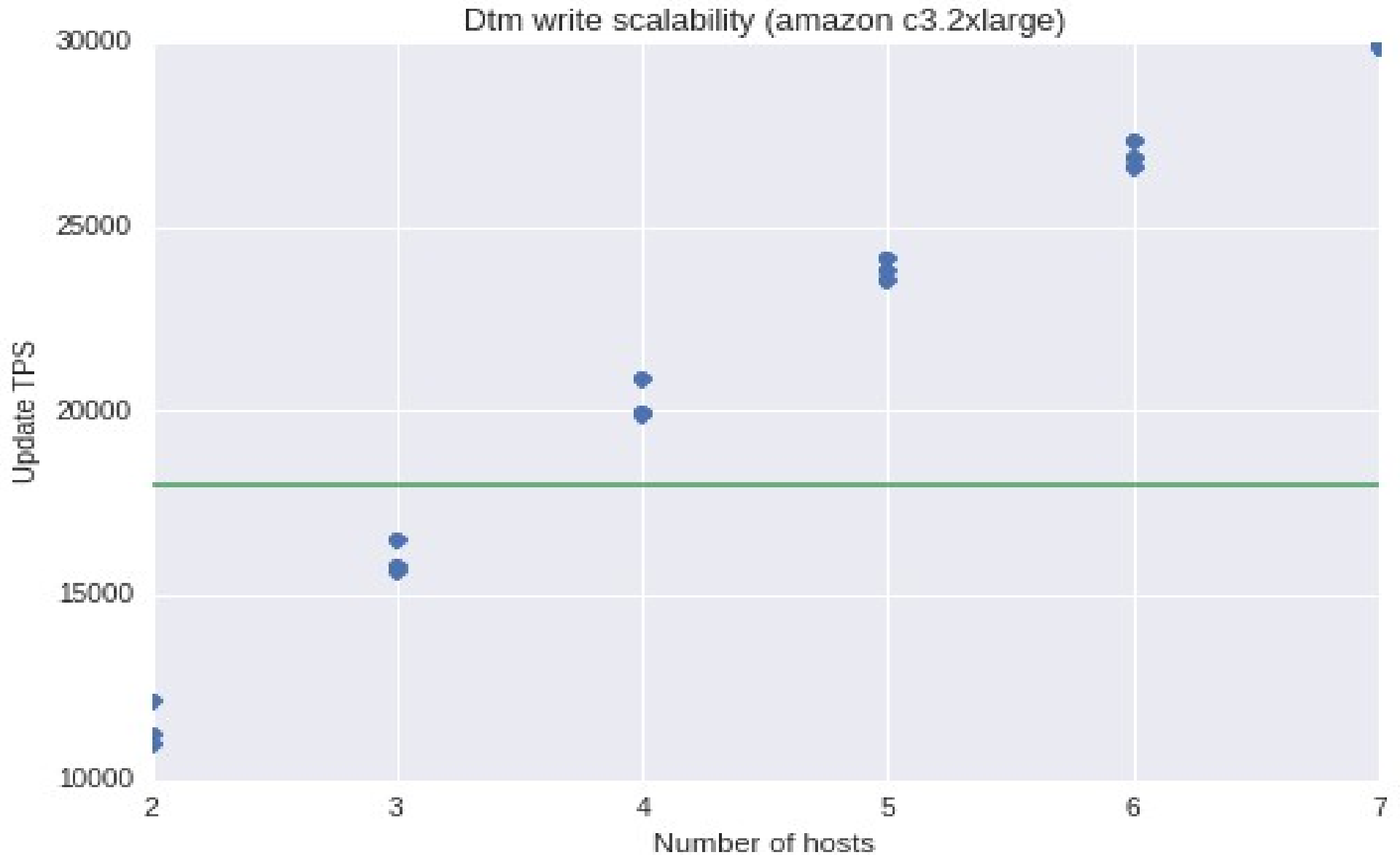# tsDTM transaction control flow

# Lightweight two-phase commit

**XactLogCommitRecord**
(flush changes in WAL)

**TransactionTreeSetCommitTsData**
(set transaction status
in CLOG)

Transaction status

**Arbiter**

**ProcArrayEndTransaction**
(mark transaction as
completed)

**ResourceOwnerRelease**
(release transaction locks)

