

# Reciprocal Consistency in Distributed Graph Databases



Jack Waudby

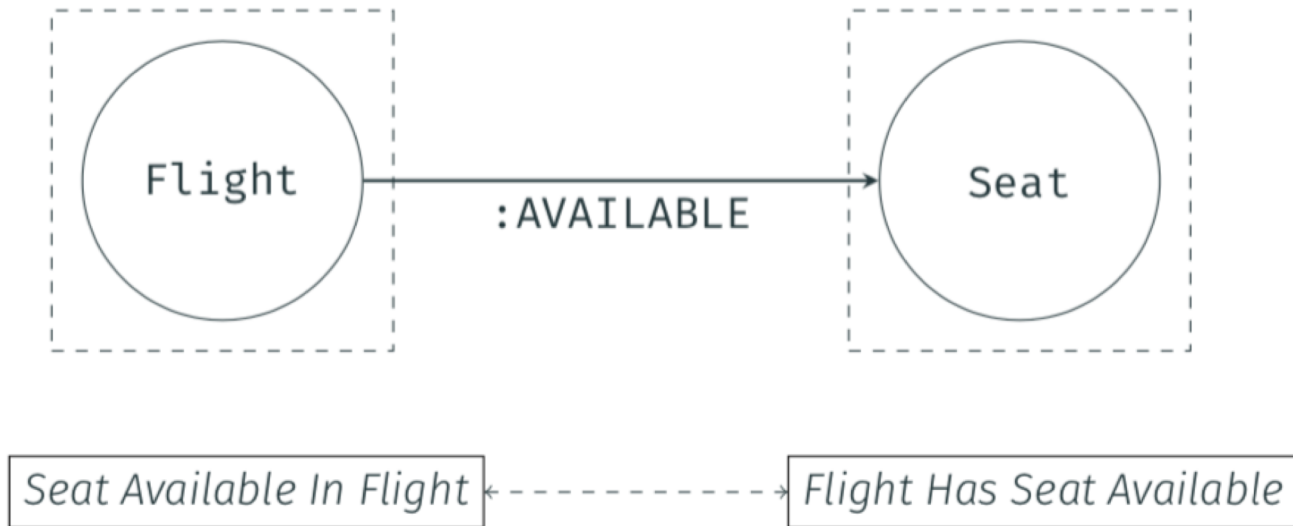


*Supervisor: Paul Ezhilchelvan*

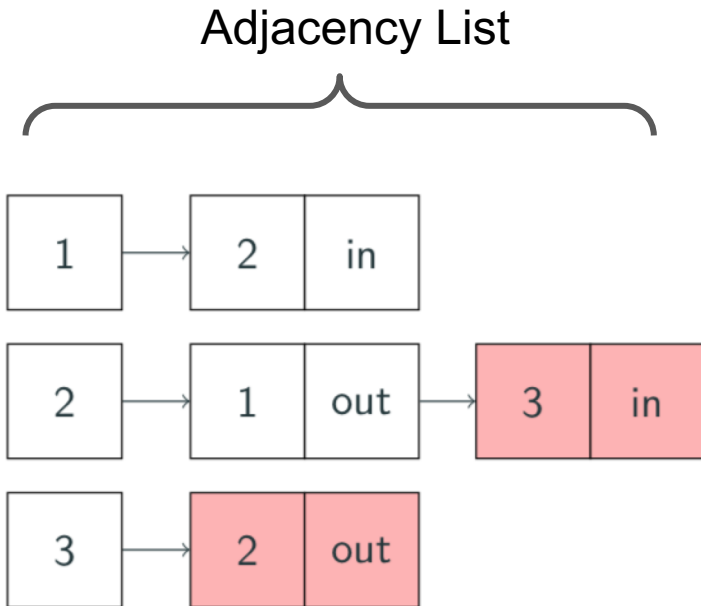
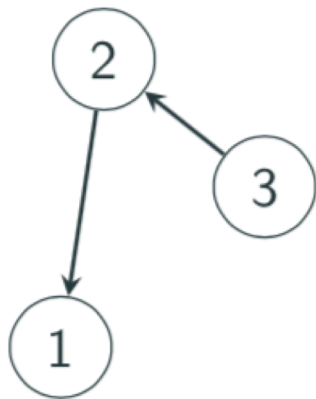
*(in collaboration with Jim Webber, Neo4j)*

# Reciprocal Consistency

An edge is represented by **two physical records**



# Reciprocal Consistency in Practice



# Distributed Graph Databases

*“Many graphs in practice are very large, often containing over a billion edges”*

- 42 survey participants reported graphs with 1-10B edges
- 17 reported graphs with 10B-100B edges
- **7 reported graphs with 100B+ edges!**



Siddhartha Sahu, Amine Mhedhbi, Semih Salihoglu, Jimmy Lin, M. Tamer Özsu  
*The Ubiquity of Large Graphs and Surprising Challenges of Graph Processing*  
VLDB 2017

# Distributed Graph Databases

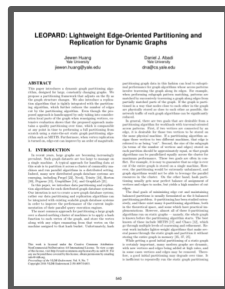
Partition data across multiple machines in a cluster

