

Schism: a Workload-Driven Approach to Database Replication and Partitioning

Carlo Curino et al.

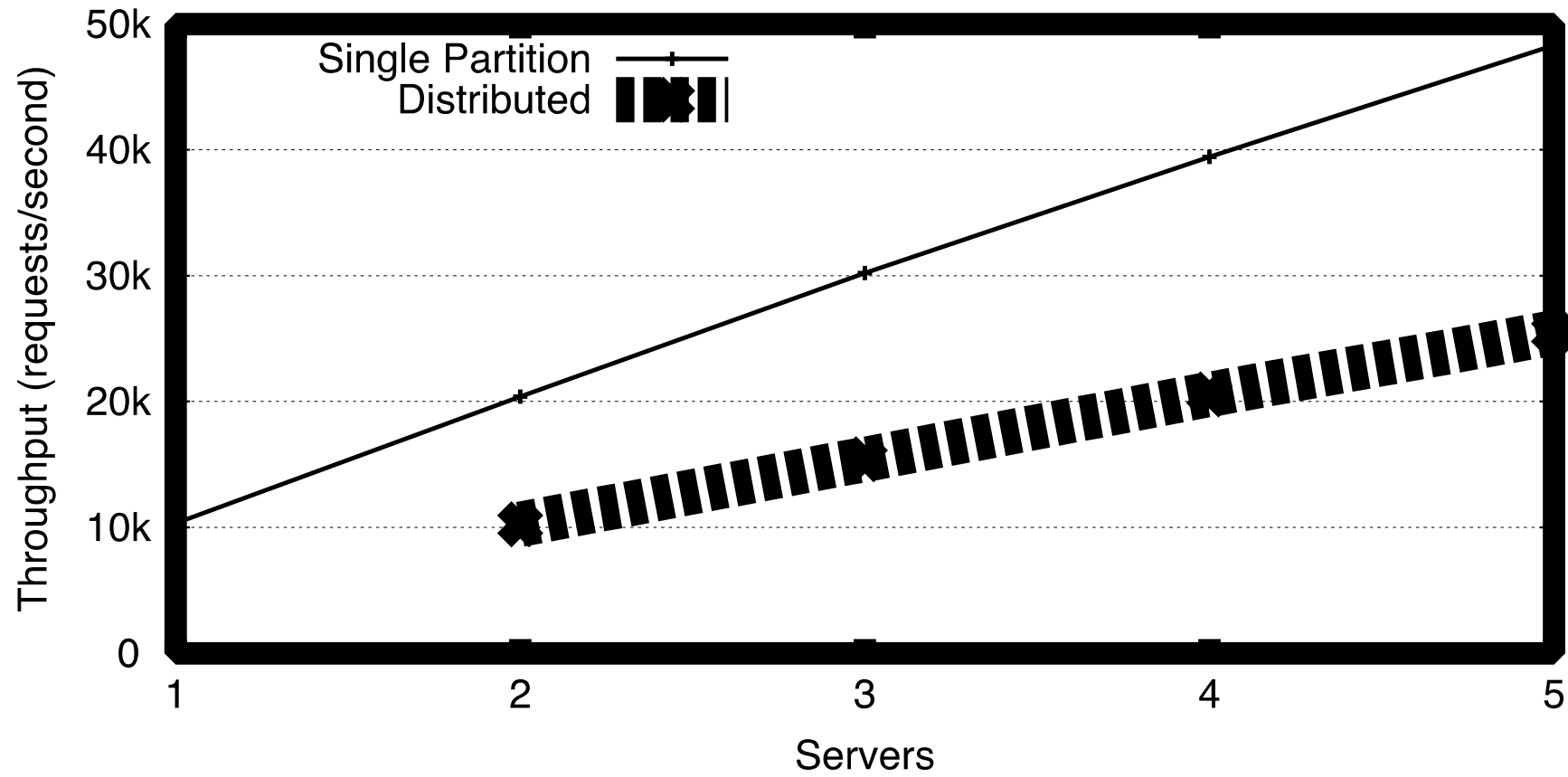
Presented by: Rui Zhang

Motivation

- Why partitioning?
 - Scalability
 - Availability
 - Manageability
- Traditional Partitioning
 - Total replication
 - Round-Robin
 - Range
 - Hash

Motivation(Cont.)

- Distributed transactions are expensive
 - Consensus protocol
 - Distributed locks
 - More communication
- Many to Many relationship
 - Hard to partition