# Multiplexing

Node 1

backends

**sockhub**

Unix domain
sockets

TCP sockets

Node 2

backends

Unix domain
sockets

**Arbiter**

# Test configuration

# DTM scalability

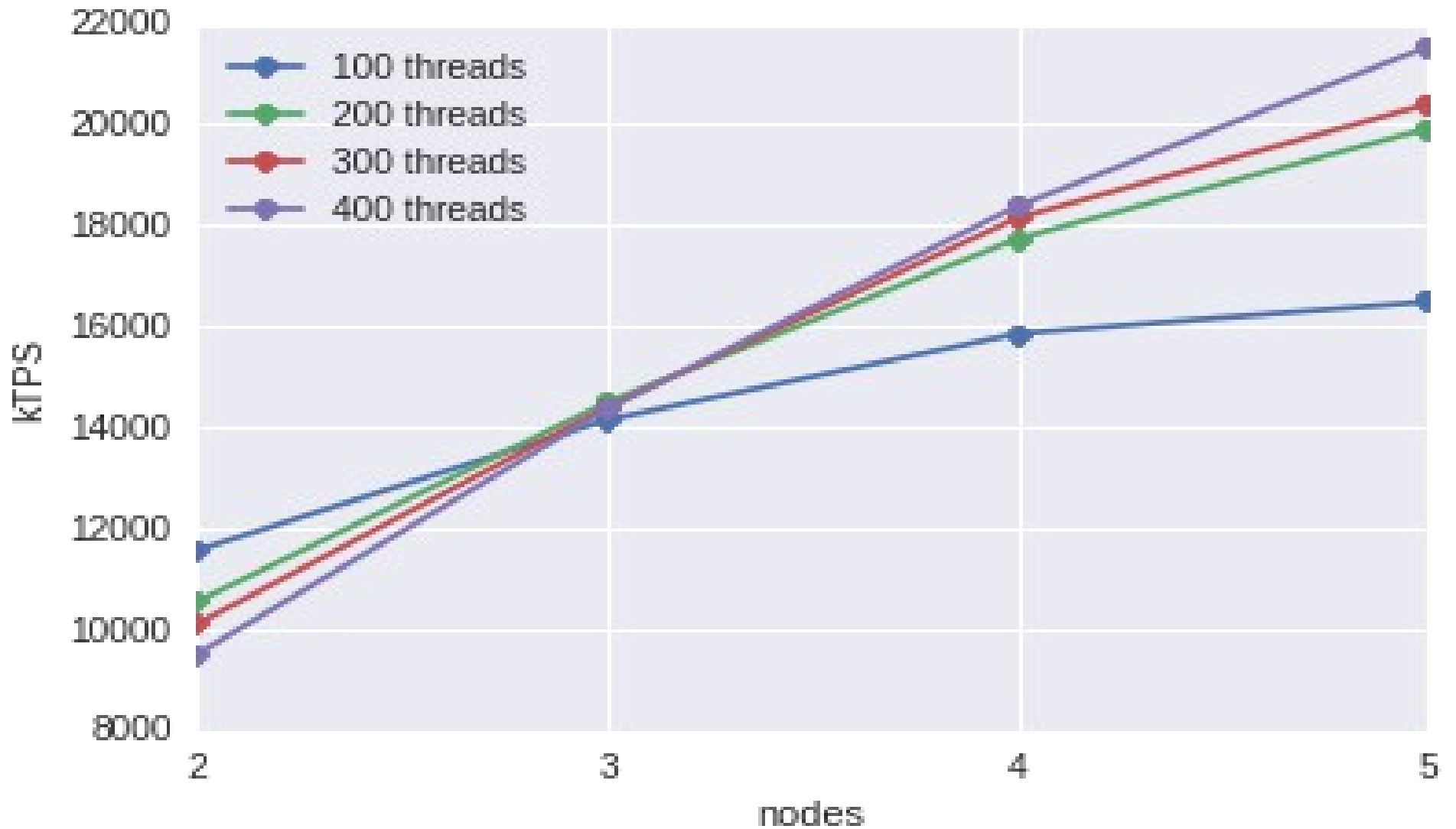

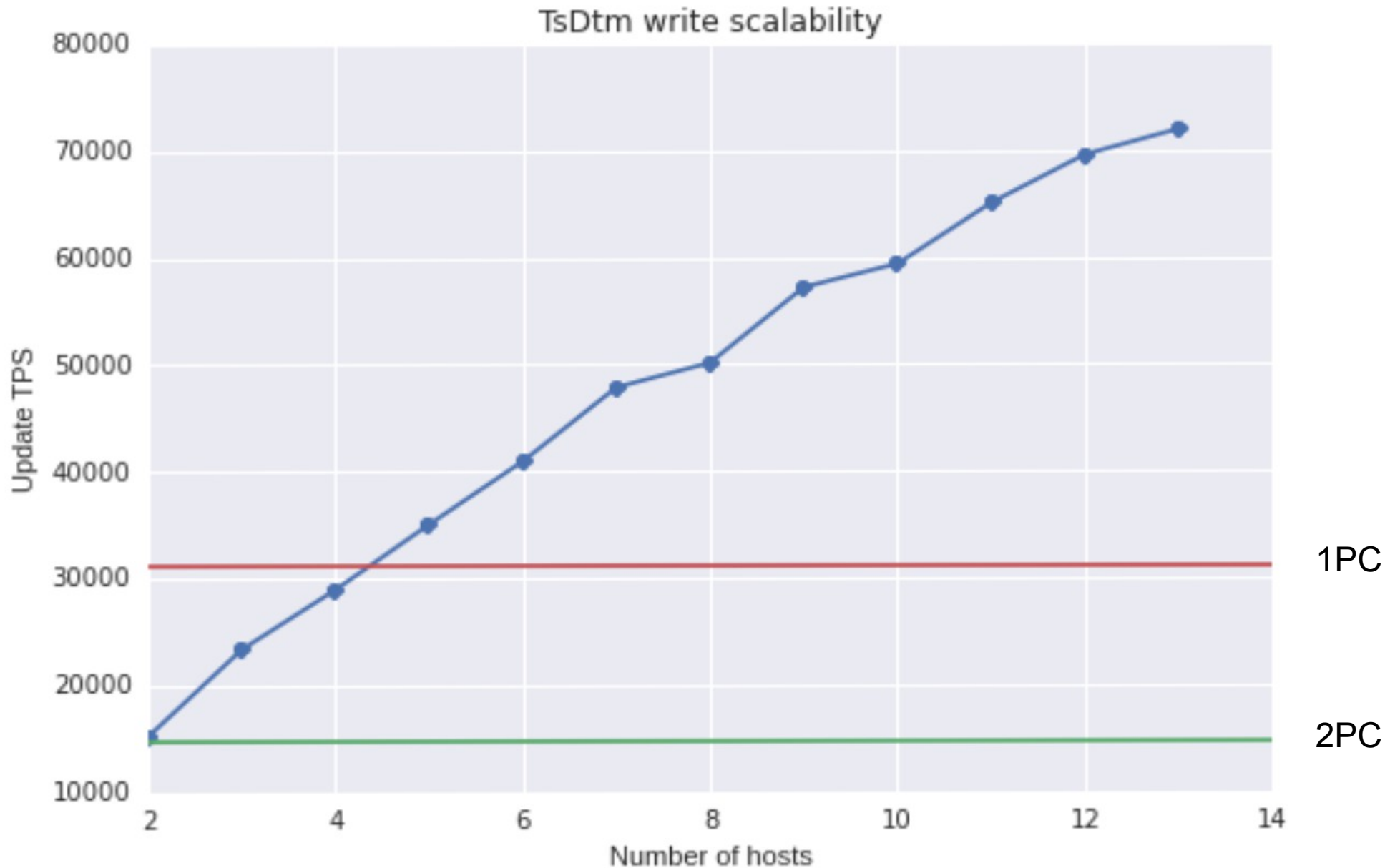Dtm write scalability (amazon c3.2xlarge)
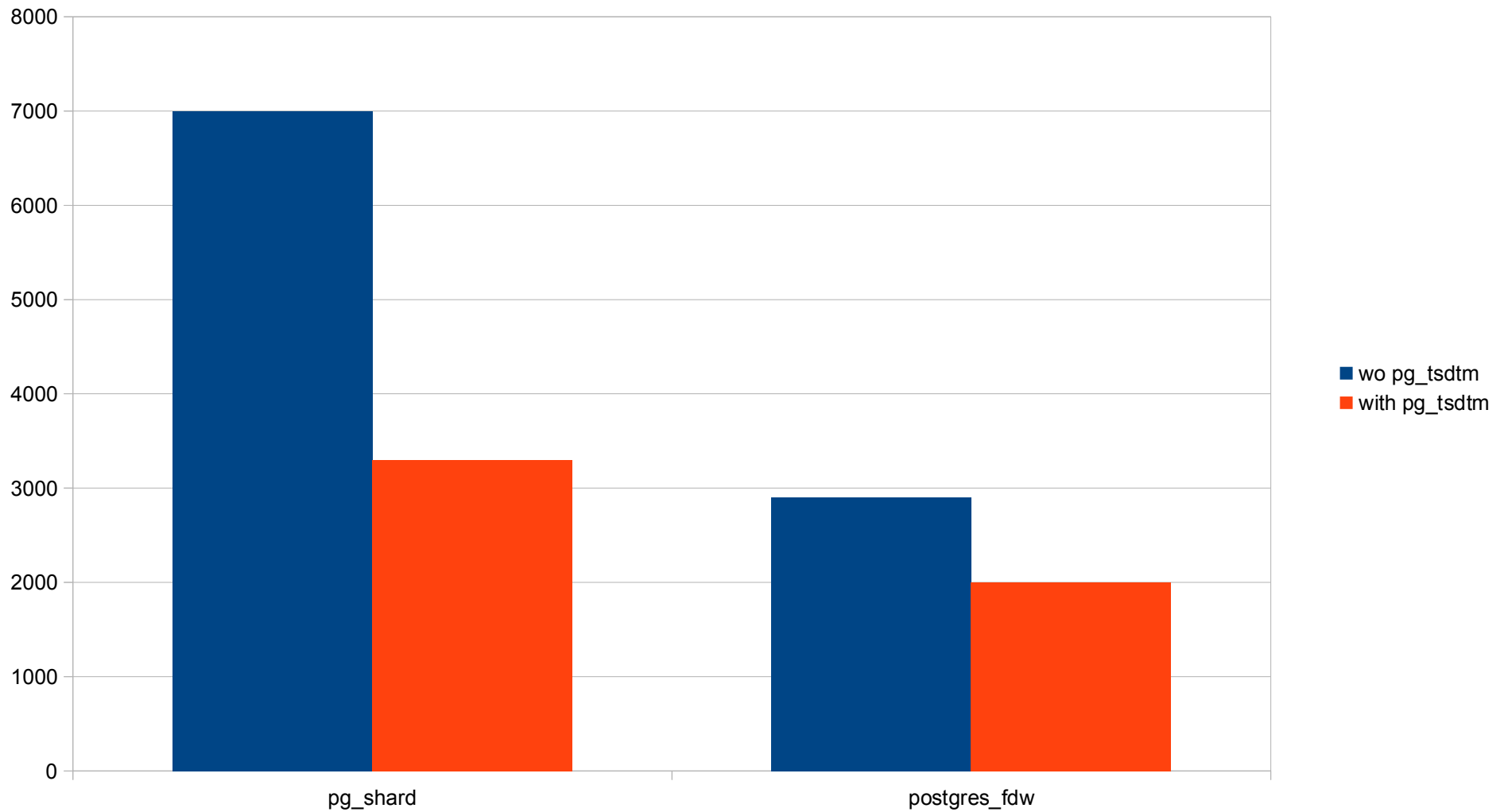
# tsDTM scalability
## (SAI cluster)