For key-value stores partitioning by keyspace does the trick

**Graph partitioning is non-trivial!**

# $k$ -Balanced Partitioning Problem

*“The dual goals of minimizing edge cut and maintaining balanced partitions”*

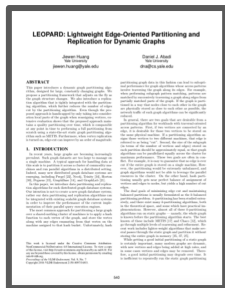
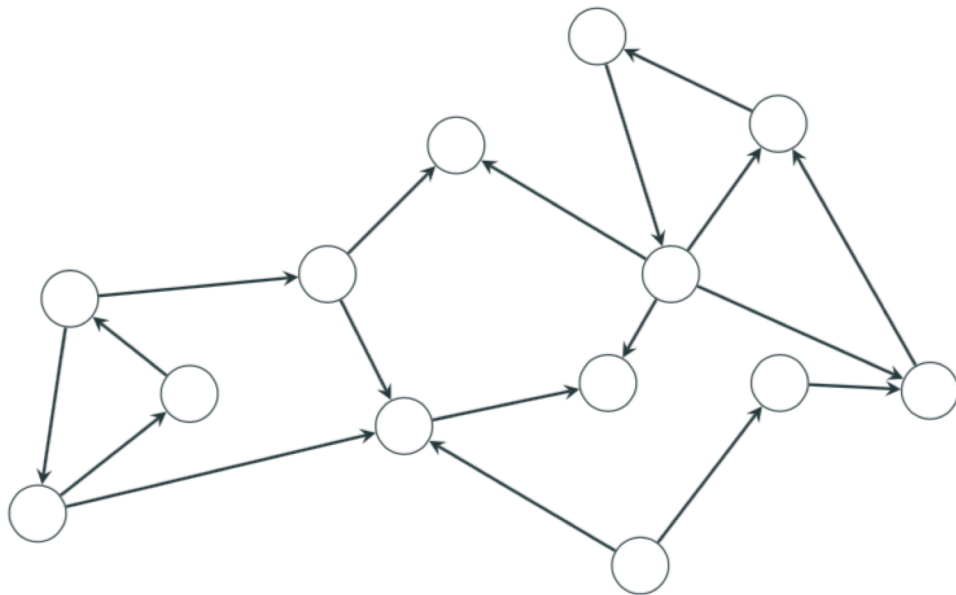


Jiewen Huang and Daniel J. Abadi

*LEOPARD: Lightweight Edge-Oriented Partitioning and Replication for Dynamic Graphs*

VLDB 2016

# $k$ -Balanced Partitioning Problem

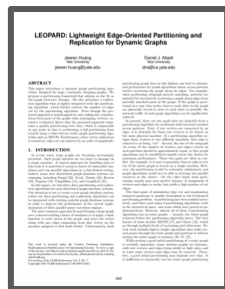
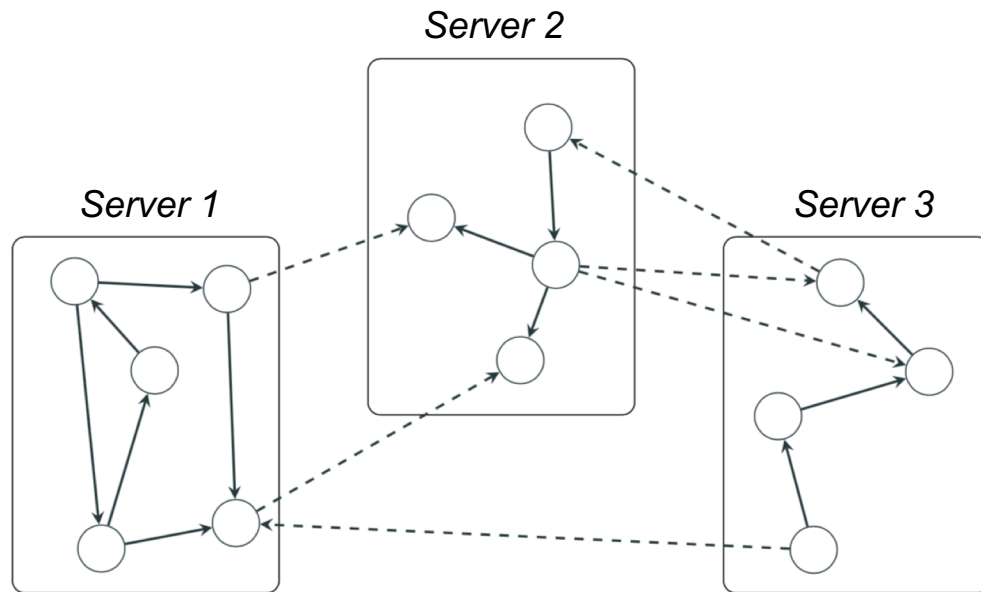


Jiewen Huang and Daniel J. Abadi

*LEOPARD: Lightweight Edge-Oriented Partitioning and Replication for Dynamic Graphs*