Solution

- Minimize distributed transactions
- Balance the workload of each partition

Schism



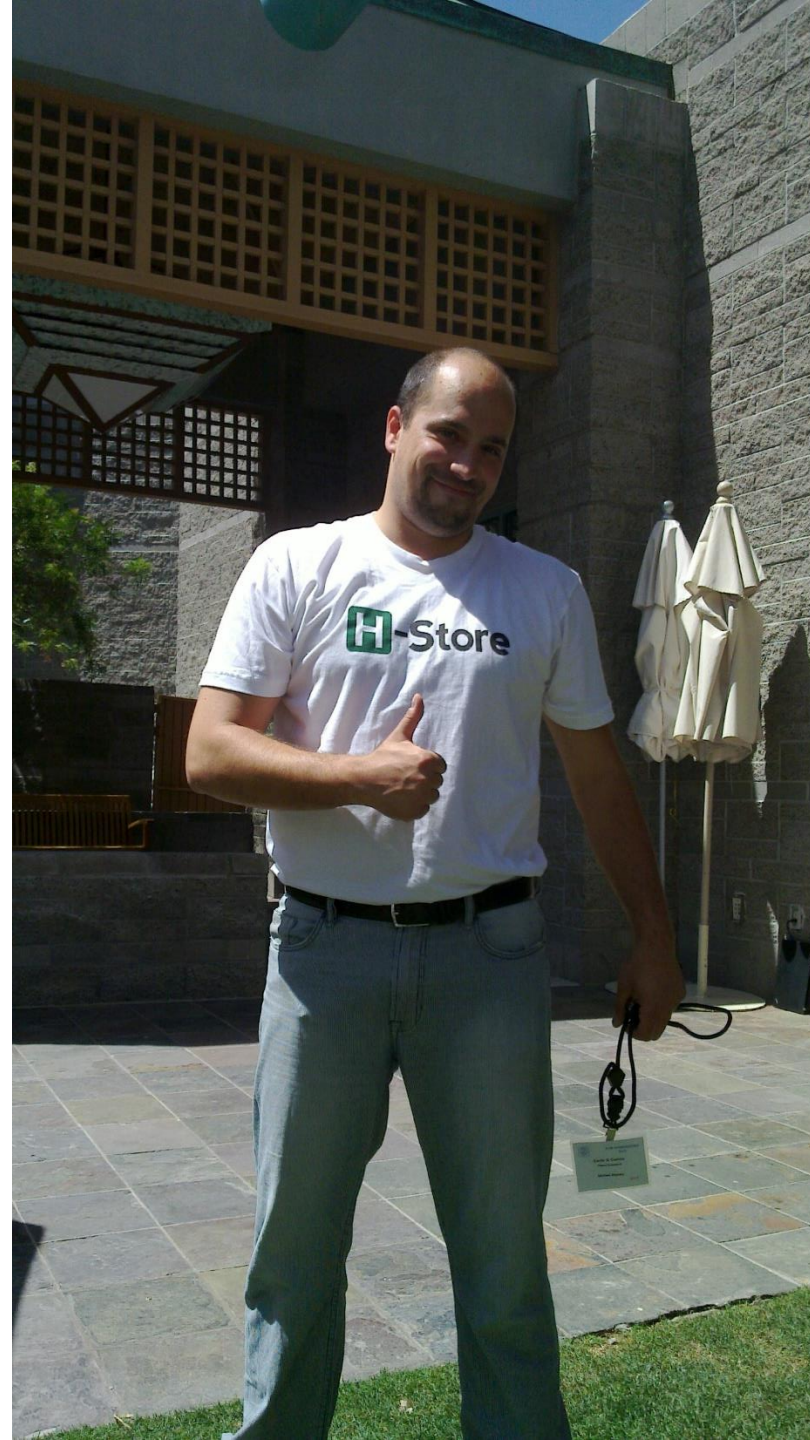
- Part of RelationalCloud System

- An MIT-based effort to investigate technologies and challenges related to Database-as-a-Service (DAAS) within cloud-computing[1]
- Researchers including Carlo, Evan and Sam participate this project

- DAAS

- Cloud computing based service is a novel market for data management
- More and more dedicated data management engines which complicates the process of choosing and deploying
- DAAS is to hide the complexity with a simple interface

[1] <http://relationalcloud.com>



Schism

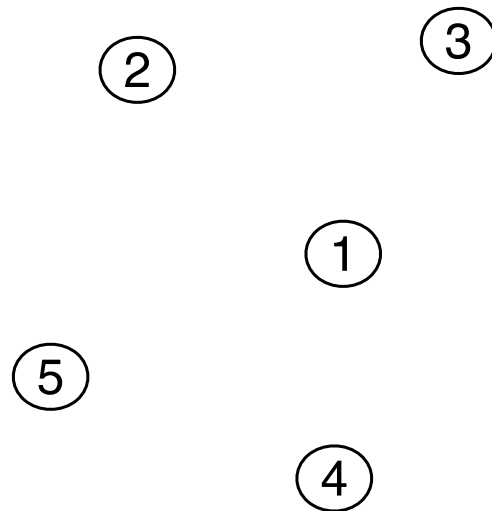
- Input
 - Workload trace
- Output
 - Partitioning and replication plan
- Five steps
 - Data pre-processing
 - Creating the graph
 - Partitioning the graph
 - Explaining the partition
 - Final Validation

Creating the graph

- Node: Tuple
- Edge: Transaction
- Weight: Number of transactions

Creating the graph

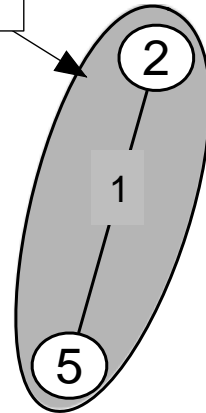
account		
id	name	bal
1	carlo	80k
2	evan	60k
3	sam	129k
4	eugene	29k
5	yang	12k
...



Creating the graph

```
BEGIN
UPDATE account SET bal=60k
  WHERE id=2;
SELECT * FROM account
  WHERE id=5;
COMMIT
```

account		
id	name	bal
1	carlo	80k
2	evan	60k
3	sam	129k
4	eugene	29k
5	yang	12k
...



transaction edges

3

1

4

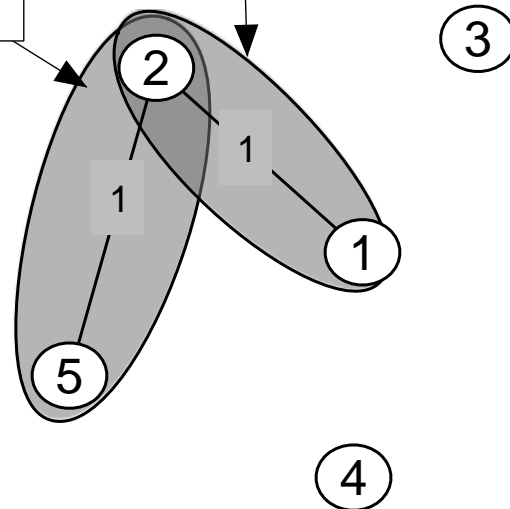
Creating the graph

```
BEGIN
UPDATE account SET bal=60k
  WHERE id=2;
SELECT * FROM account
  WHERE id=5;
COMMIT
```

account		
id	name	bal
1	carlo	80k
2	evan	60k
3	sam	129k
4	eugene	29k
5	yang	12k
...

```
BEGIN
UPDATE account SET bal=bal-1k WHERE name="carlo";
UPDATE account SET bal=bal+1k WHERE name="evan";
COMMIT
```

