

YAHOO!



Omid: A Transactional Framework for HBase

Francisco Perez-Sorrosal
Ohad Shacham



Outline

- Background
- Basic Concepts
- Use cases
- Architecture
- Transaction Management
- High Availability
- Performance
- Summary

Background

- New Big data apps → new requirements:
 - Low-latency
 - Incremental data processing
 - e.g. Percolator
- Multiple clients updating same data concurrently
 - **Problem:** Conflicts/Inconsistencies may arise
 - **Solution:** **Transactional Access to Data**

Background

- Transaction → Abstract UoW to manage data with certain guarantees
 - **ACID**
 - Relational databases
- Big data → NoSQL datastores → Transactions in NoSQL
 - Hard to Scale
 - Data partition
 - Data replication
 - Relaxed Guarantees:
 - e.g. Atomicity, Consistency

Omid is a...

- **Flexible**
- **Reliable**
- **High Performant**
- **Scalable**

...**OLTP framework** that allows **BigData** apps to execute **ACID transactions** on top of **HBase**

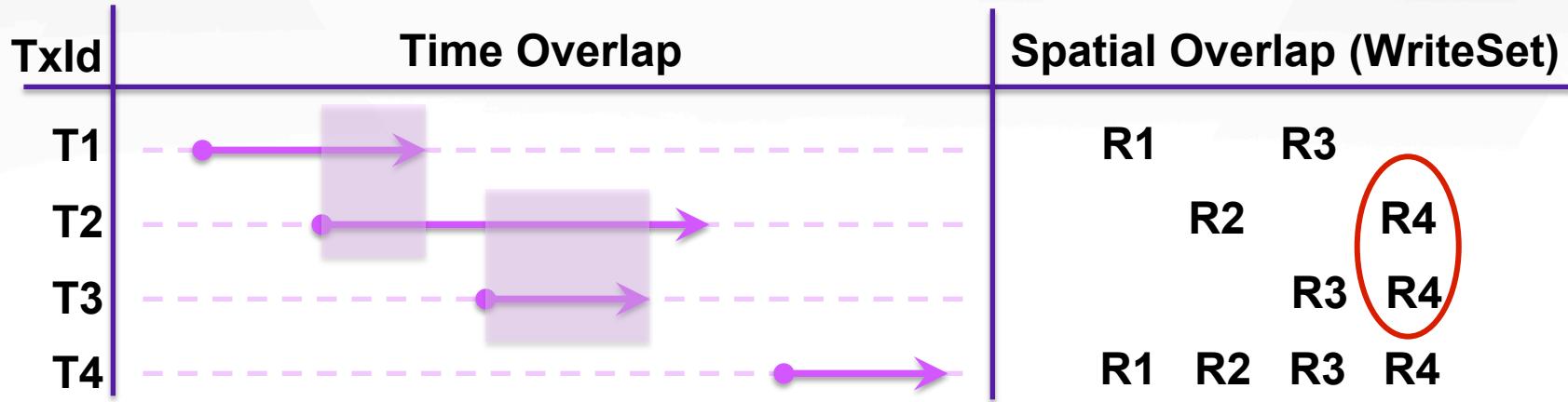


Why use Omid?

- **Simplifies development** of apps requiring consistency
 - Multi-row/multi-table transactions on HBase
 - Simple & well-known interface
- **Good performance & reliability**
- **Lock-free**
- **Snapshot Isolation**
- **HBase is a blackbox**
 - No HBase code modification
 - No changes on table schemas
- **Used successfully at Yahoo**