VLDB 2016

# $k$ -Balanced Partitioning Problem



Jiewen Huang and Daniel J. Abadi

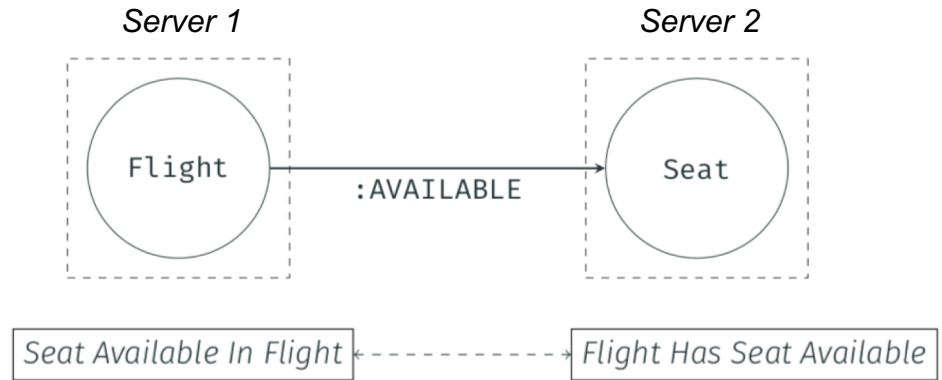
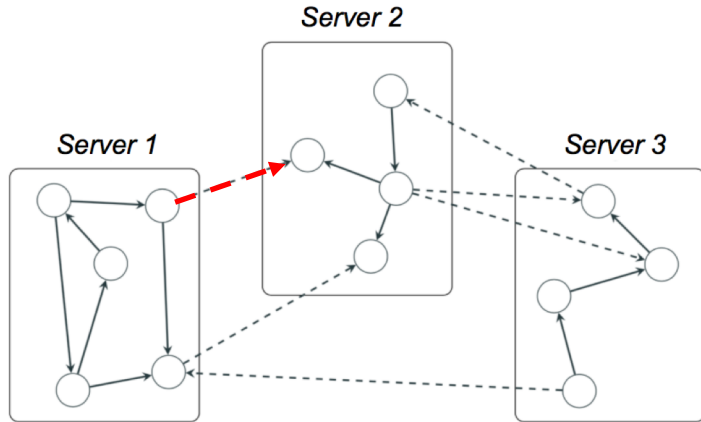
*LEOPARD: Lightweight Edge-Oriented Partitioning and Replication for Dynamic Graphs*

VLDB 2016



# Distributed Edges

Reciprocal information resides on different servers



# A Distributed Graph Database Architecture



Graph Layer



NoSQL Storage Engine

# Reciprocal Consistency Violation

- NoSQL datastores provide weak isolation between concurrent transactions
- Distributed edges can become **reciprocal inconsistent**

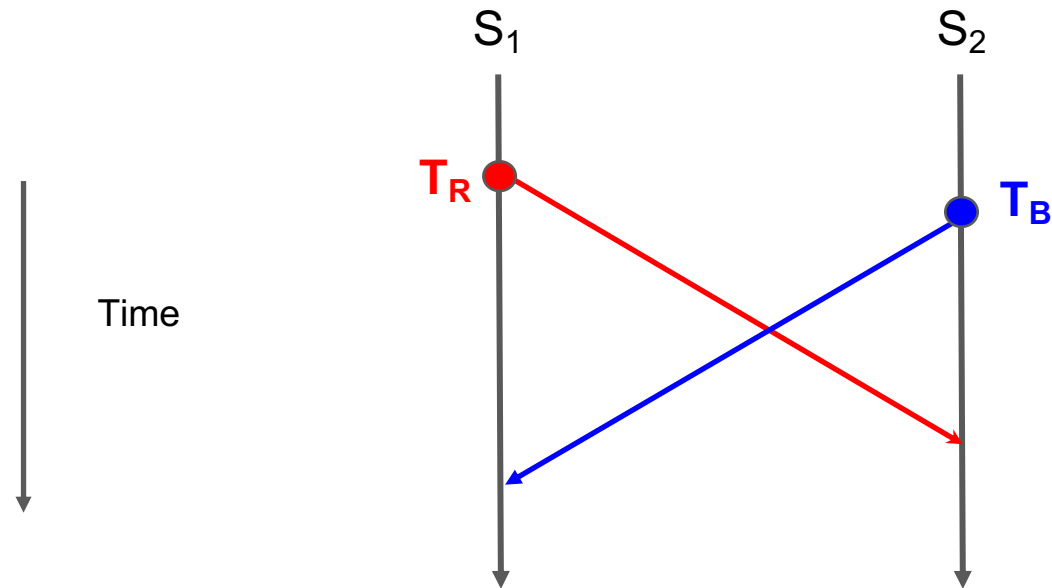
# Reciprocal Consistency Violation: Example

Two concurrent transactions:

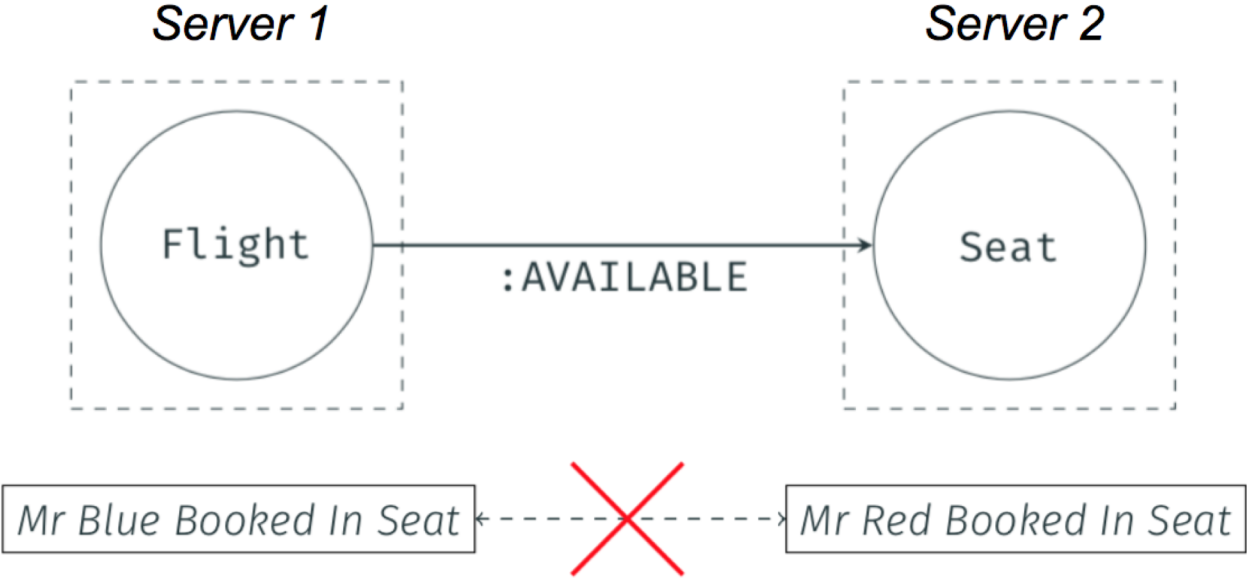
- Mr. Red requests to book the seat,  $T_R$
- Mr. Blue requests to book the seat,  $T_B$