▬ transaction edges



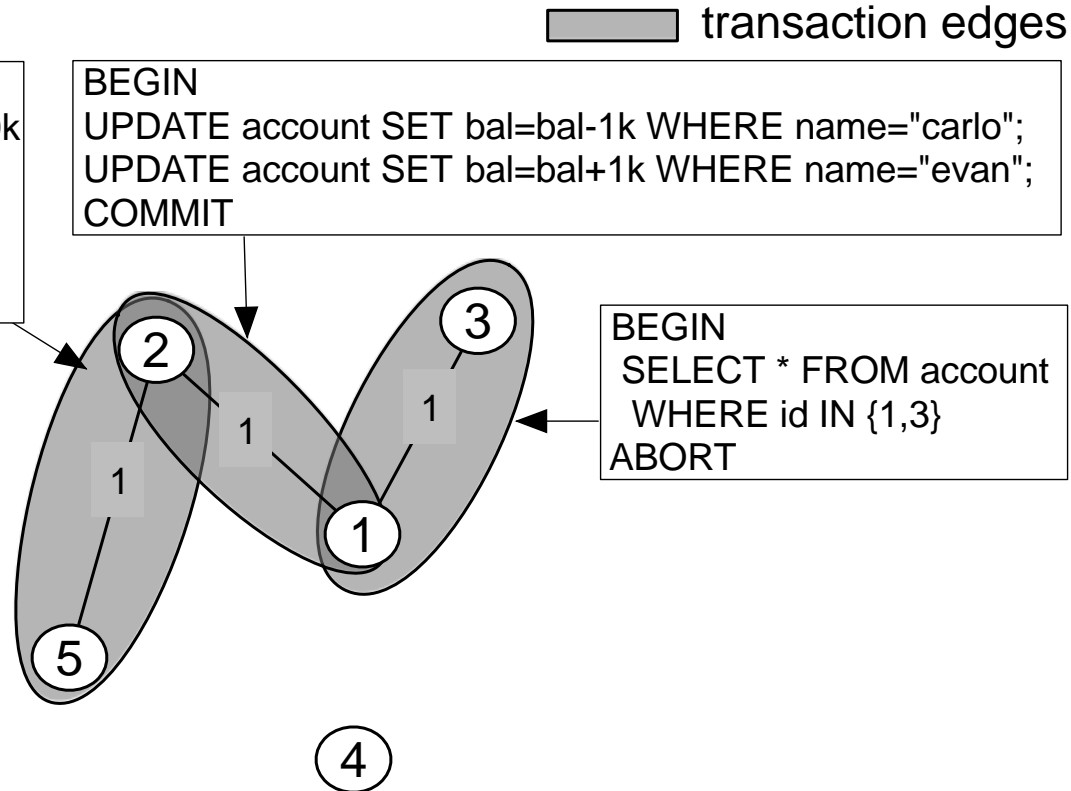
Creating the graph

```
BEGIN
UPDATE account SET bal=60k
WHERE id=2;
SELECT * FROM account
WHERE id=5;
COMMIT
```

account		
id	name	bal
1	carlo	80k
2	evan	60k
3	sam	129k
4	eugene	29k
5	yang	12k
...

```
BEGIN
UPDATE account SET bal=bal-1k WHERE name="carlo";
UPDATE account SET bal=bal+1k WHERE name="evan";
COMMIT
```

```
BEGIN
SELECT * FROM account
WHERE id IN {1,3}
ABORT
```



Creating the graph

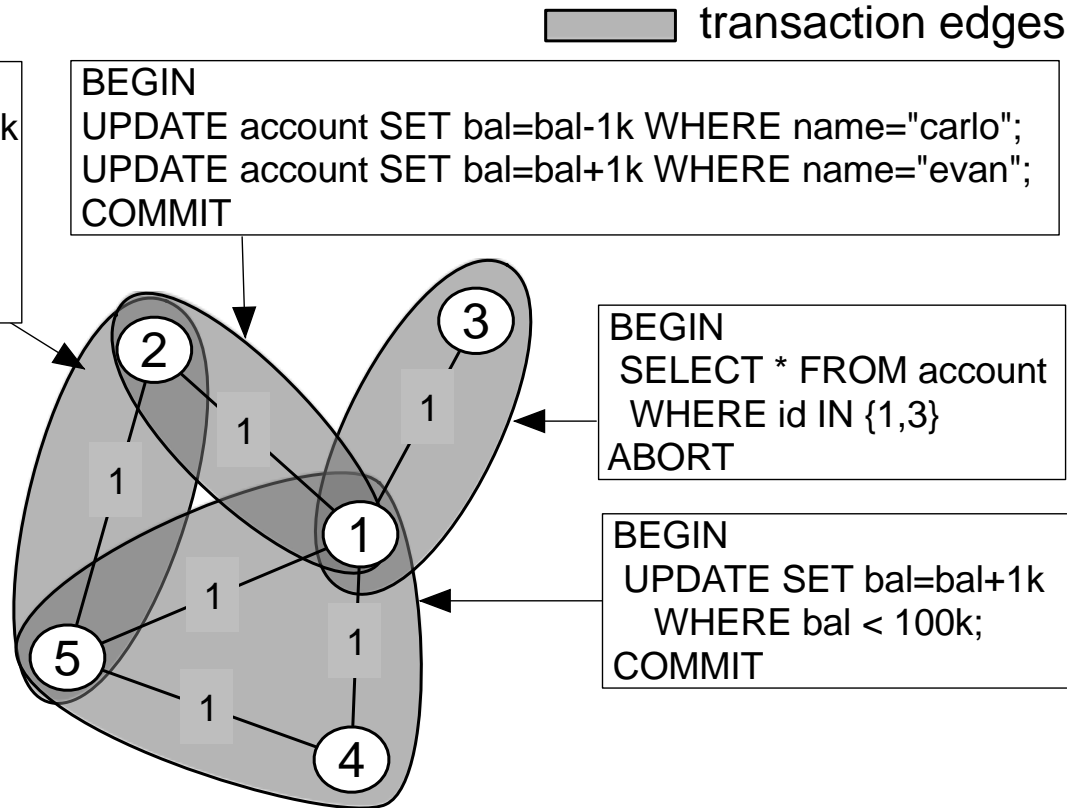
```
BEGIN
UPDATE account SET bal=60k
  WHERE id=2;
SELECT * FROM account
  WHERE id=5;
COMMIT
```

account		
id	name	bal
1	carlo	80k
2	evan	60k
3	sam	129k
4	eugene	29k
5	yang	12k
...