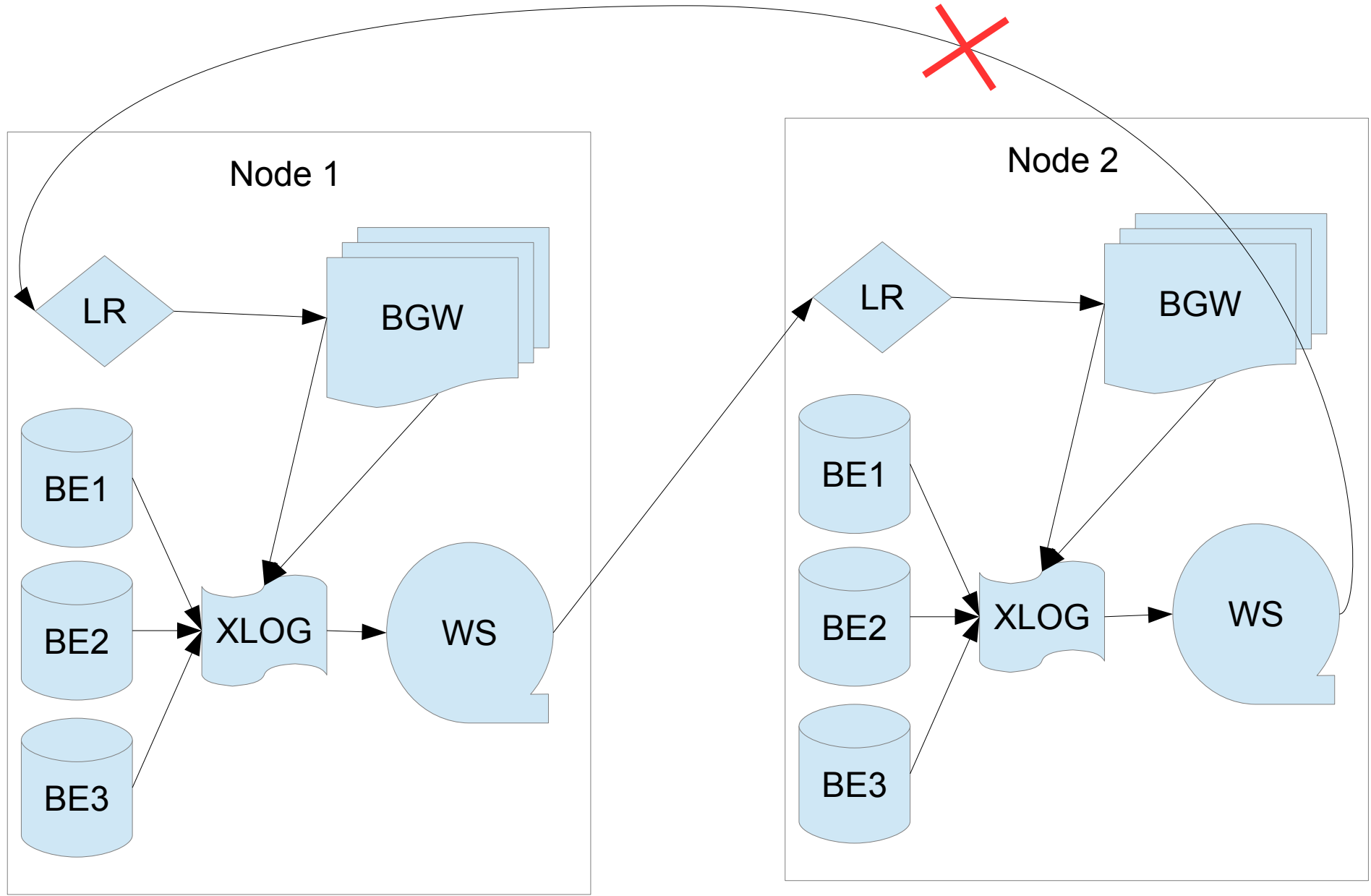# tsDTM scalability
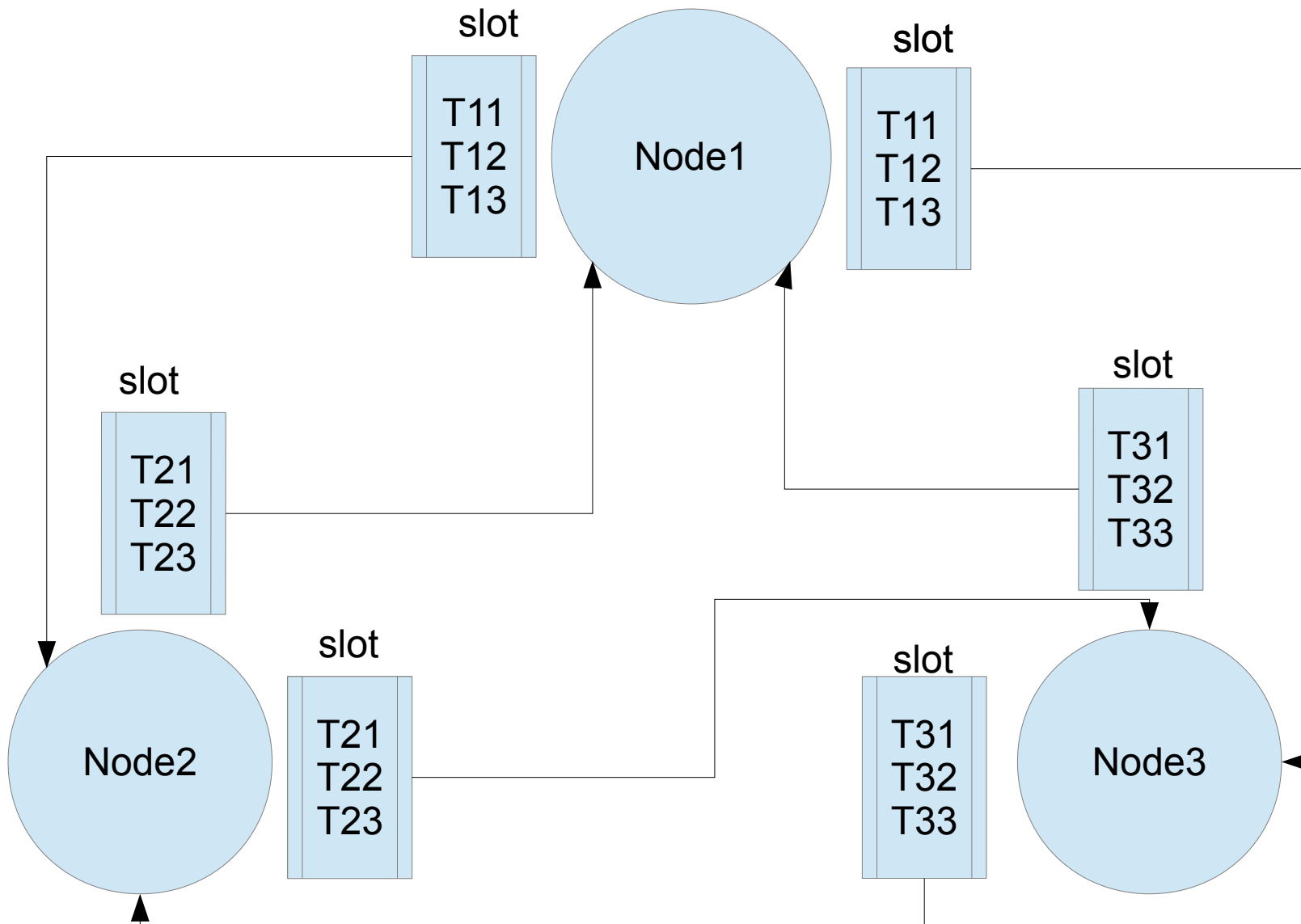
## (SAI cluster)