# Reciprocal Consistency Violation



# Reciprocal Consistency Violation



# Spread of Corruption

1. Subsequent transactions read corrupt edge(s)
2. Performs write based on read
3. Corruption spreads through the database

**How quickly does this happen?**

# Spread of Corruption

*Model the rate of corruption in distributed graph databases with weak isolation*

Parameters:

- 1 billion nodes, 11 billion edges (30% distributed)
- 5 edge types with varying access probabilities
- Reads per query geometrically distributed with avg. 15 edges
- 10% queries read-write transactions



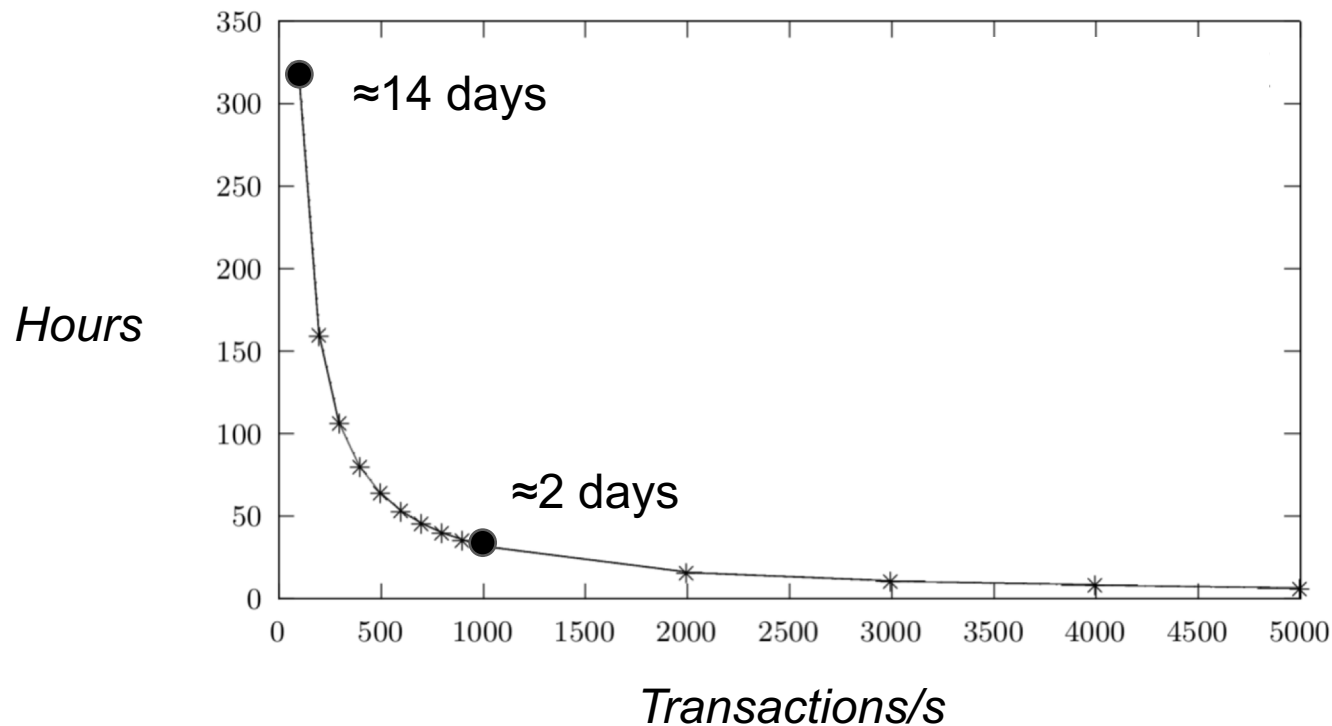
Paul Ezhilchelvan, Isi Mitrani, Jim Webber,

*On the degradation of distributed graph databases with eventual consistency*,  
EPEW 2018



# Simulation Results

Time taken until 10% edges were in a corrupted state



# Preserving Reciprocal Consistency

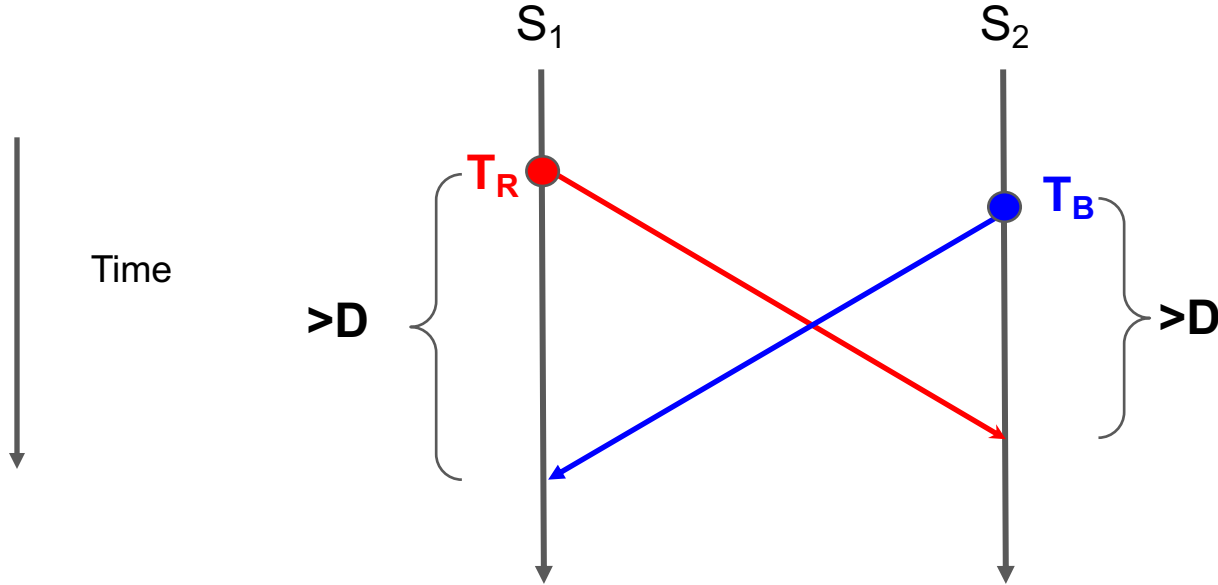
1. Delta Protocol
2. Collision Detection Protocol