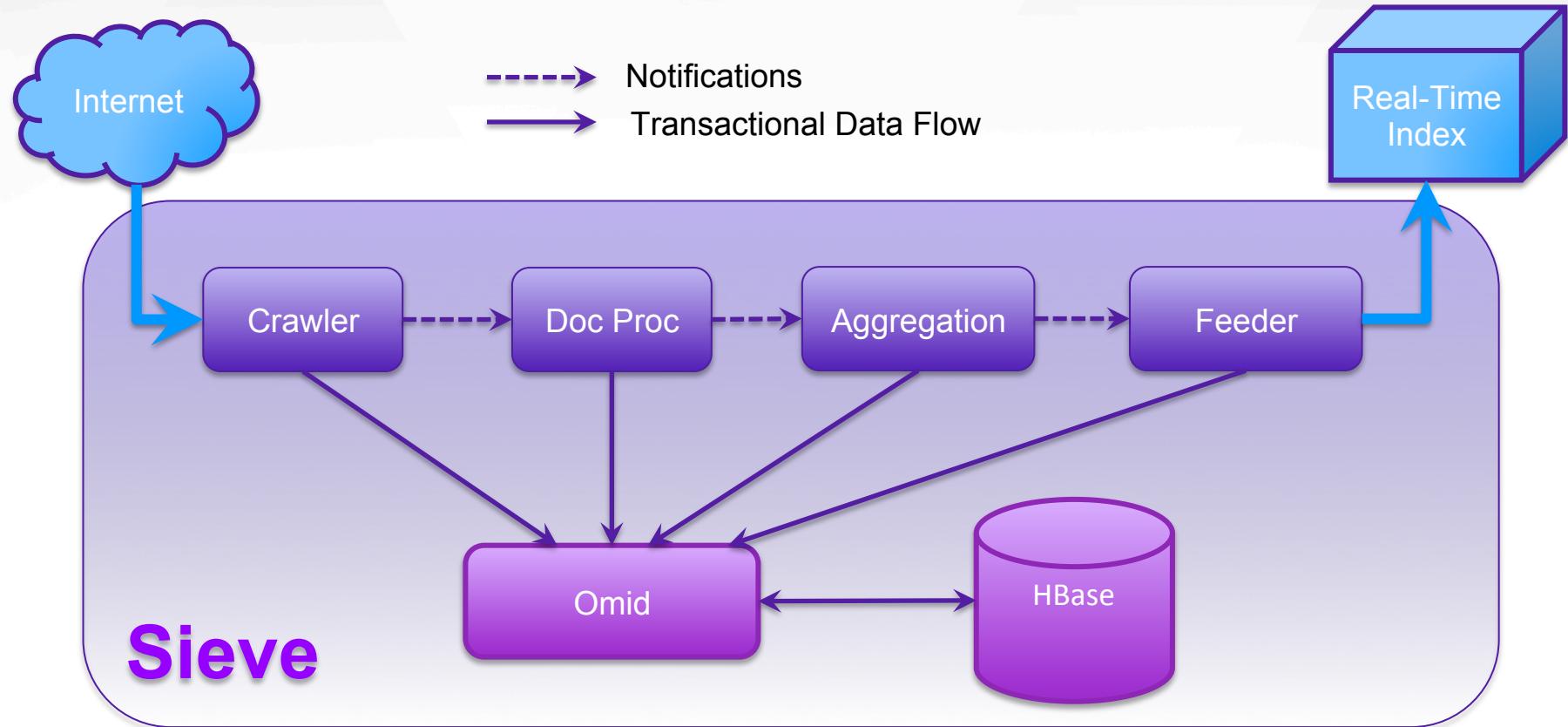


Snapshot Isolation

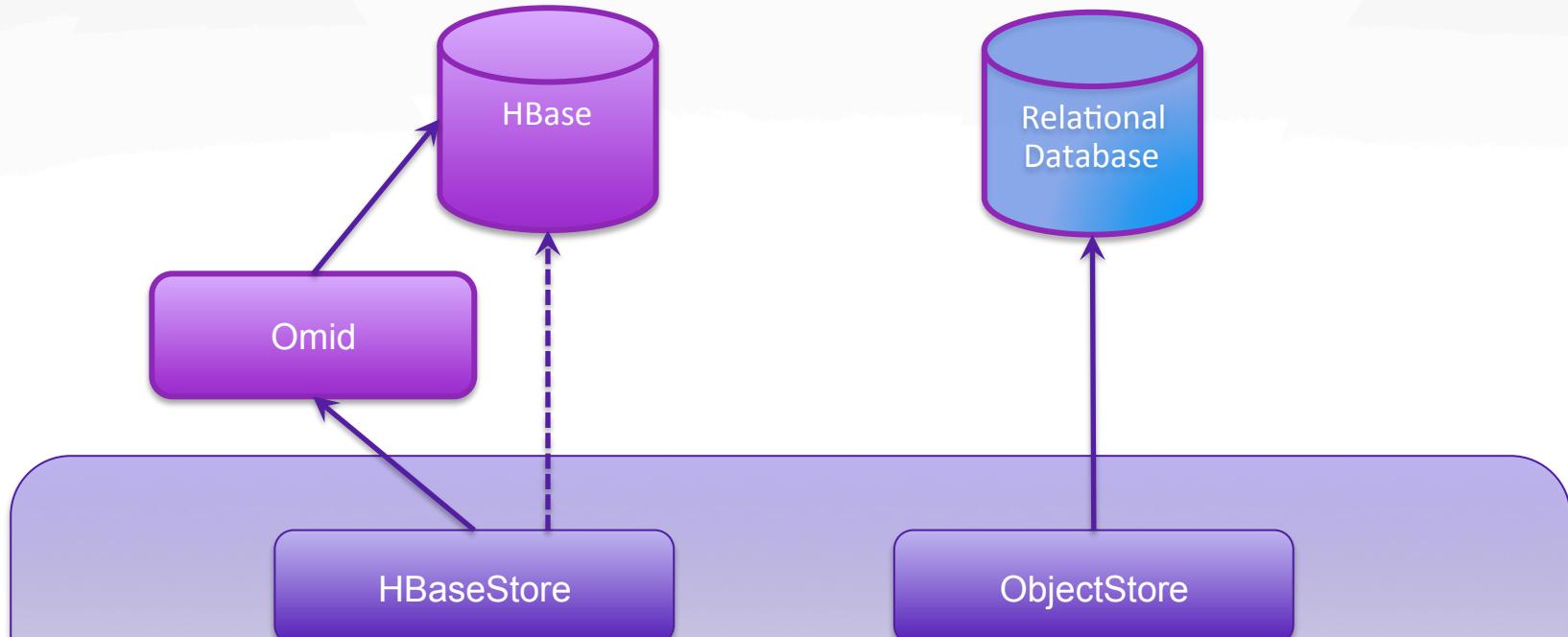


- Transaction T2 overlaps in time with T1 & T3, but spatially:
 - $T1 \cap T2 = \emptyset$
 - $T2 \cap T3 = \{ R4 \}$ Transactions T2 and T3 conflict
- Transaction T4 does not have conflicts

Use Cases: Sieve @ Yahoo

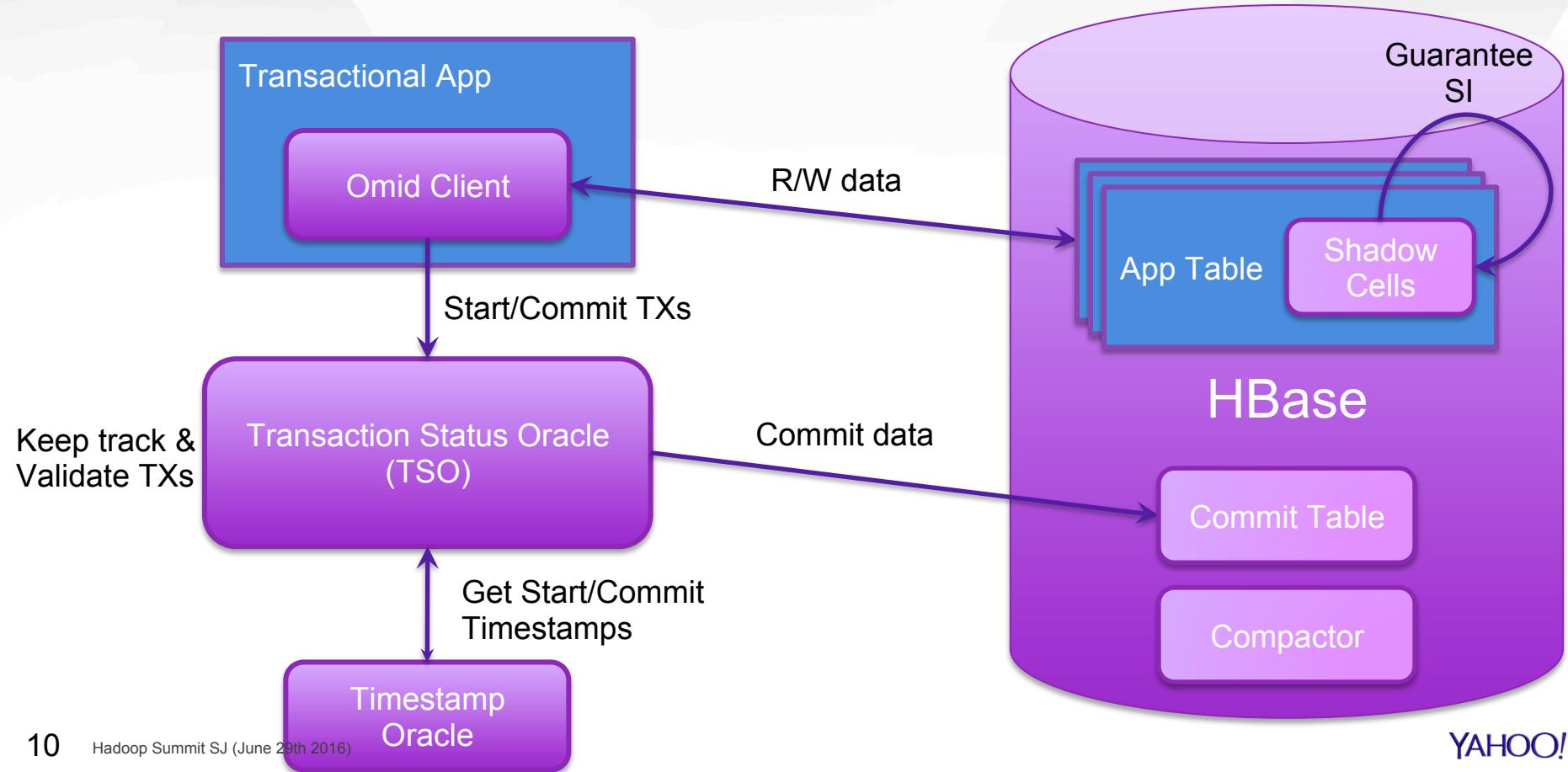


Use Cases:



Hive Metastore Thrift Server

Architectural Components



Client APIs

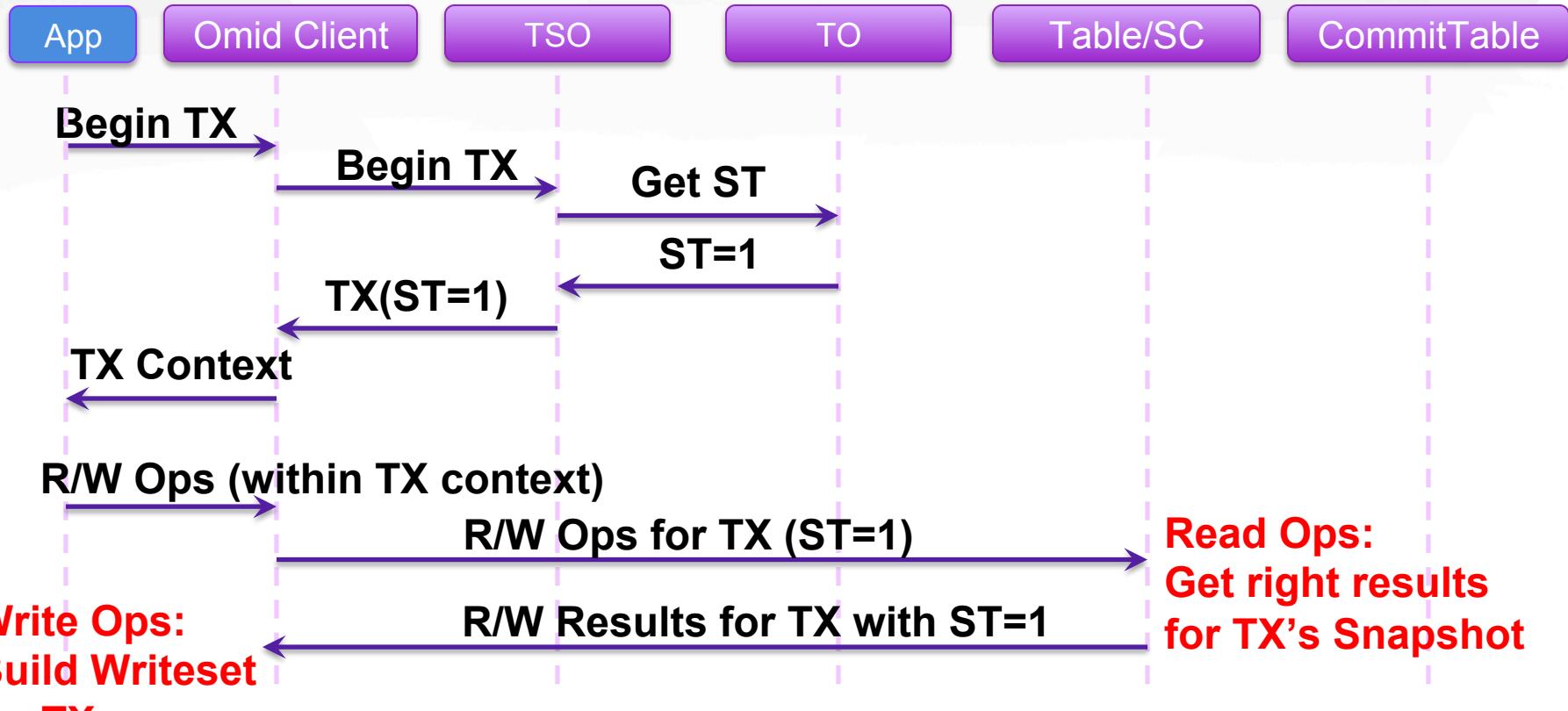
- Transaction Manager → Create Transactional contexts

```
Transaction begin();  
void commit(Transaction tx);  
void rollback(Transaction tx);
```

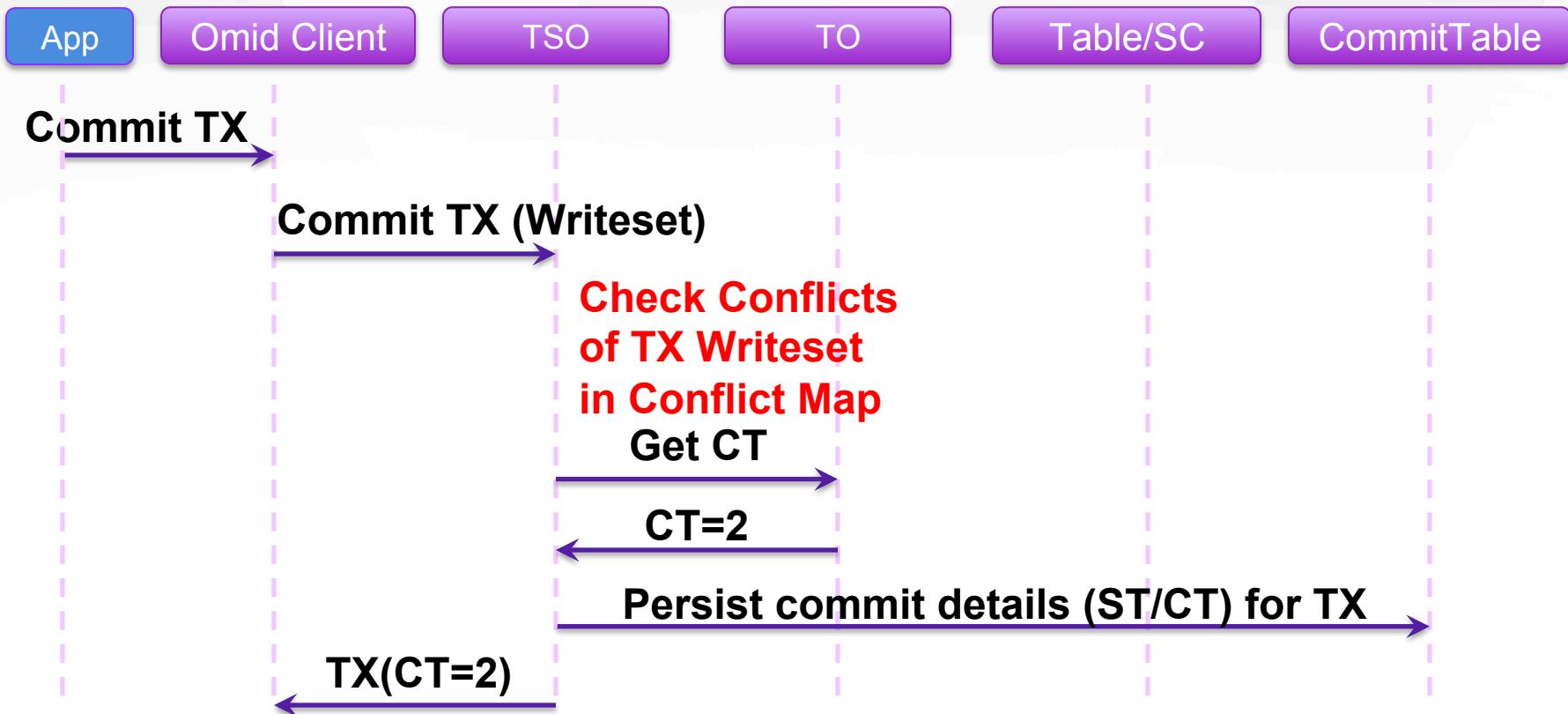
- Transactional Tables (TTable) → Data access

```
Result get(Transaction tx, Get g);  
void put(Transaction tx, Put p);  
ResultScanner getScanner(Transaction tx, Scan s);
```