# tsDTM write scalability

# Multimaster formula

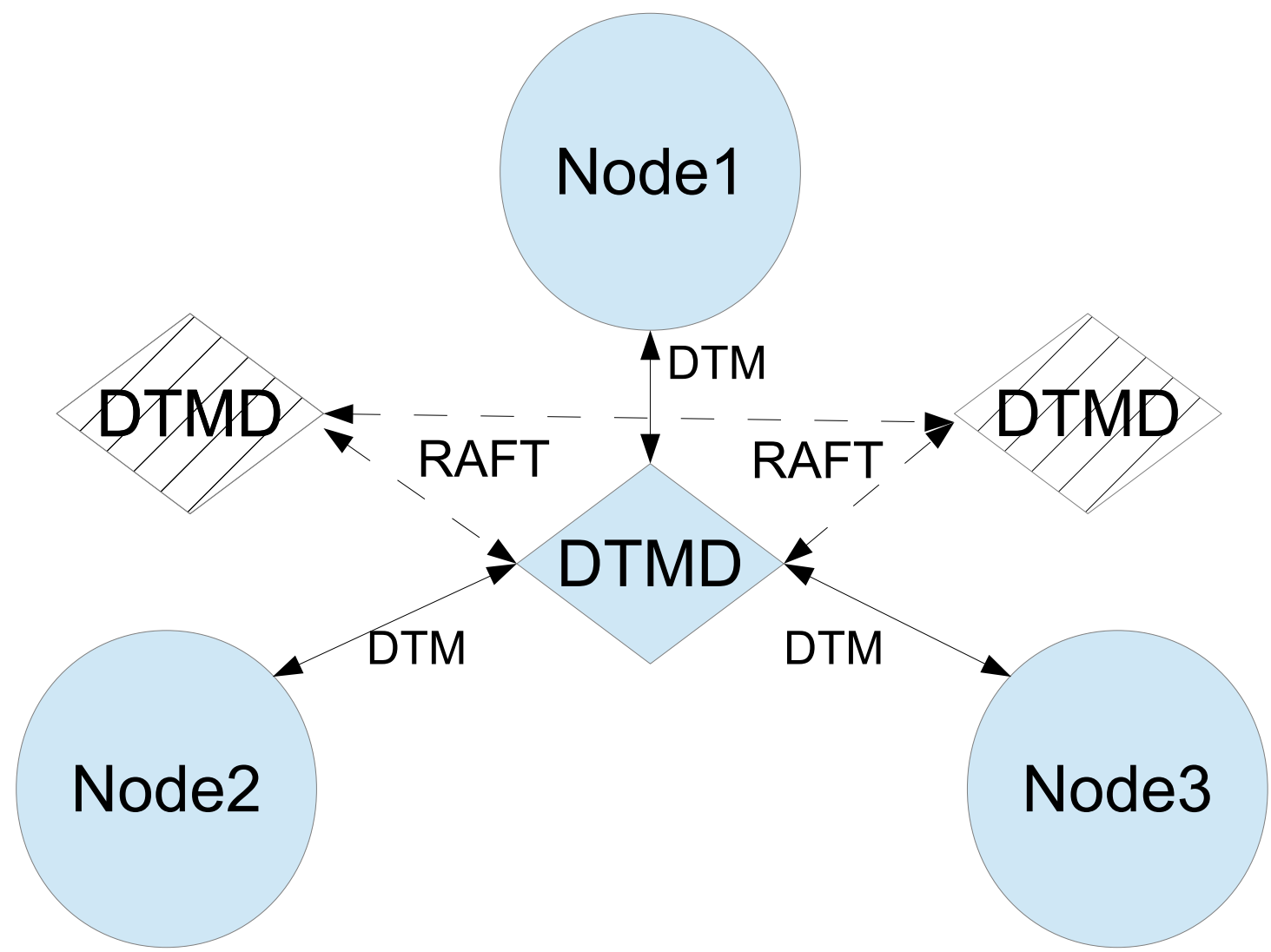$$MM = BDR + DTM$$

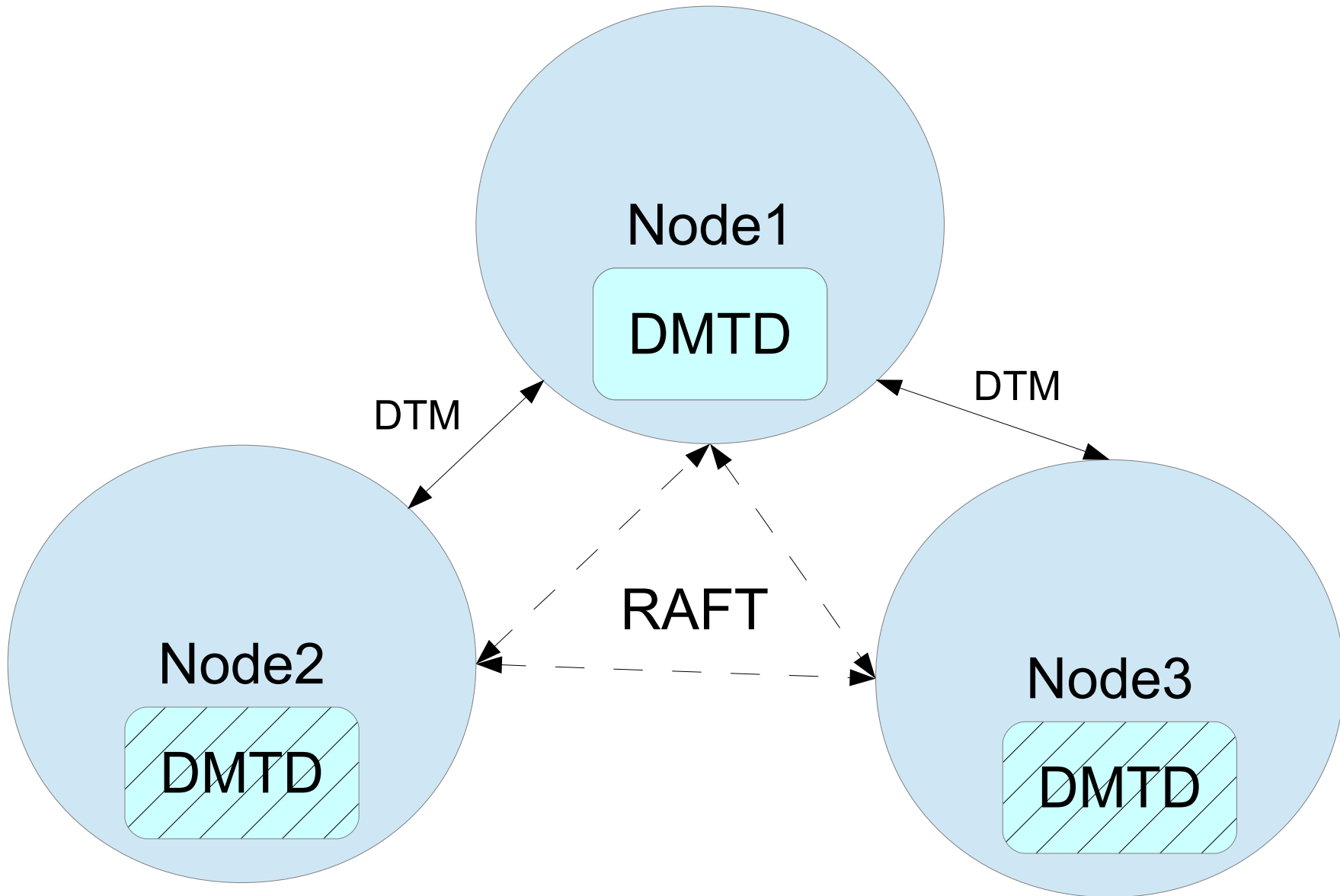# Multimaster based on logical replication