# Delta Protocol

- A write to a distributed edge consists of 2 writes
- Assume a bound  $D$
- $D$  reflects time taken by a transaction to write a distributed edge
- Rule: when attempting to write if another write exists within  $D$ , abort

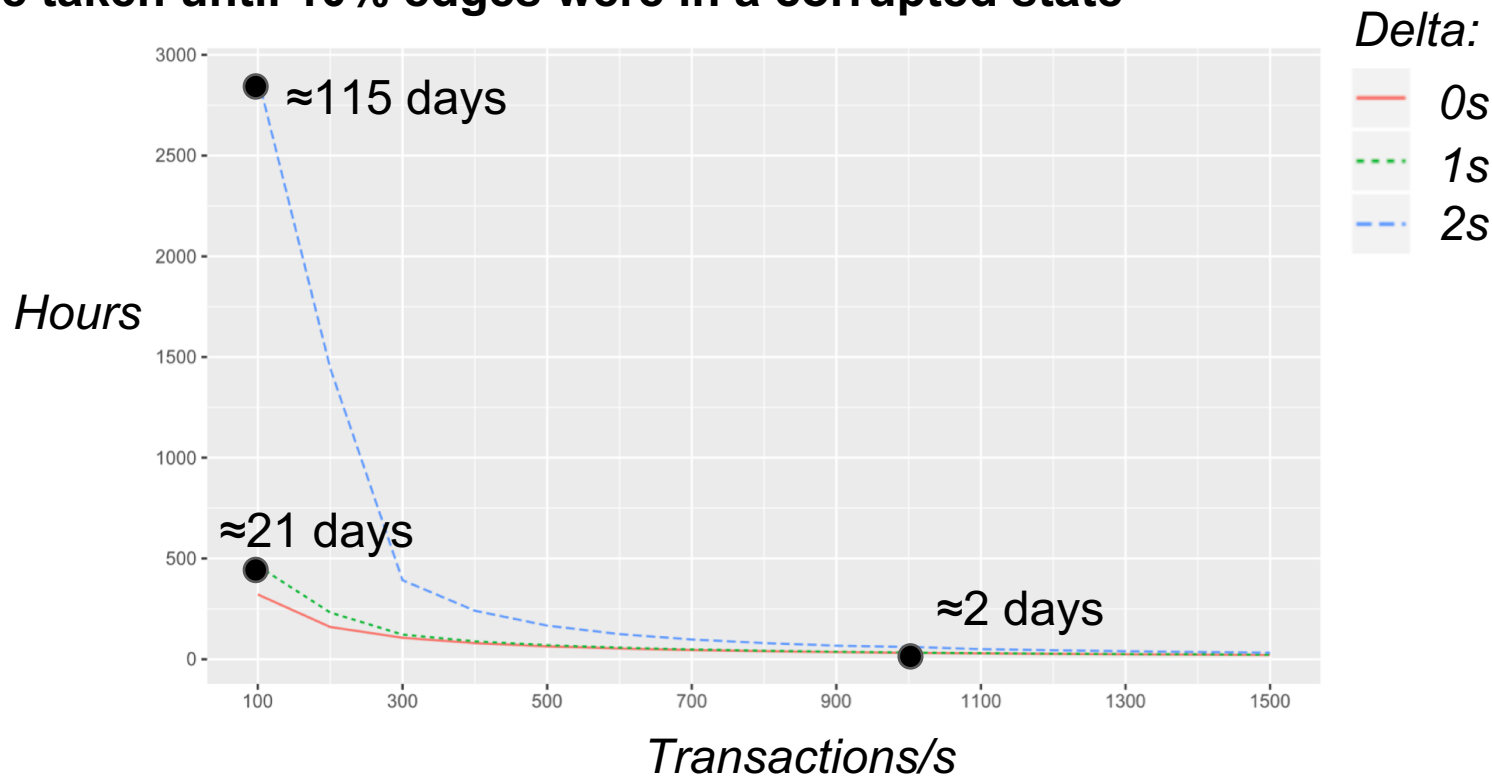
**Scenarios exist when corruption can still occur**