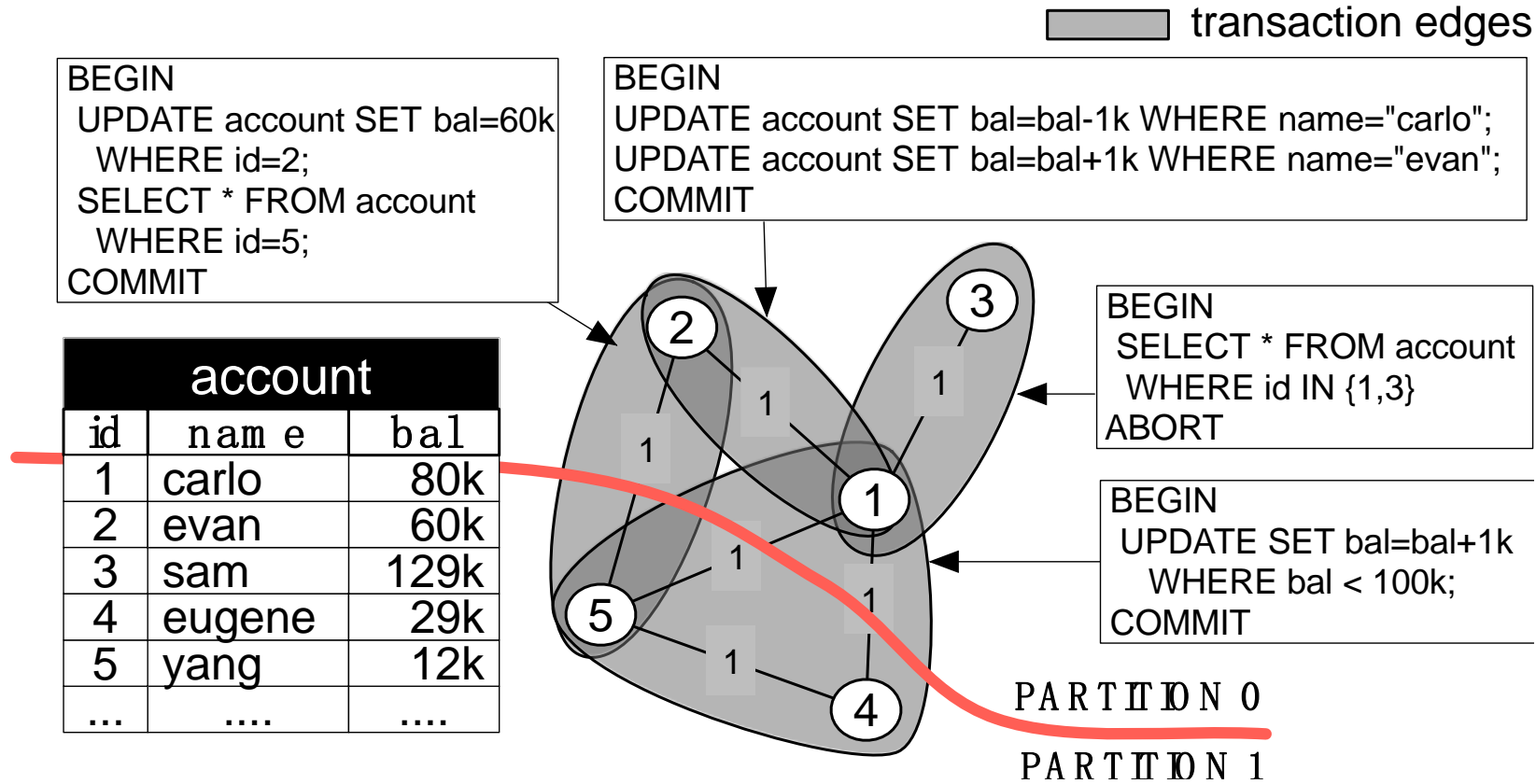
```
BEGIN
UPDATE account SET bal=bal-1k WHERE name="carlo";
UPDATE account SET bal=bal+1k WHERE name="evan";
COMMIT
```

```
BEGIN
SELECT * FROM account
  WHERE id IN {1,3}
ABORT
```