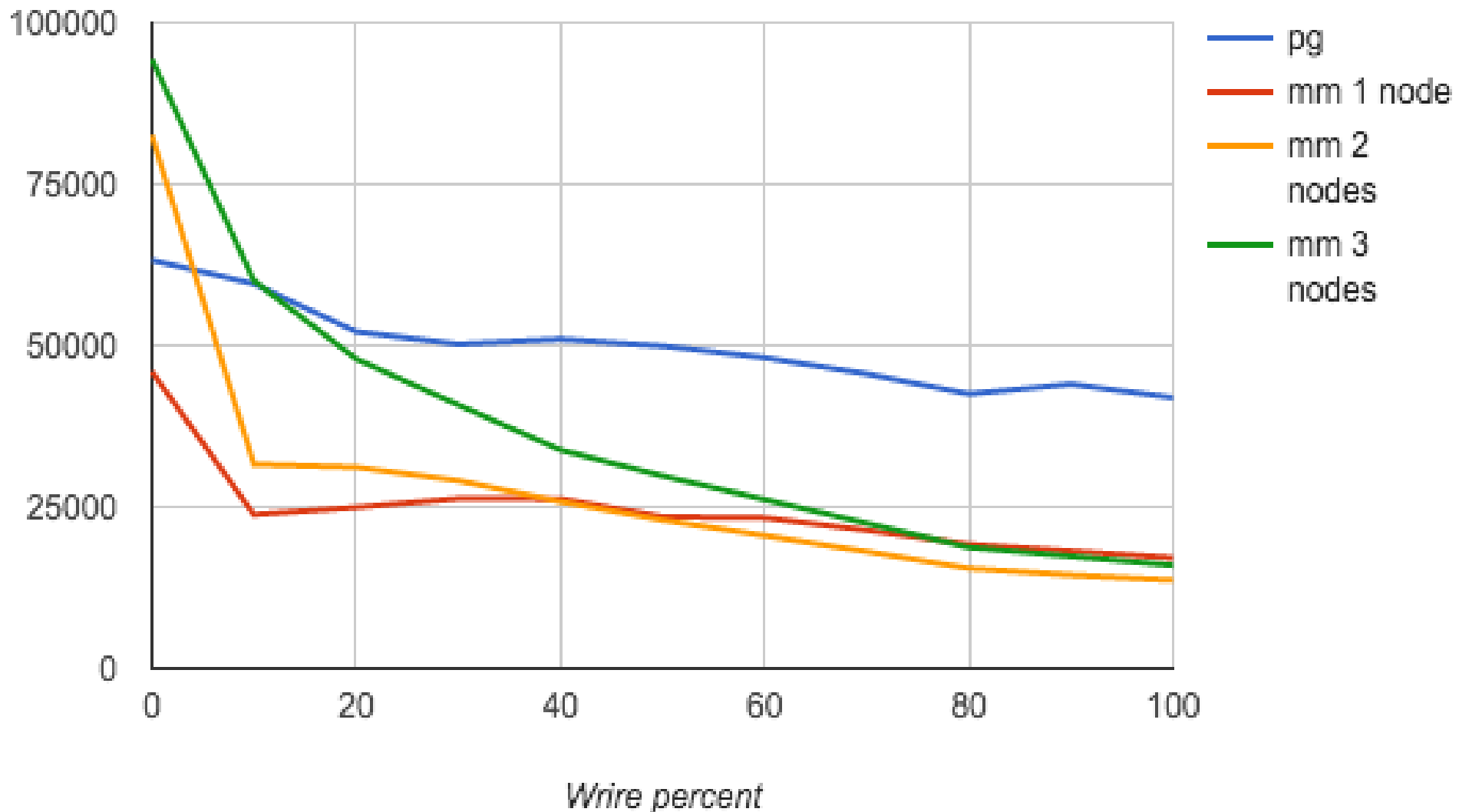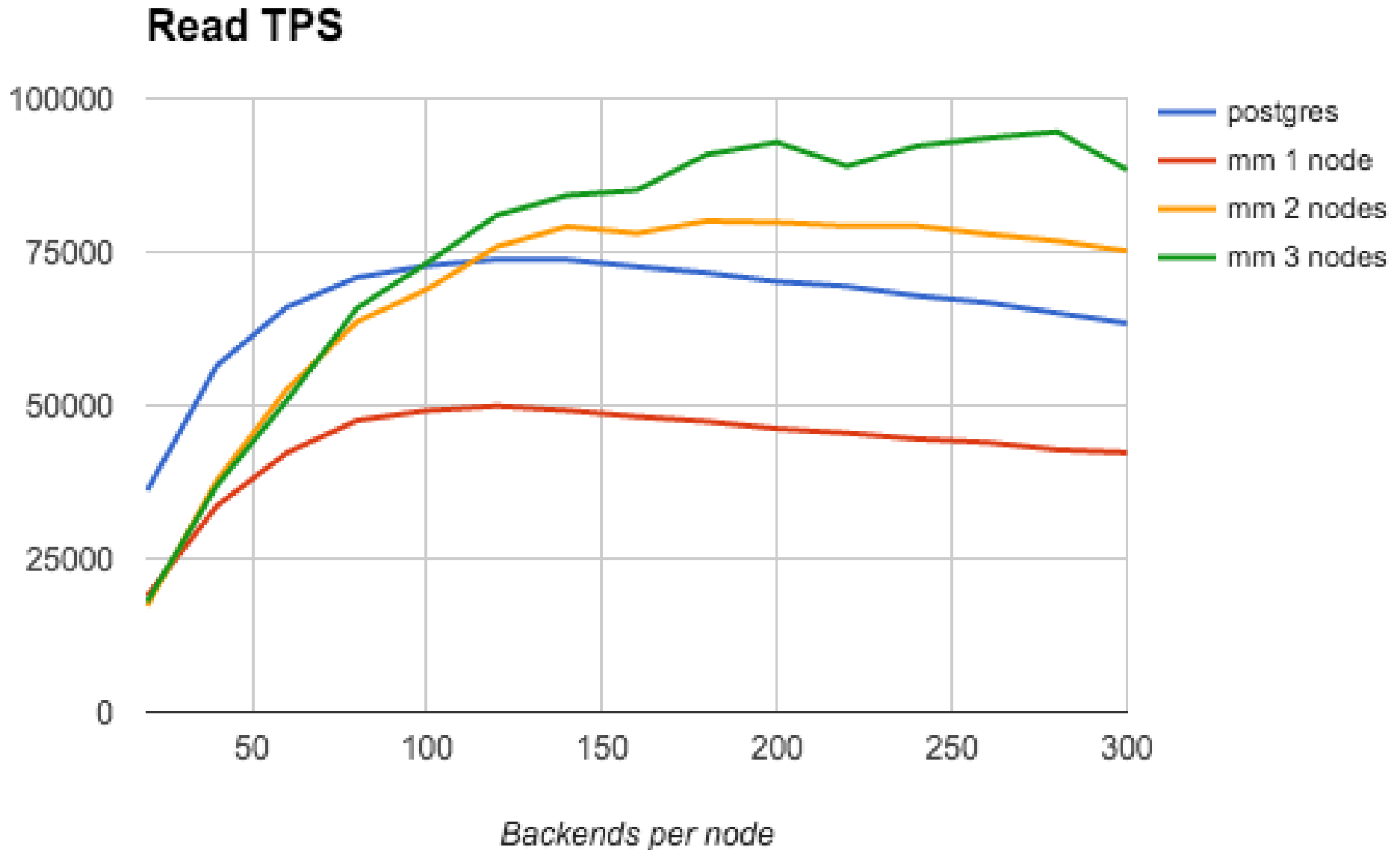# Logical repliction slots