# Reciprocal Consistency Violation



# Delta Protocol: Simulation Results

Time taken until 10% edges were in a corrupted state



# Collision Detection Protocol

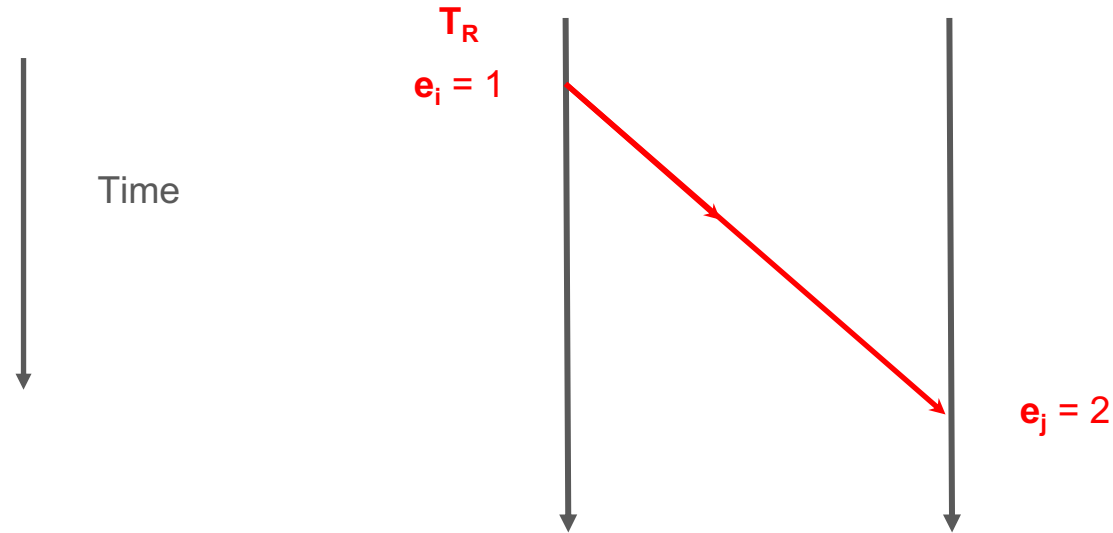
- Introduce **tracking metadata**
- Transactions annotate their writes to distributed edges with a “1” and “2”

**Rule: For any “1” seen, there must be a “2” in the opposite end.**

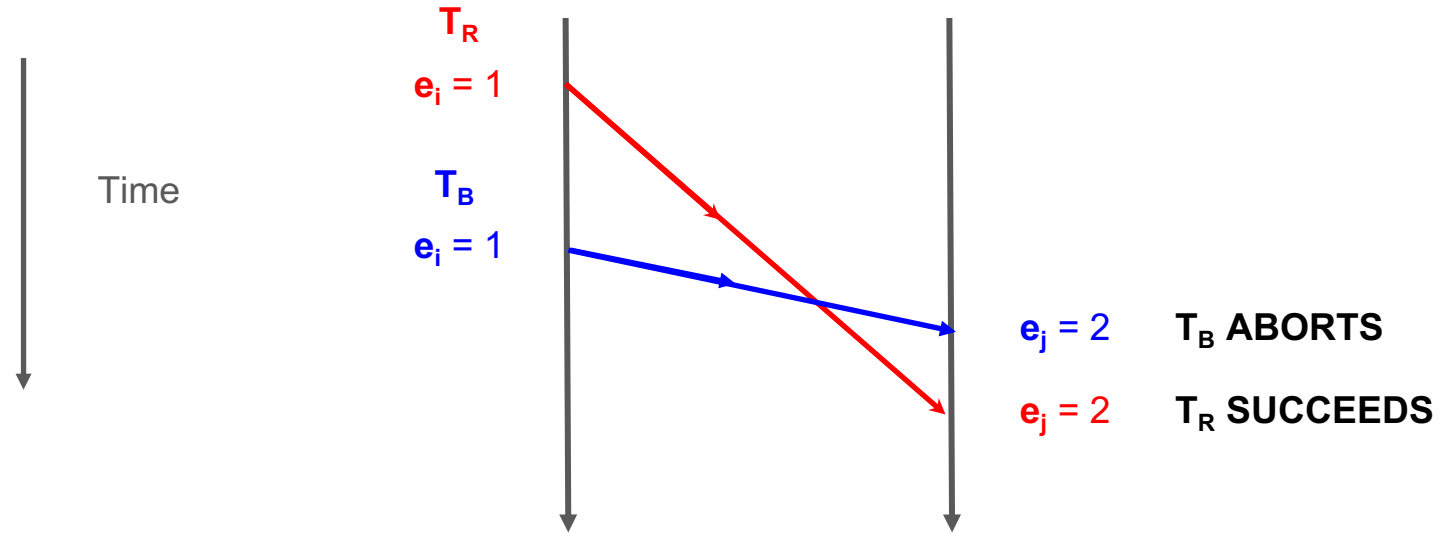


Paul Ezhilchelvan, Isi Mitrani, Jack Waudby, Jim Webber,  
*Design and Evaluation of an Edge Concurrency Control Protocol for  
Distributed Graph Databases*  
EPEW 2019