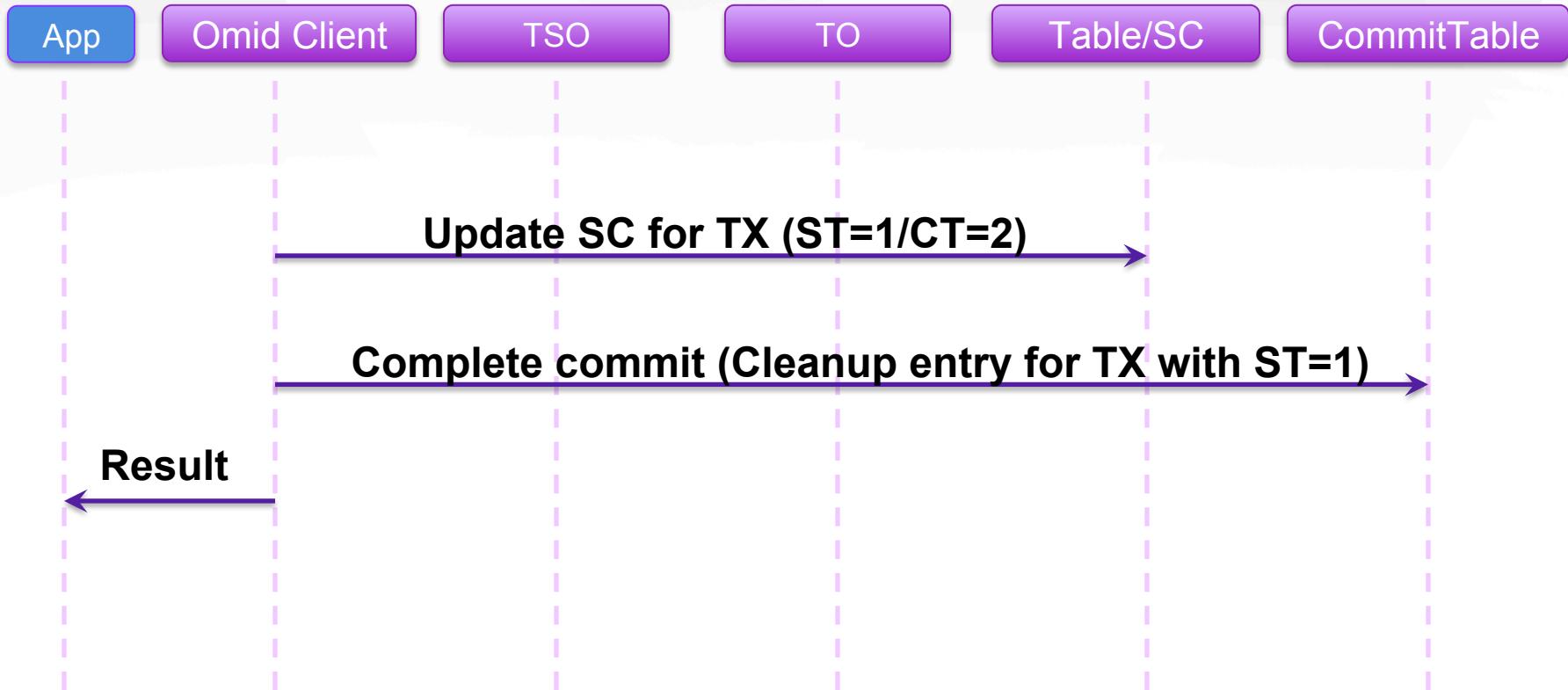
TX Management (Begin TX phase)



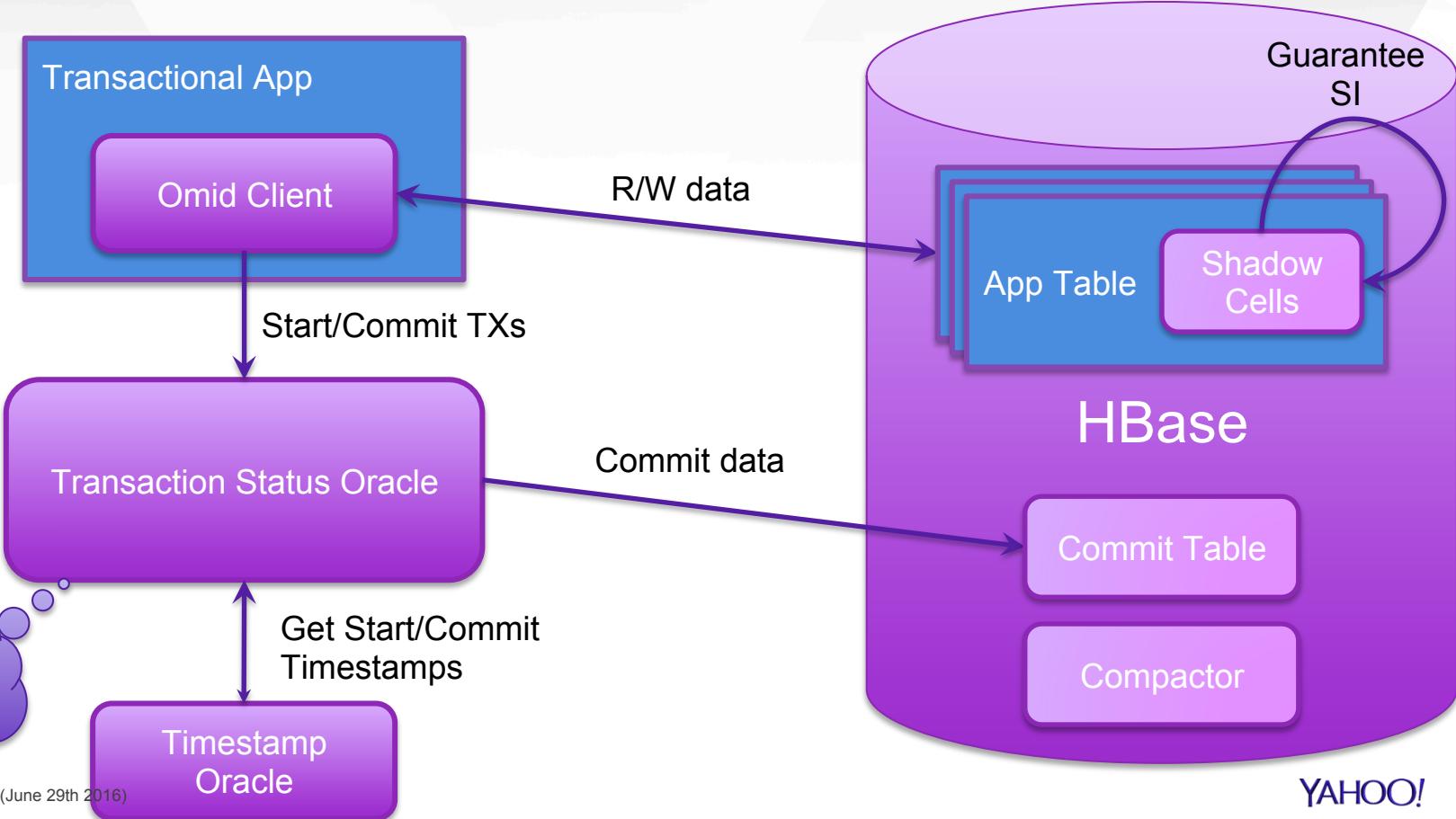
TX Management (Commit TX Phase)