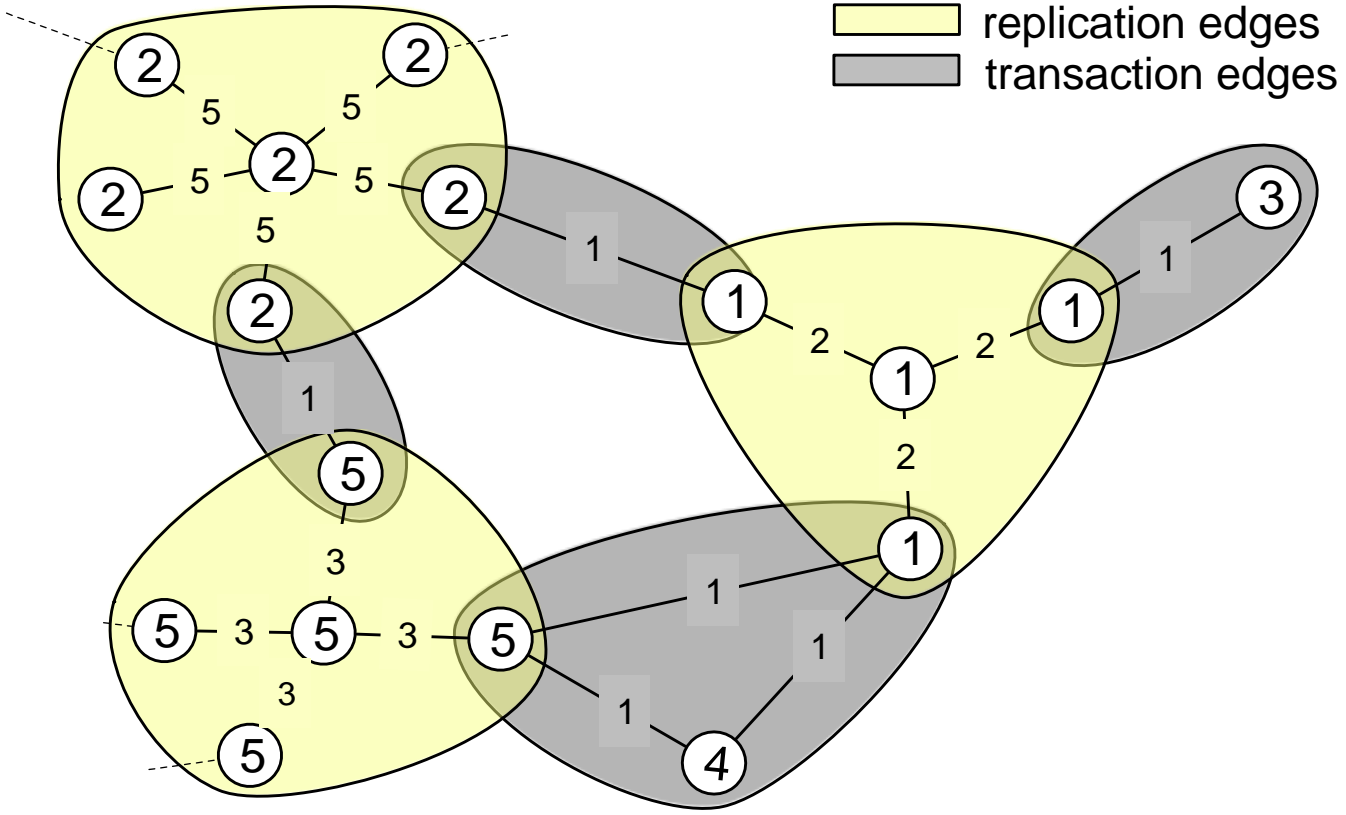
```
BEGIN
UPDATE SET bal=bal+1k
  WHERE bal < 100k;
COMMIT
```