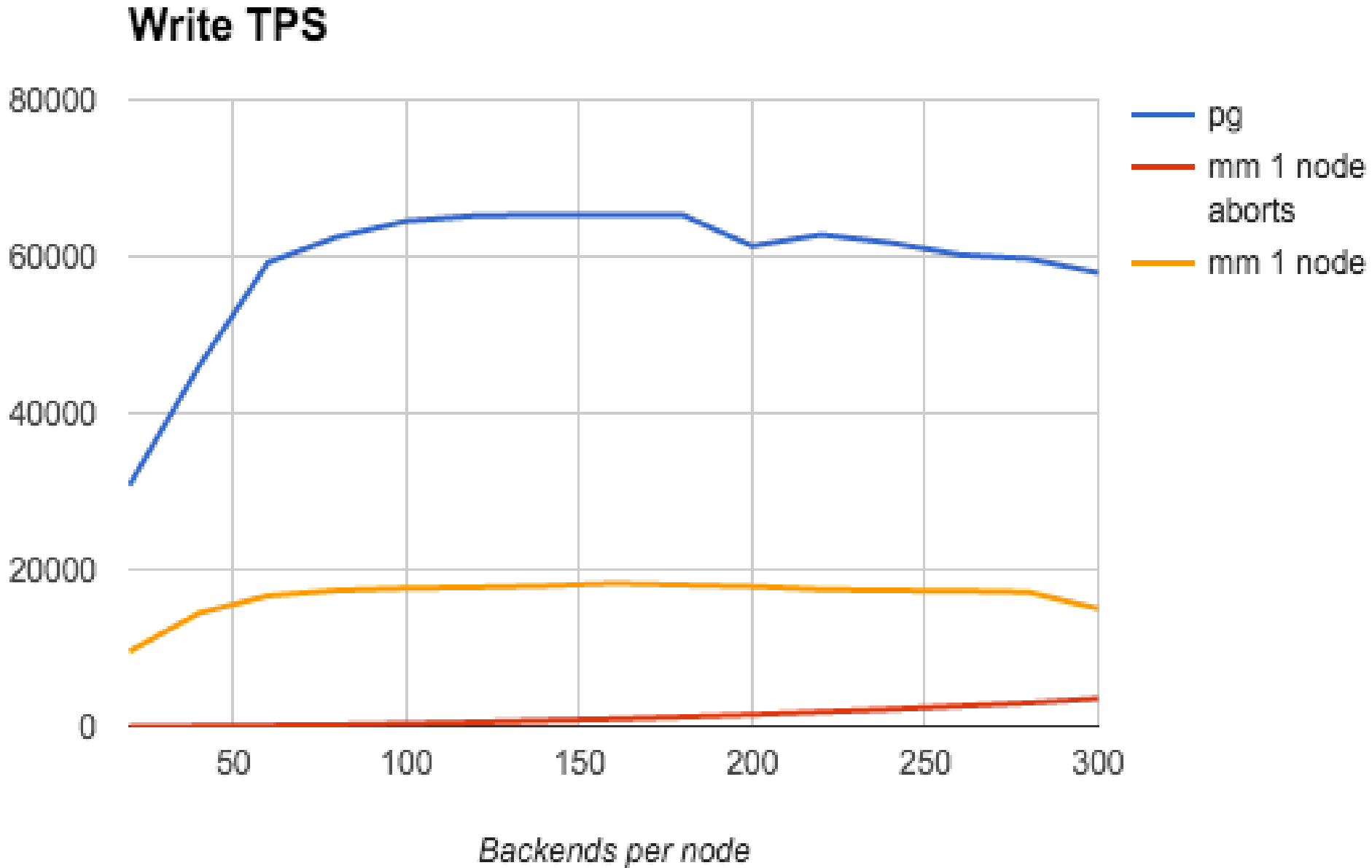# Current multimaster topology

HA topology

# Multimaster performance



TPS with different select/update ration

# Multimaster read scaling



**Read TPS**

Legend:
- postgres (blue)
- mm 1 node (red)