# Collision Detection Protocol: Example

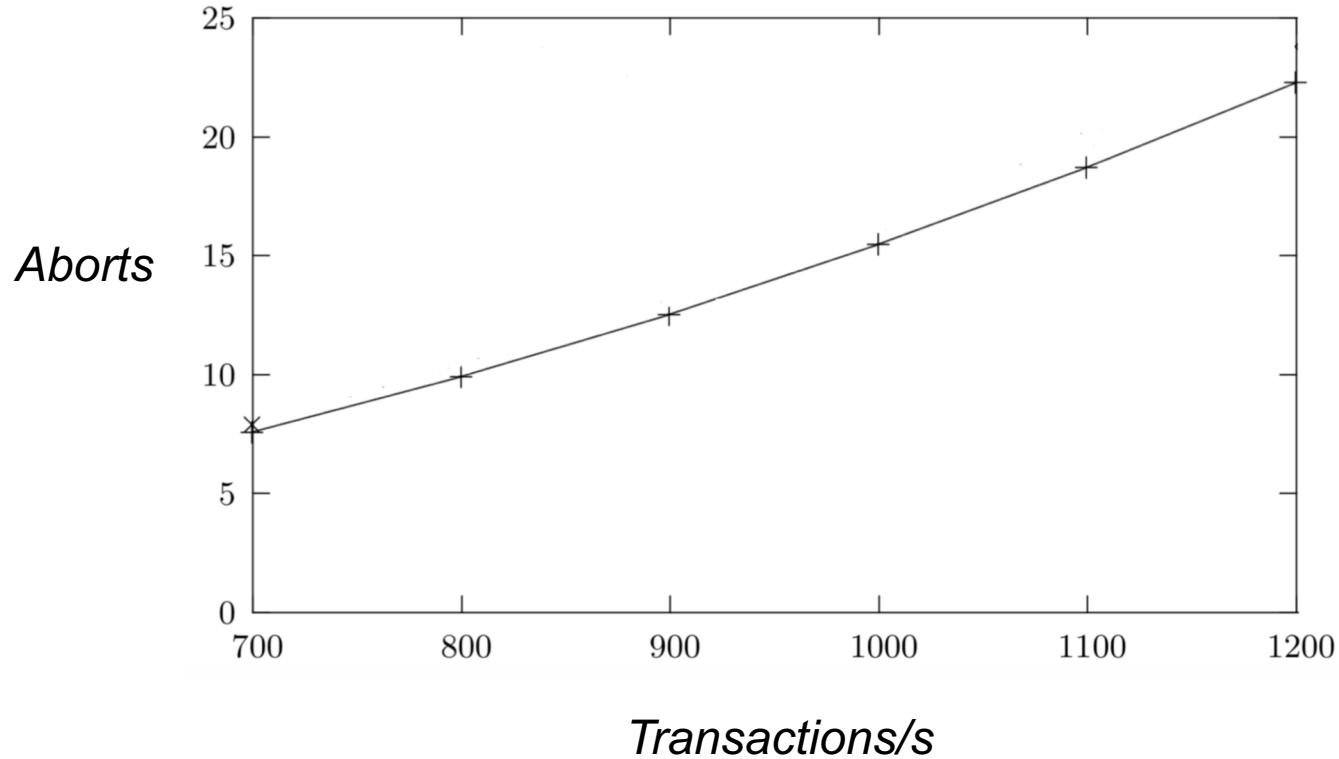


# Collision Detection Protocol: Example





# Collision Detection Protocol: Simulation Results





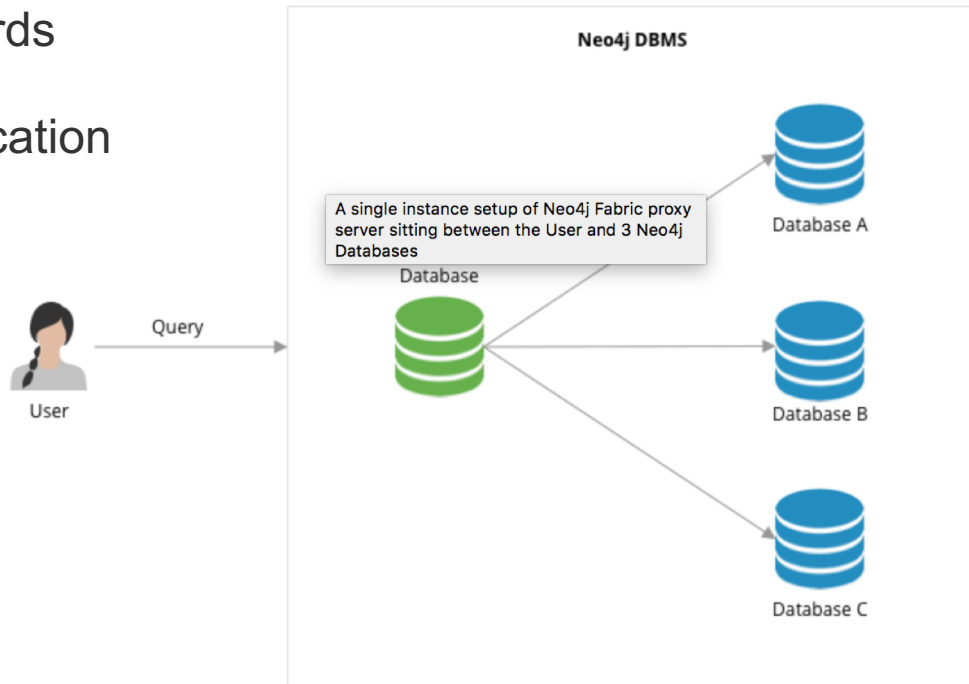
# Neo4j Recent Developments

- **Fabric** introduced in 4.0 (January 2020)
- Store and retrieve data in multiple databases using a single Cypher query
- **Data Sharding**: the ability to access data available in distributed sources in the form of a **common graph partitioned on multiple databases**.