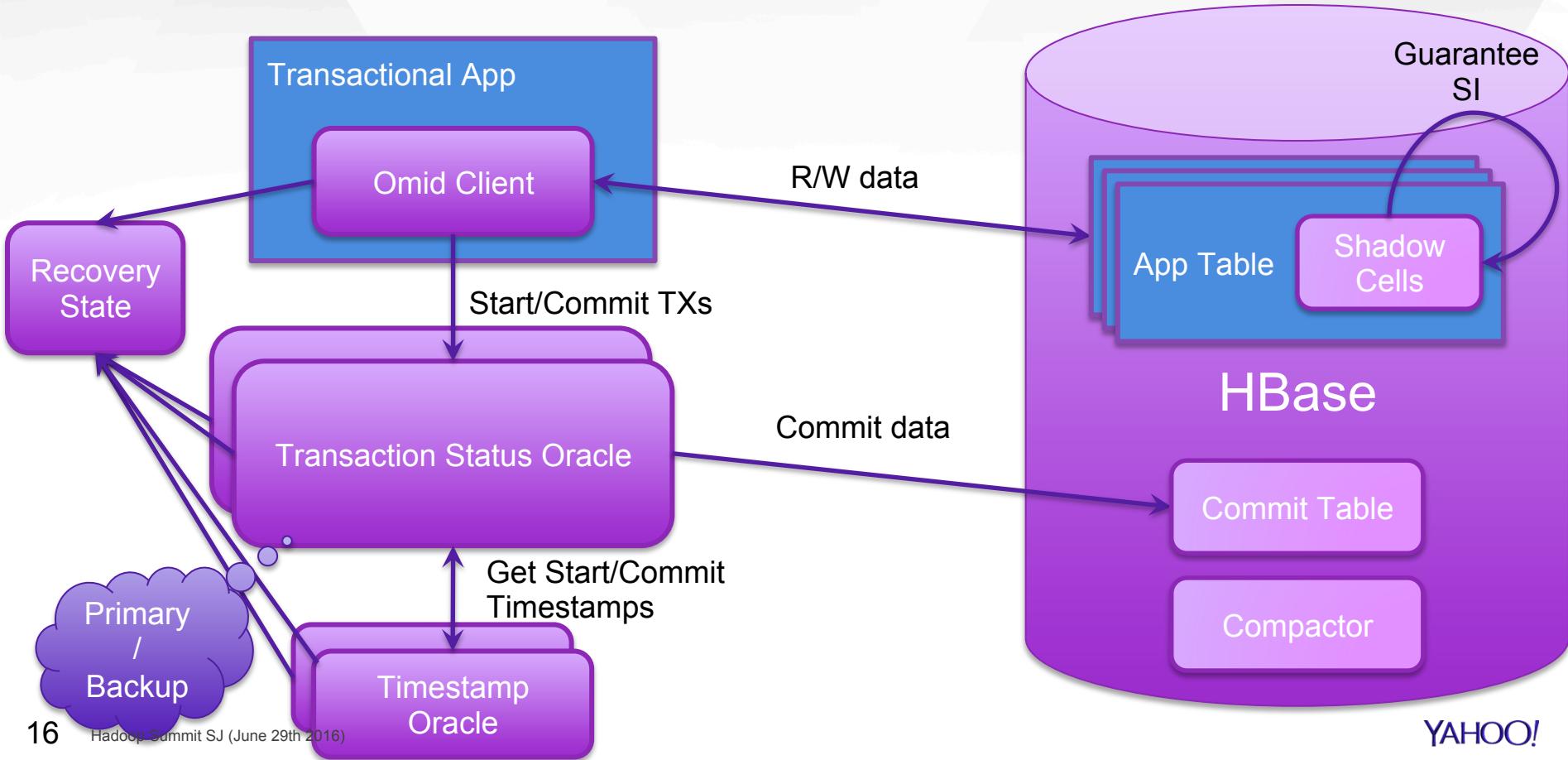
TX Management (Complete TX Phase)



High Availability



High Availability



High Availability – Failing Scenario

App

Omid Client

TSO P

TSO B

TO

Table/SC

CommitTable

Begin TX

Begin TX

Get ST

ST=1

TX(ST=1)

TX 1

TX 1 Write(k1, v1)

Write(k1, v1) (ST=1)



High Availability – Failing Scenario

App

Omid Client

TSO P

TSO B

TO

Table/SC

CommitTable

Write(k2, v2)

Write(k2, v2) (ST=1)

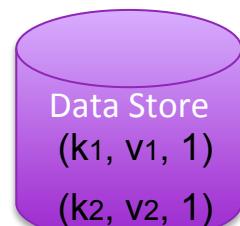
Commit TX 1

Commit TX 1{k1, k2}

Get CT

CT=2

Persist commit details for TX 1



High Availability – Failing Scenario

App

Omid Client

TSO B

TO

Table/SC

CommitTable

Begin TX

Begin TX

Get ST

ST=3

TX(ST=3)

TX 3

TX 3 Read(k1)

Read(k1) (ST=3)

(k1, v1, 1)

Data Store
(k1, v1, 1)
(k2, v2, 1)

Commit Table

High Availability – Failing Scenario

App

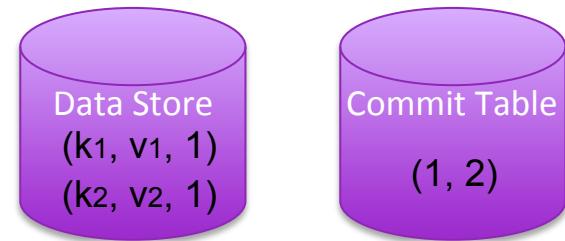
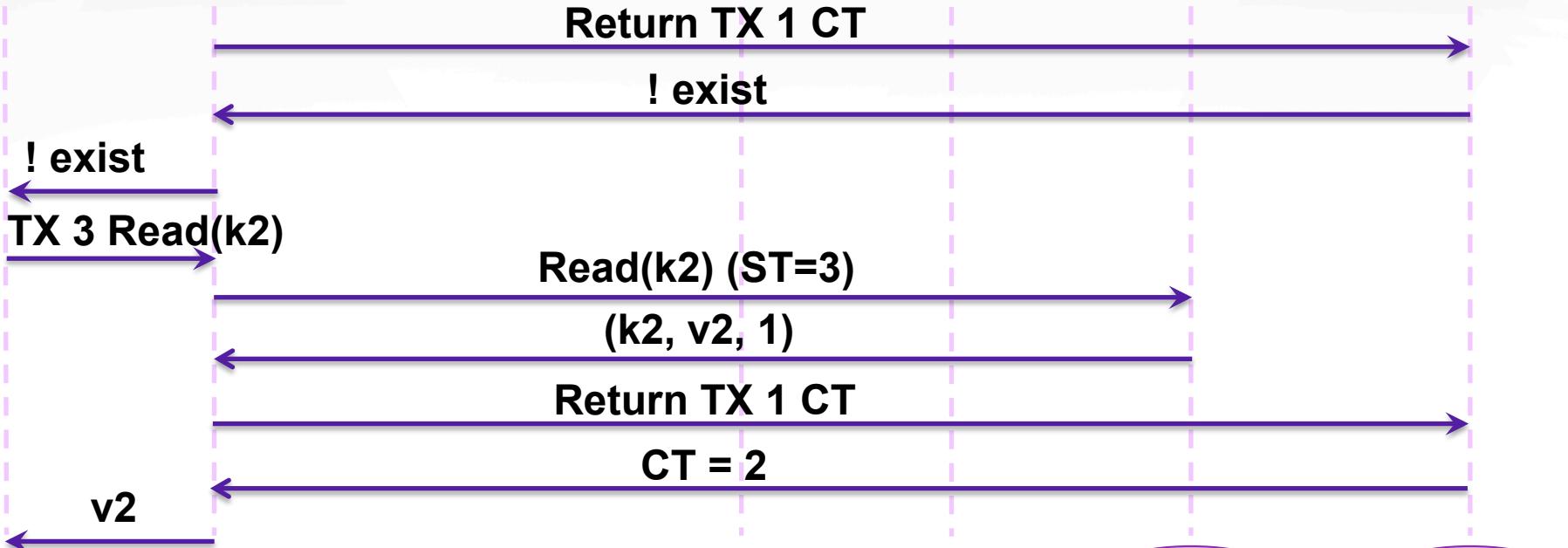
Omid Client