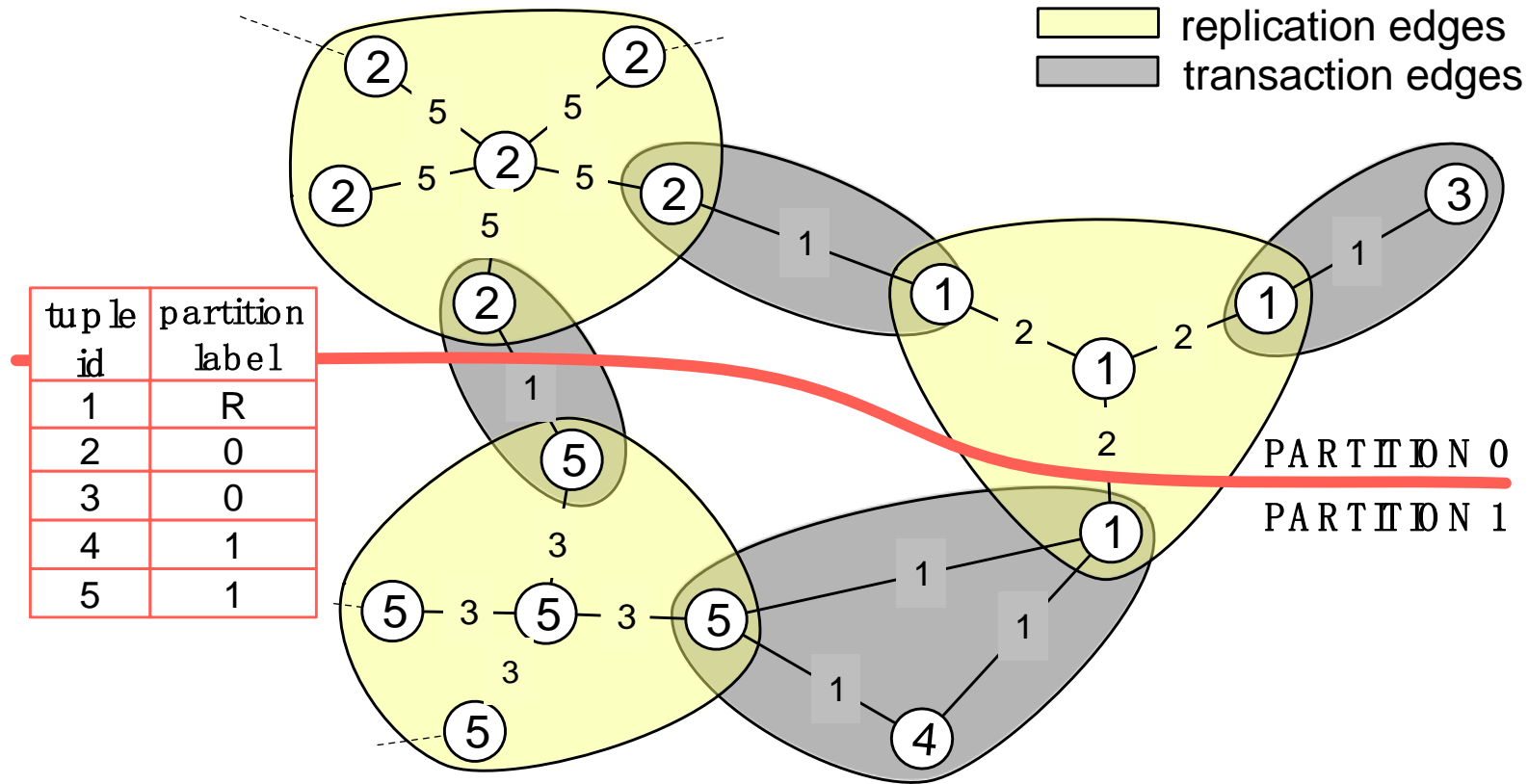
Creating the graph



Tuple Replication



Tuple Replication



Graph Partitioning

- Metis algorithm

- Coarsening
- Initial partitioning
- Refinement

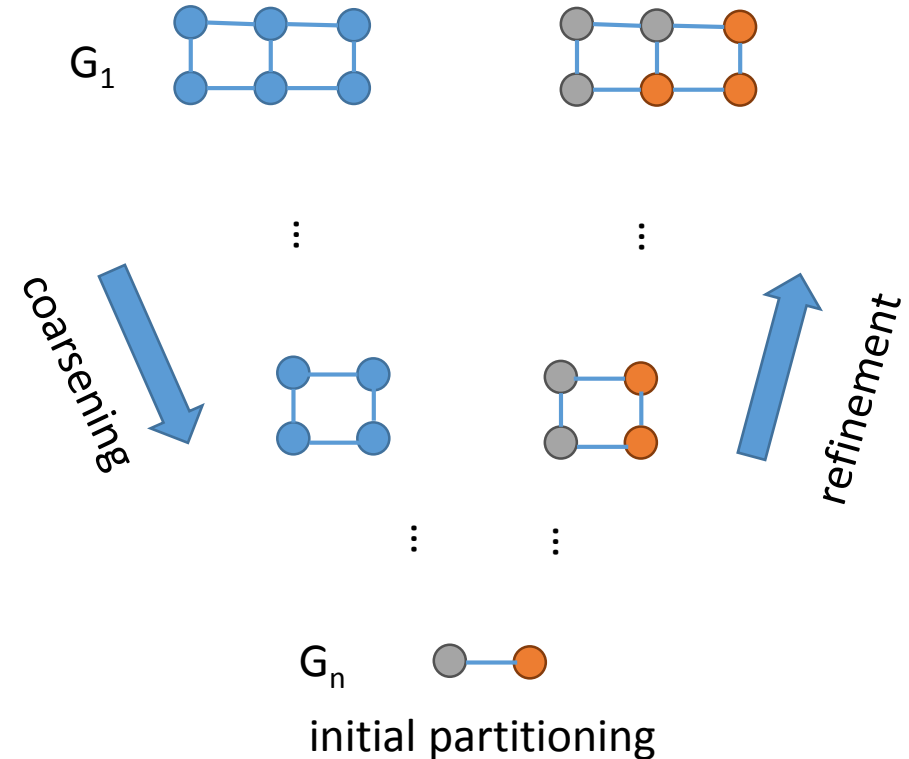
See www.cs.umn.edu/~metis for more detail

- Goal

- Minimum cut and balanced partitions

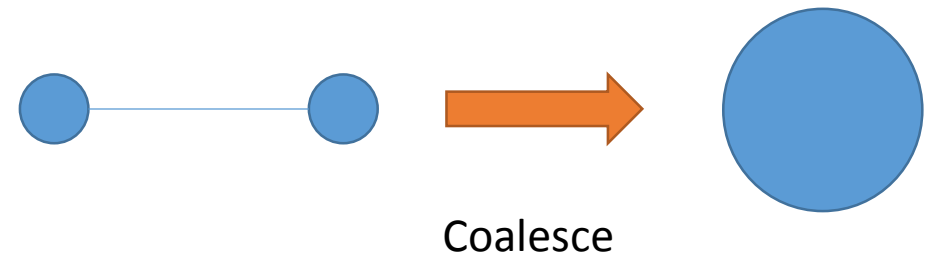
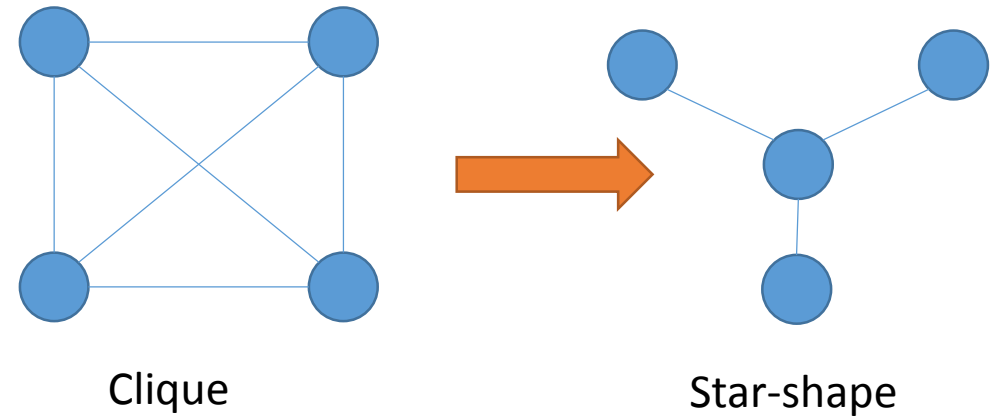
- Output

- Lookup table



Graph Size Reduction

- Sampling
 - Transaction level
 - Tuple level
- Blanket-statement filtering
- Relevance filtering
- Star-shaped replication
- Tuple coalescing



Explanation Phase

- Method
 - Decision Tree
- Input
 - Tuple->partition mappings
- Output
 - Compact Model

tuple id	partition label
1	R
2	0
3	0
4	1
5	1



(id = 1) -> partitions = {0, 1}
(2 <= id < 4) -> partition = 0
(id >= 4) -> partition = 1