- mm 2 nodes (orange)
- mm 3 nodes (green)

Y-axis: 0, 25000, 50000, 75000, 100000

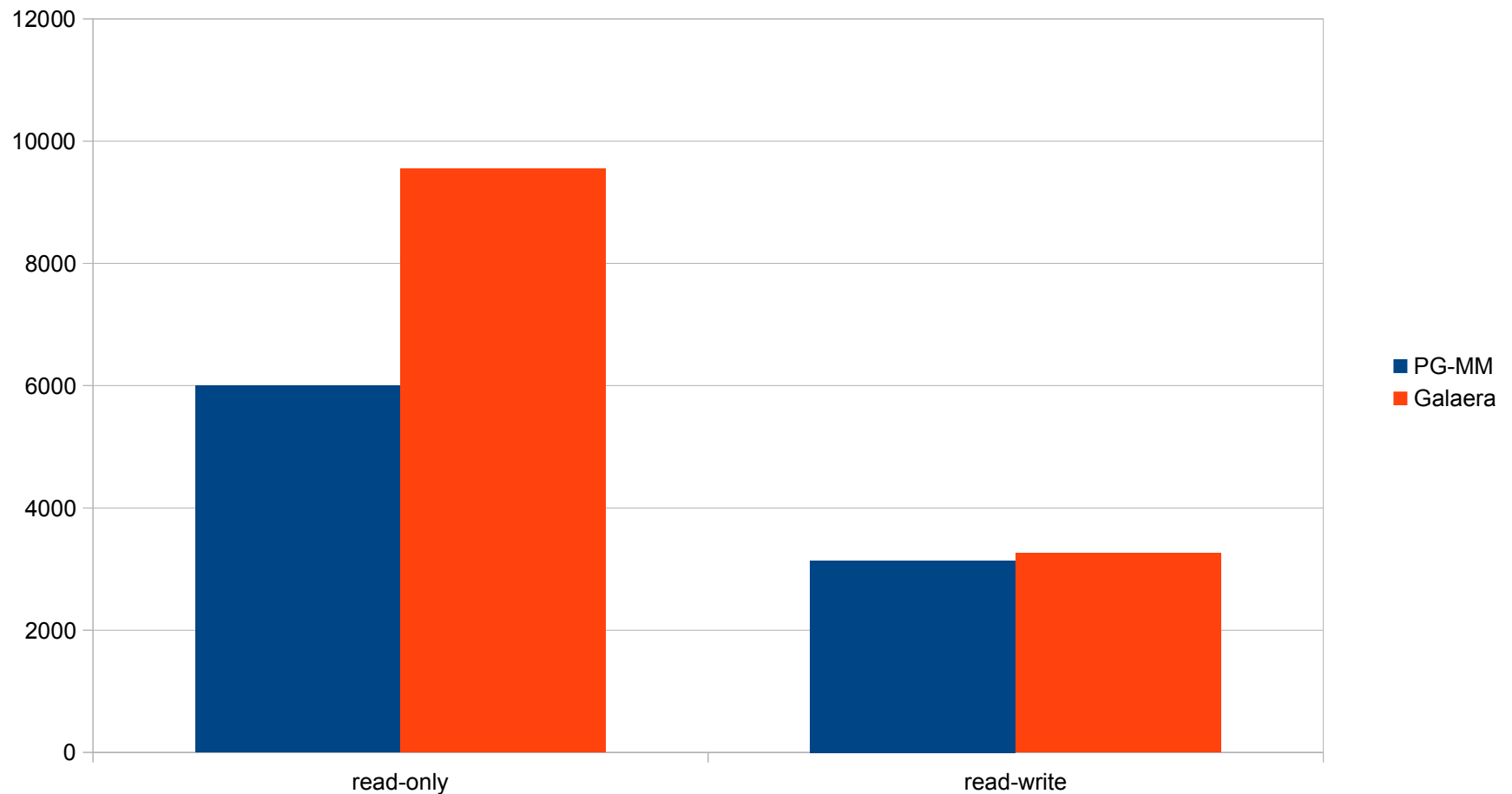X-axis (Backends per node): 50, 100, 150, 200, 250, 300
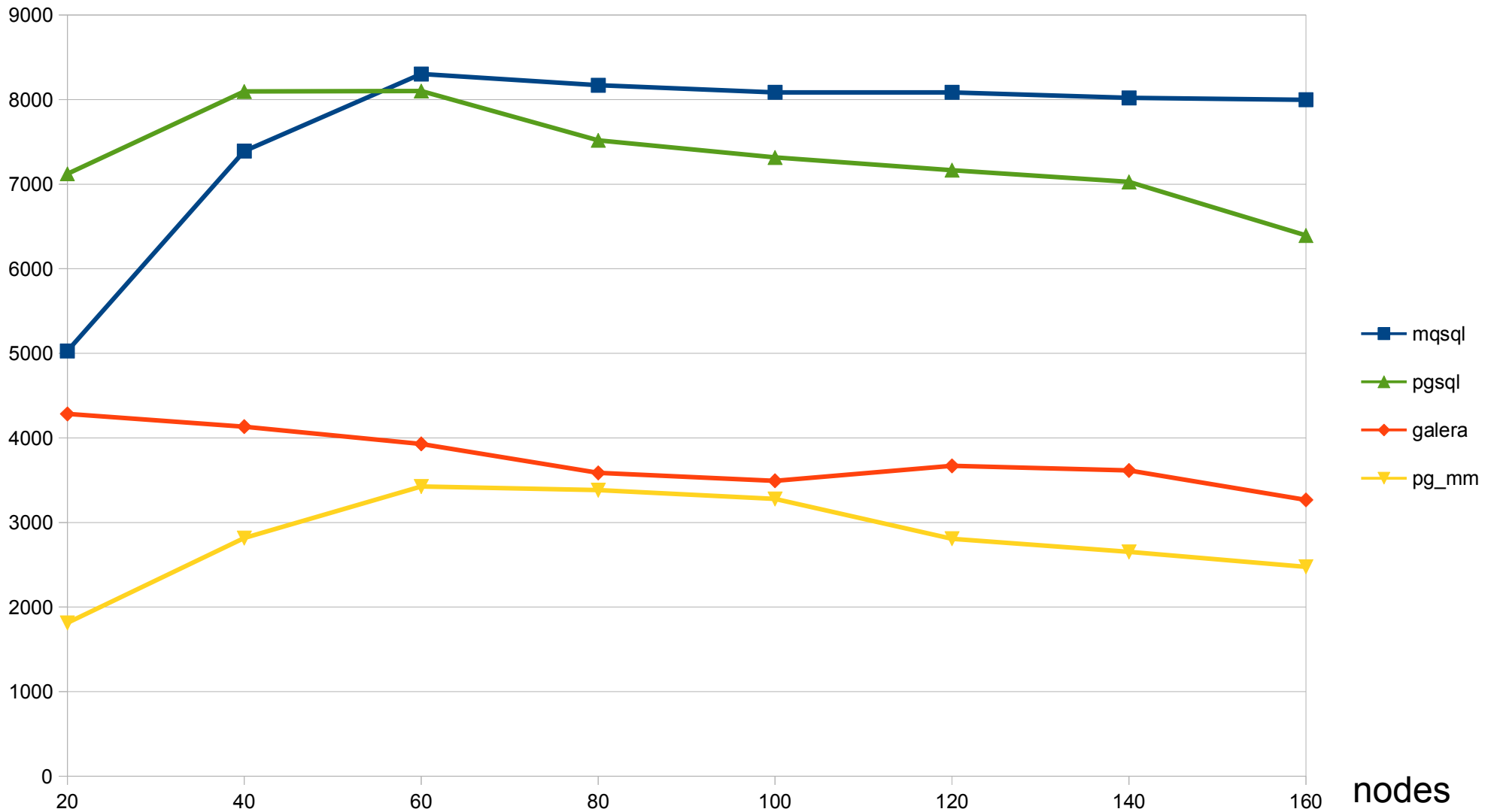
# Multimaster writer "scaling"