TSO B

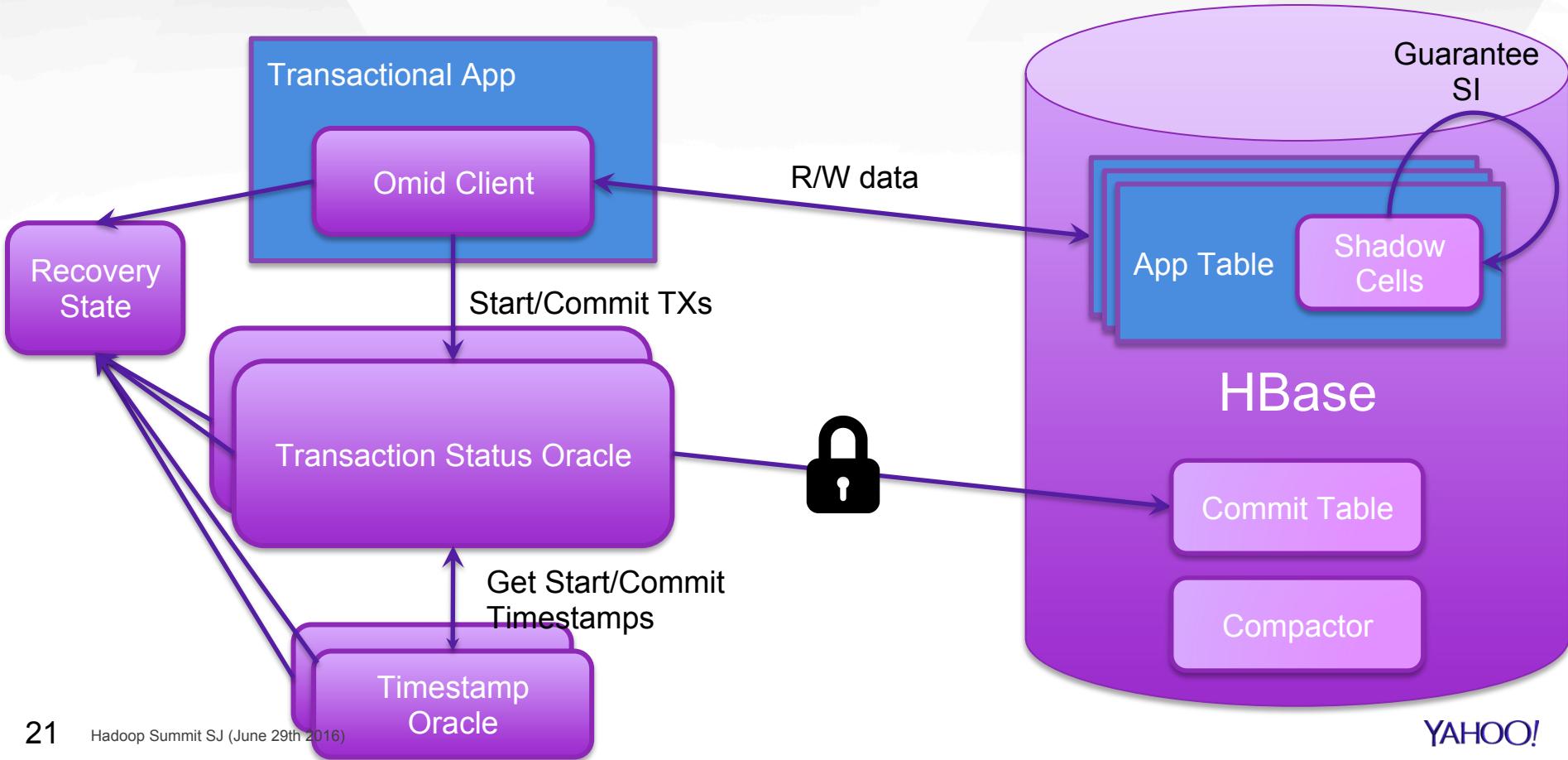
TO

Table/SC

CommitTable



High Availability



High Availability – Solution

App

Omid Client

TSO P

TSO B

TO

Table/SC

CommitTable

Begin TX

Begin TX

Get ST

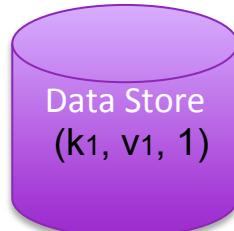
ST=1

TX(ST=1,E=1)

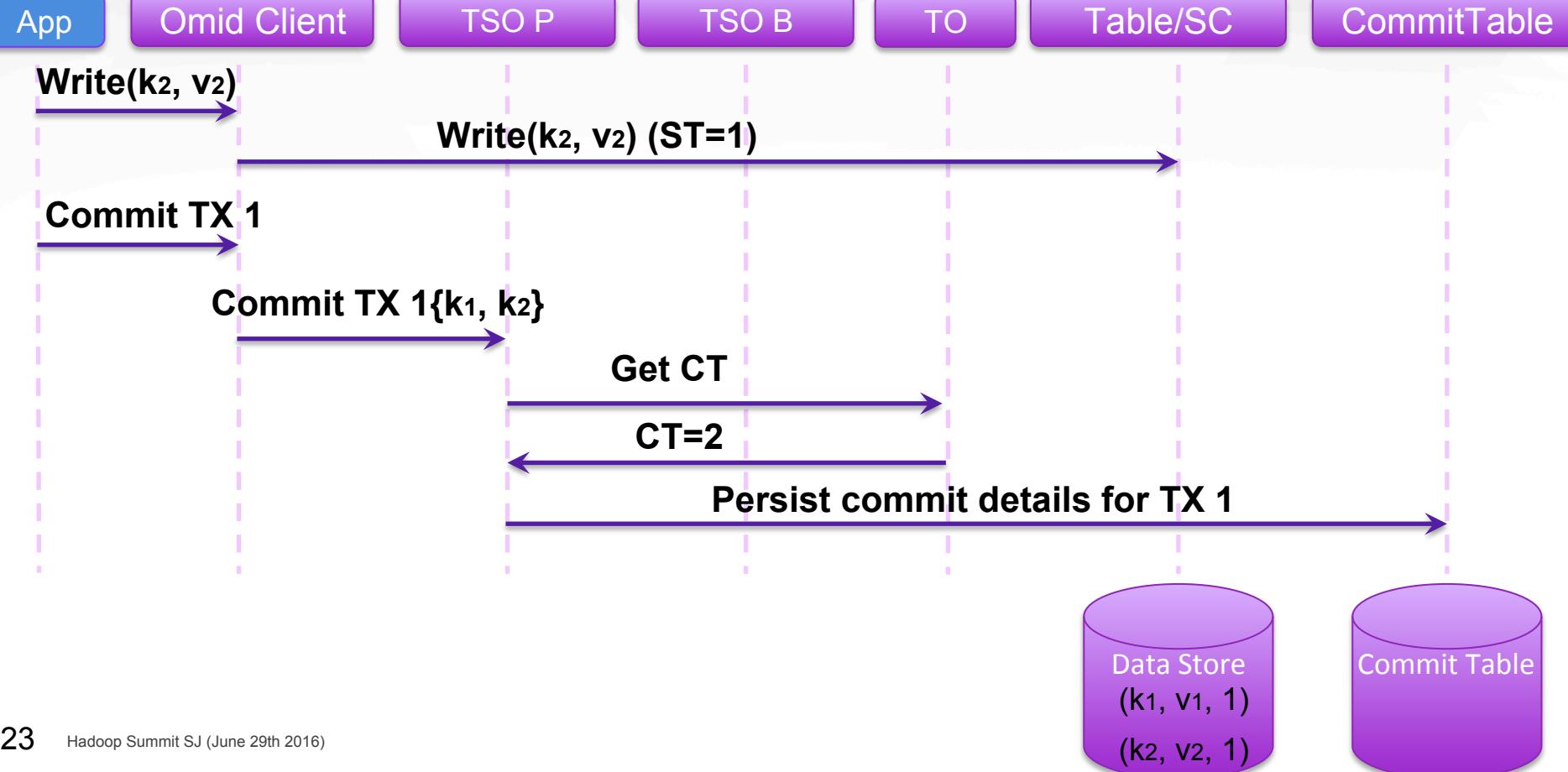
TX 1, 1

TX 1 Write(k1, v1)