Explanation Phase

- Candidate attributes
 - Frequent used attributes in WHERE clauses
 - Item ids and warehouse ids for “stock”
 - Join predicates
 - Co-location of the joined tuples
 - Limited use (perform the join) but supported

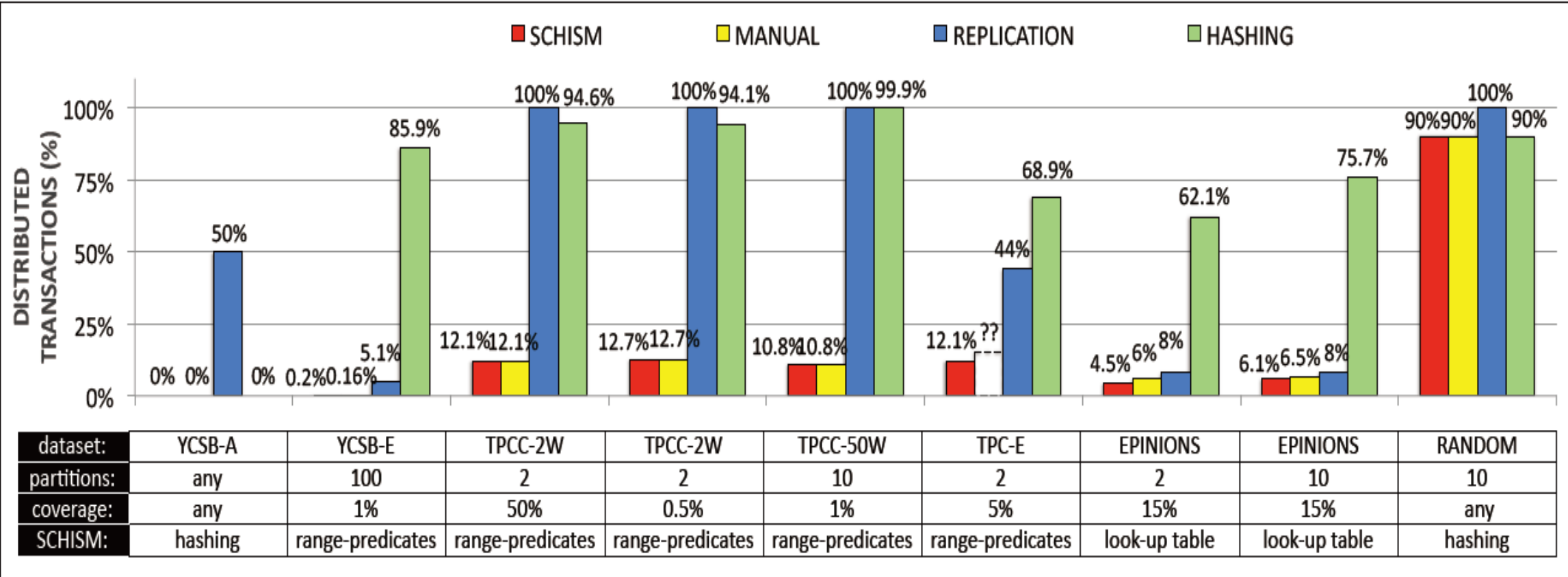
Query/Tuple Routing

- Centralized Router
 - A decision tree
 - Determine which partition the query or tuple belongs to
 - Spare partition that store tuples don't hit
 - If full, rerun the partitioning algorithm and retrain the decision tree
 - After retraining, some tuples should be migrated from one partition to another via data migration SQL scripts
 - Broadcast to all partitions for attributes that are not partitioning attributes