# Neo4j Fabric

Queries cannot traverse across shards

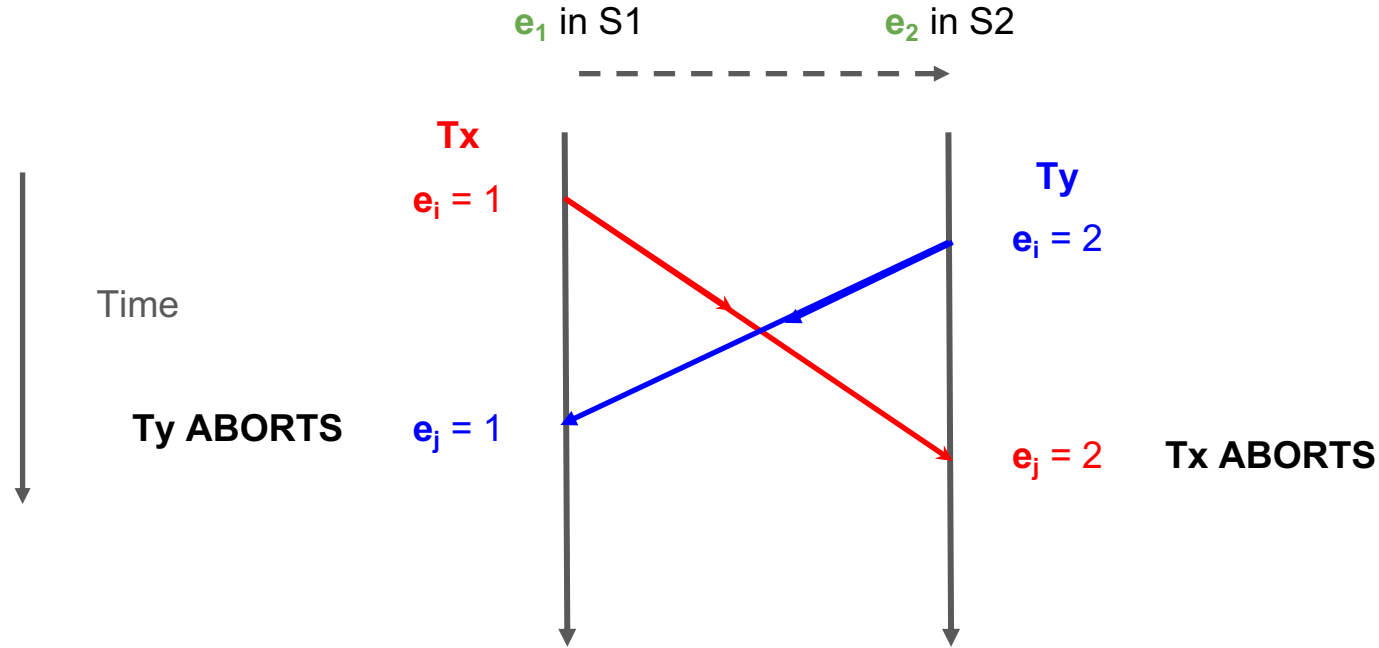
Needs a larger degree of data duplication



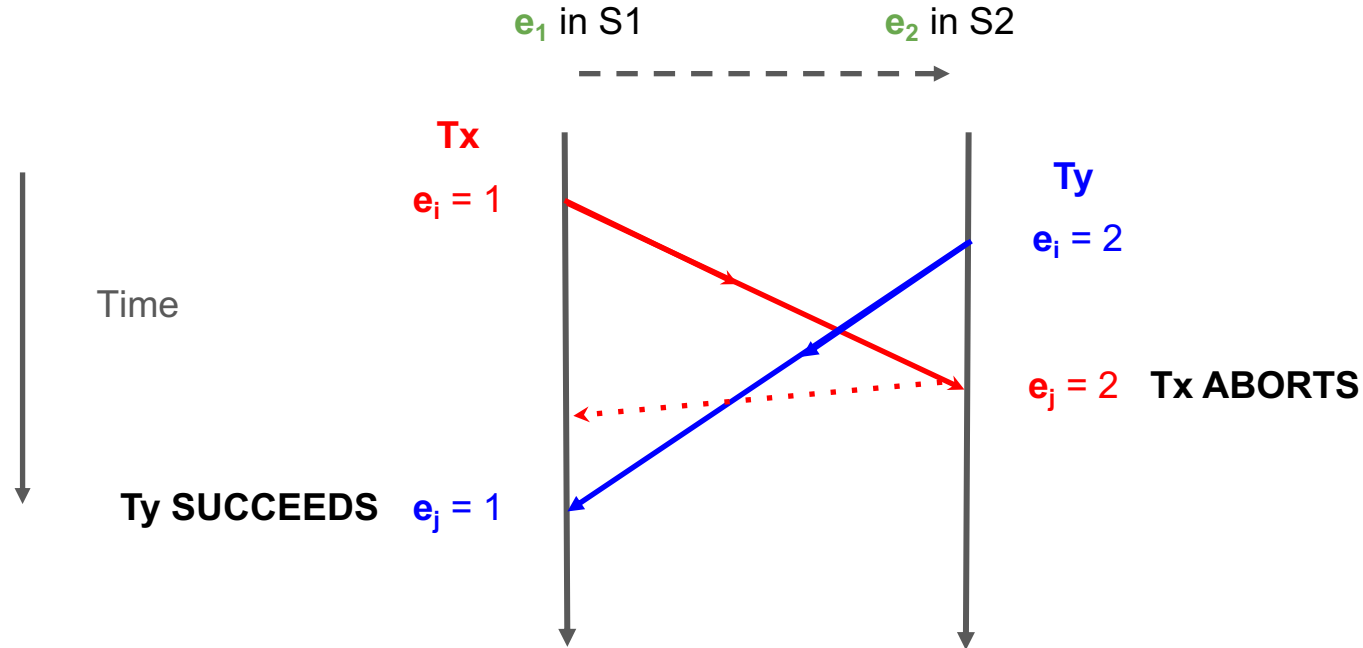
# Summary

1. Data corruption happens and **spreads quick**
2. Probabilistic solutions reduce the time until corruption
3. **Imperative reciprocal consistency is preserved**
4. Deterministic mechanism that preserves reciprocal consistency at all times
5. Neo4j Fabric supports sharding at the cost of higher data redundancy

# Collision Detection Protocol: Example 2



# Collision Detection Protocol: Example 3



# Neo4j Causal Clustering Architecture

- Default isolation: *Read Committed*
- Support for *Serializability* via explicit locking
- Clients get *read-your-own-writes* via bookmarking parameter for reads
- Replication semantics: *Causal Consistency*