Write(k1, v1) (ST=1)



High Availability – Solution



High Availability – Solution

App

Omid Client

TSO B

TO

Table/SC

CommitTable

Begin TX

Begin TX

TSO B

TO

Table/SC

CommitTable

Get ST

ST=3

TX(ST=3,E=3)

TX 3,3

TX 3 Read(k1)

Read(k1) (ST=3)

(k1, v1, 1)

Data Store
(k1, v1, 1)
(k2, v2, 1)

Commit Table

High Availability – Solution

App

Omid Client

TSO B

TO

Table/SC

CommitTable

Return TX1 CT

! exist

Try invalidate

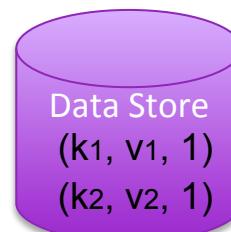
Invalid

! exist

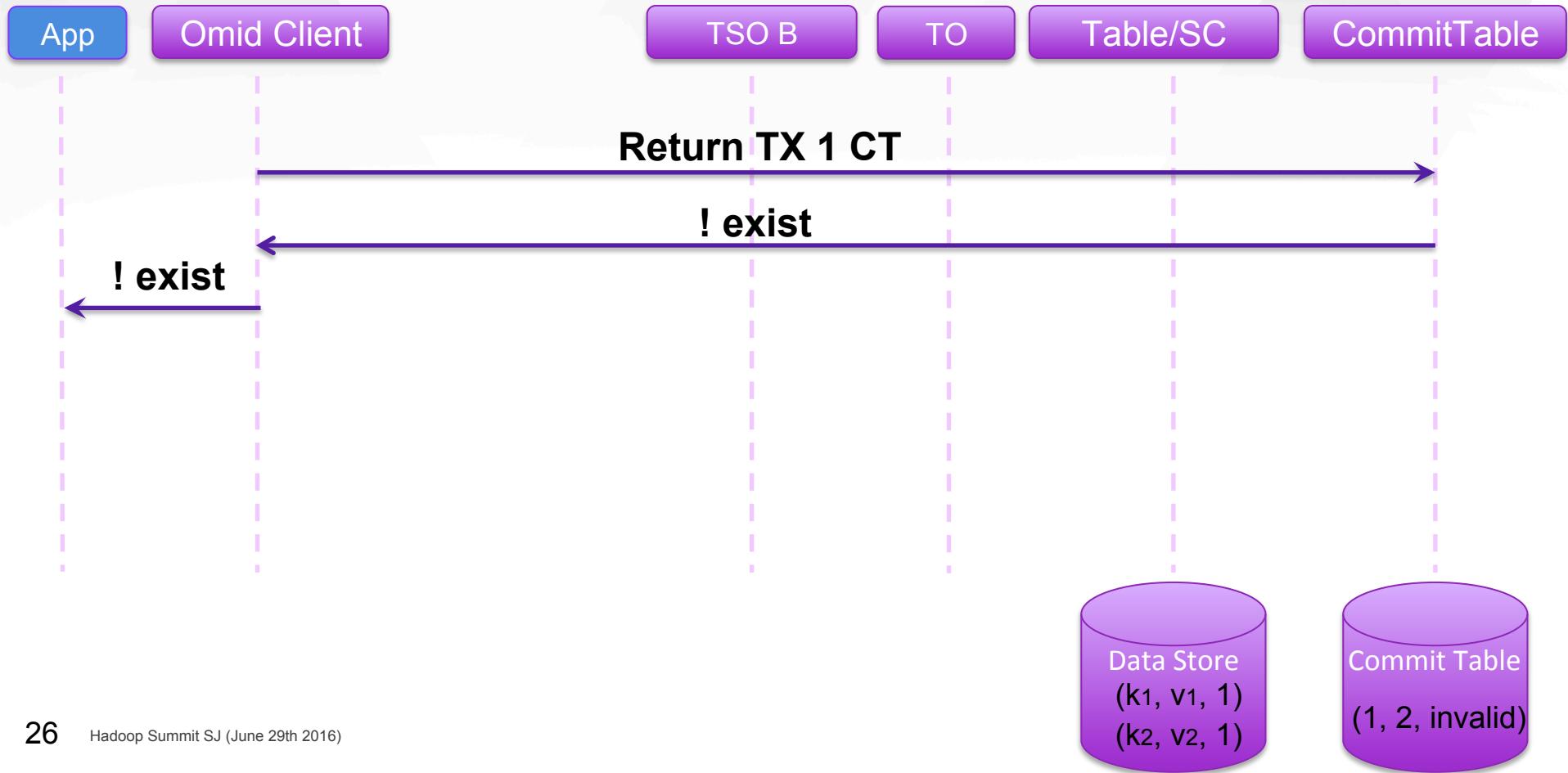
TX 3 Read(k2)

Read(k2) (ST=3)

(k2, v2, 1)