# Galera vs. Postgres Multimaster
## (sysbench OLAP-complex, 10M records, 3 nodes)

# Galera vs. PostgreSQL multimaster
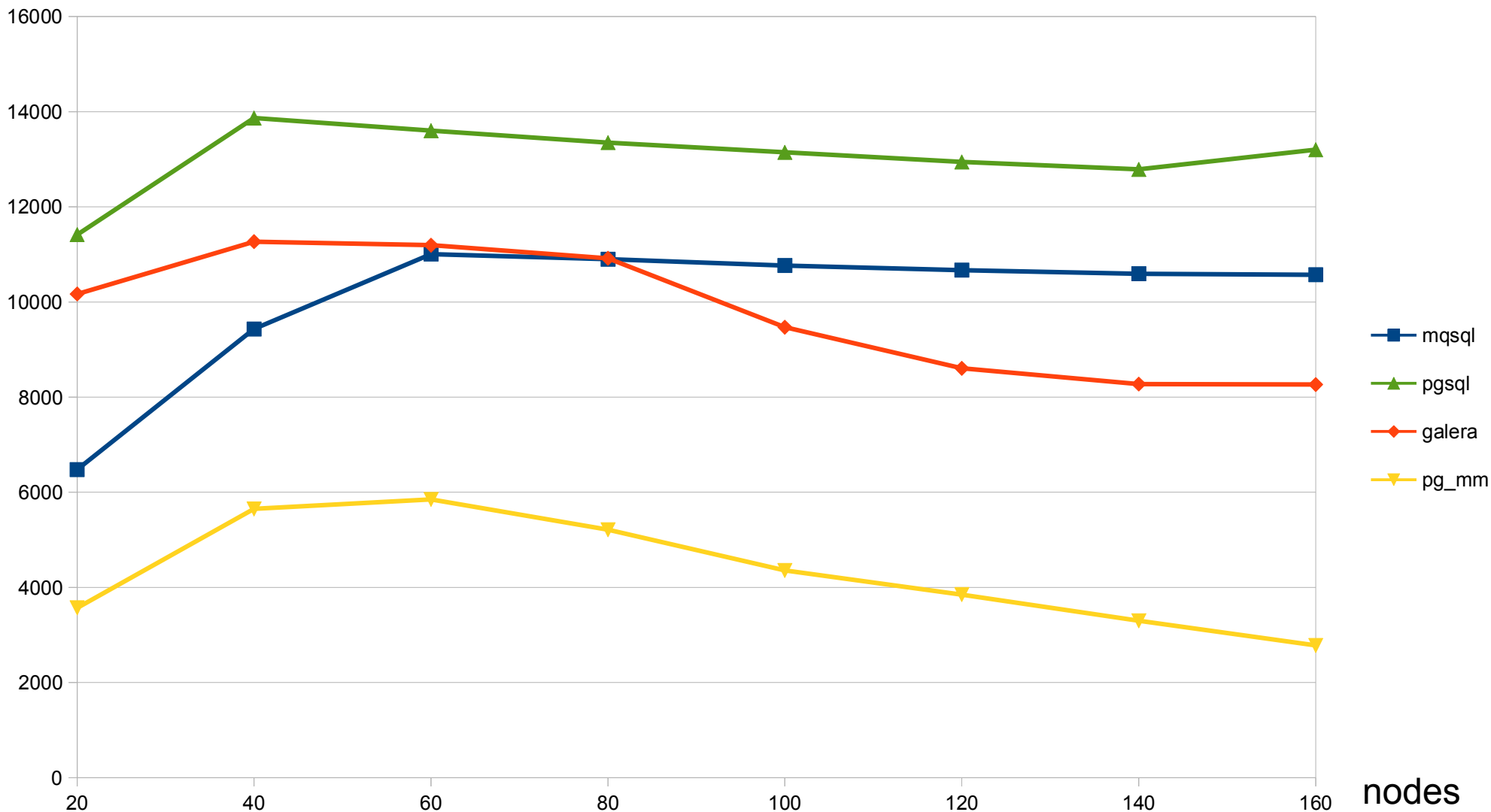## (sysbench, read-write complex OLAP)

TPS

nodes

- mqsql
- pgsql
- galera
- pg_mm

Galera vs PostgreSQL multimaster
(sysbench, read-only complex OLAP)

# Current multimaster limitations

- Table should have primary key.

- DDL is not handled by logical replication and requires separate replication channel which currently is not using 2PC.

- Subtransactions are not supported (limitation of DTM).

- Explicit locks are not distributed.

- Number of concurrent transactions is limited by number of BG workers

# Questions?

https://github.com/postgrespro/postgres_cluster.git