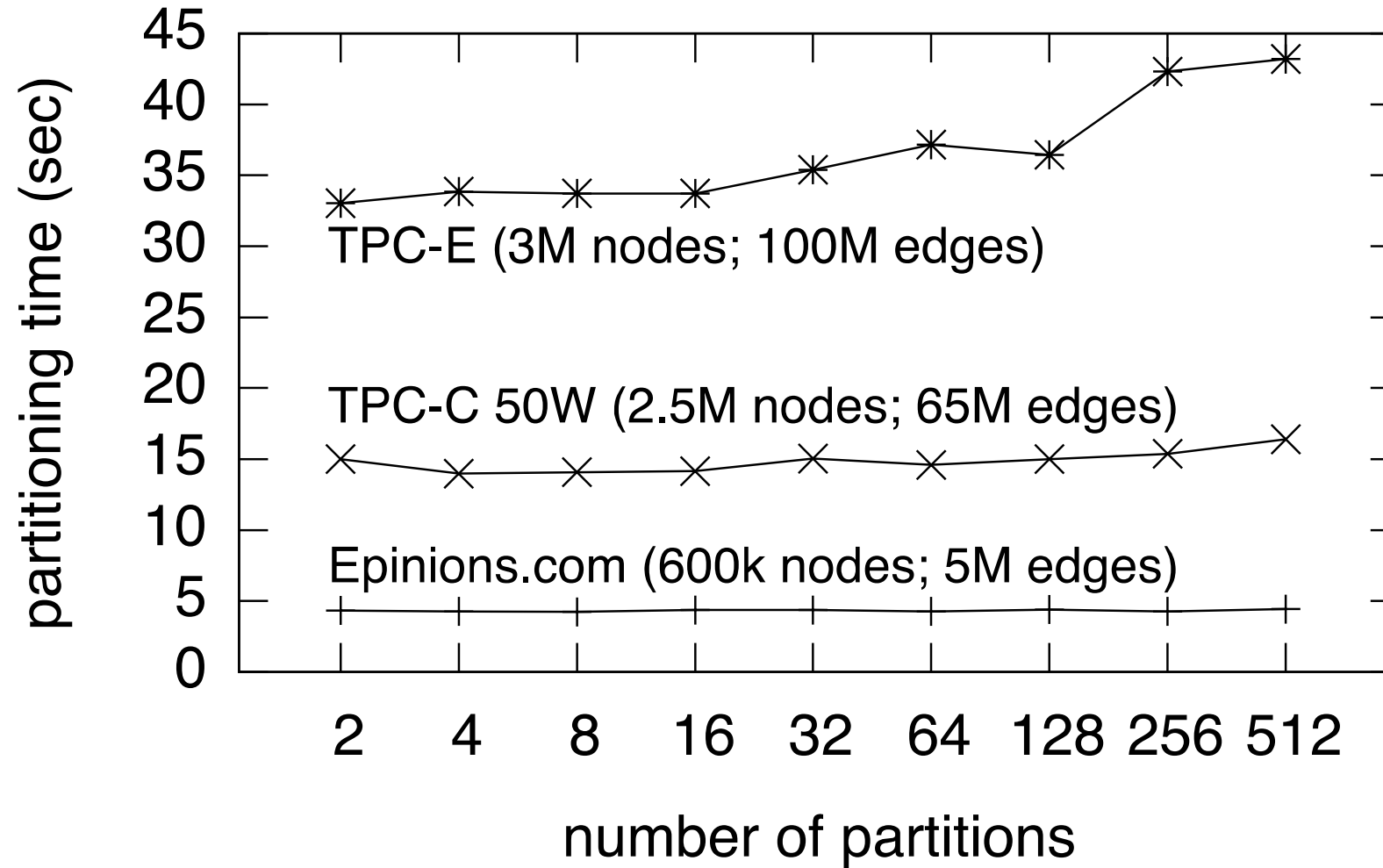
Evaluation

- Schism
- Manual
- Total Replication
- Hash partitioning

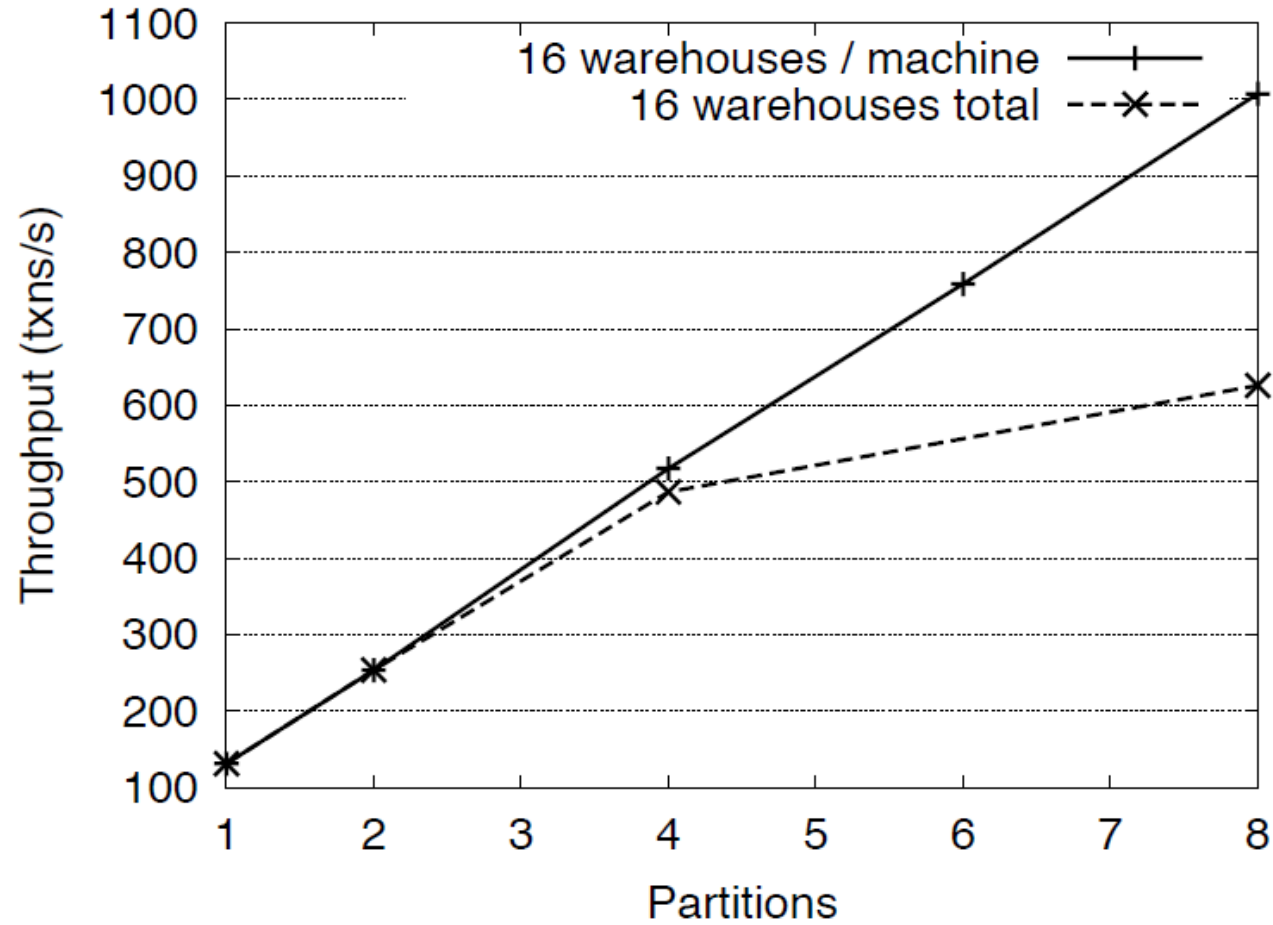
Evaluation



Evaluation



Evaluation



Limitation

- Not suitable for insert-heavy workloads
 - Involves many data migrations
- Not suitable for databases with large size
 - Partitioning time will be long
- There will be false positive in decision tree
 - Overhead to correct this error
- Compared with Horticulture, it does not consider temporal skew