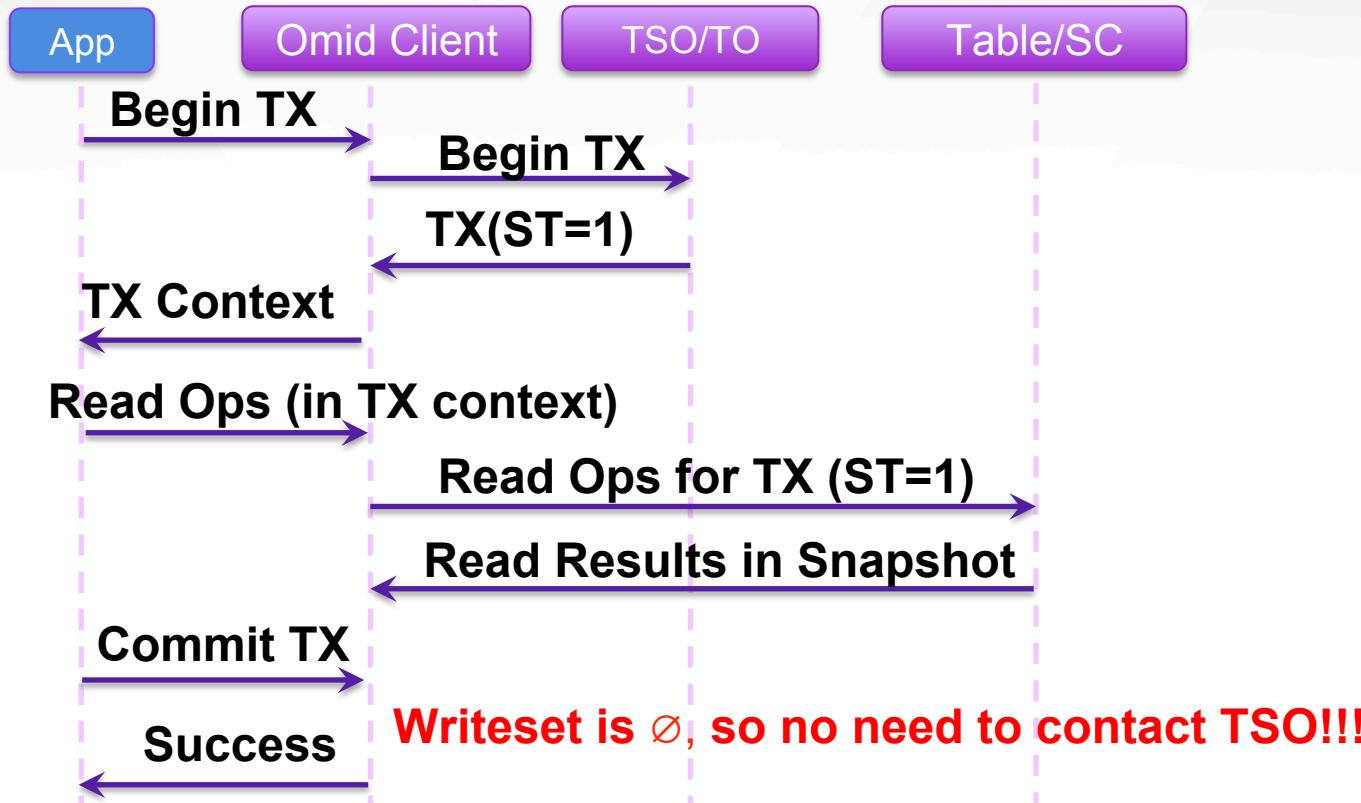
High Availability – Solution



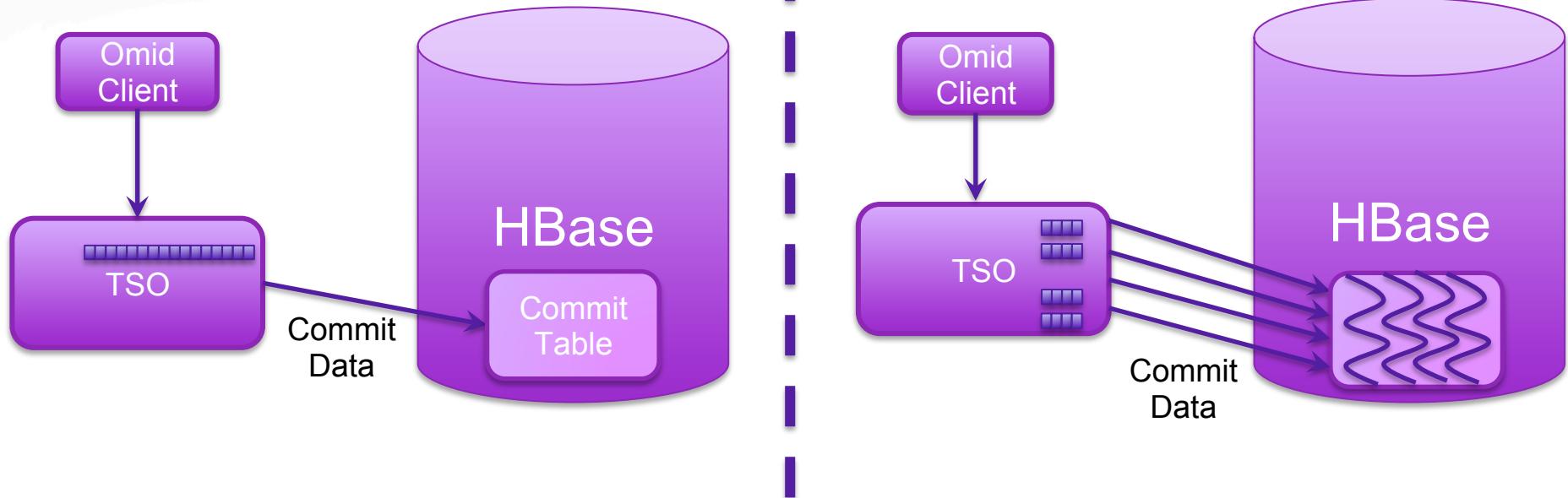
High Availability

- **No runtime overhead** in mainstream execution
 - Minor overhead after failover
- TSO uses **regular writes**
- **Leases** for leader election
 - Lease status check before/after writing to Commit Table

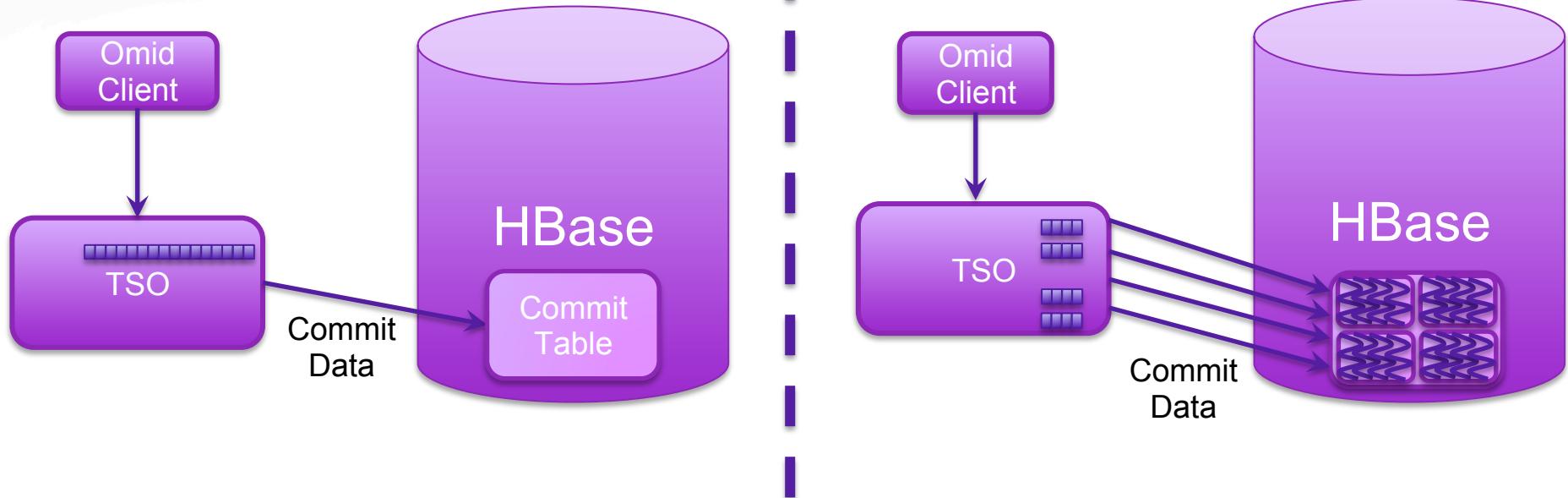
Perf. Improvements: Read-Only Txns



Perf. Improvements: Commit Table Writes



Perf. Improvements: Commit Table Writes



Omid Throughput with Improvements





Summary

- **Transactions in NoSQL**
 - Use cases in incremental big data processing
 - *Snapshot Isolation*: Scalable consistency model
- **Omid**
 - Web-scale TPS for HBase
 - Reliable and performant
 - Battle-tested



<http://omid.incubator.apache.org/>

Questions?



[@ApacheOmid](#)
[Apache Omid Incubator Page](#)