



# Distributed Transactions Without Atomic Clocks

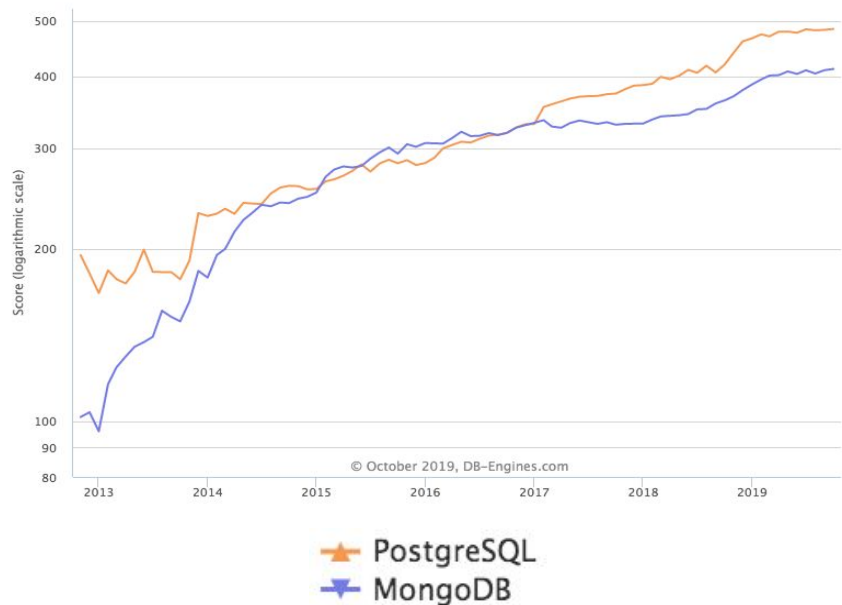
Sometimes, it's all just about *good timing*

*Karthik Ranganathan, co-founder & CTO*

# Introduction

# Designing the Perfect Distributed SQL Database

## Skyrocketing adoption of PostgreSQL for cloud-native applications



## Google Spanner

*The first horizontally scalable, strongly consistent, relational database service*

PostgreSQL is not highly available or horizontally scalable

Spanner does not have the RDBMS feature set

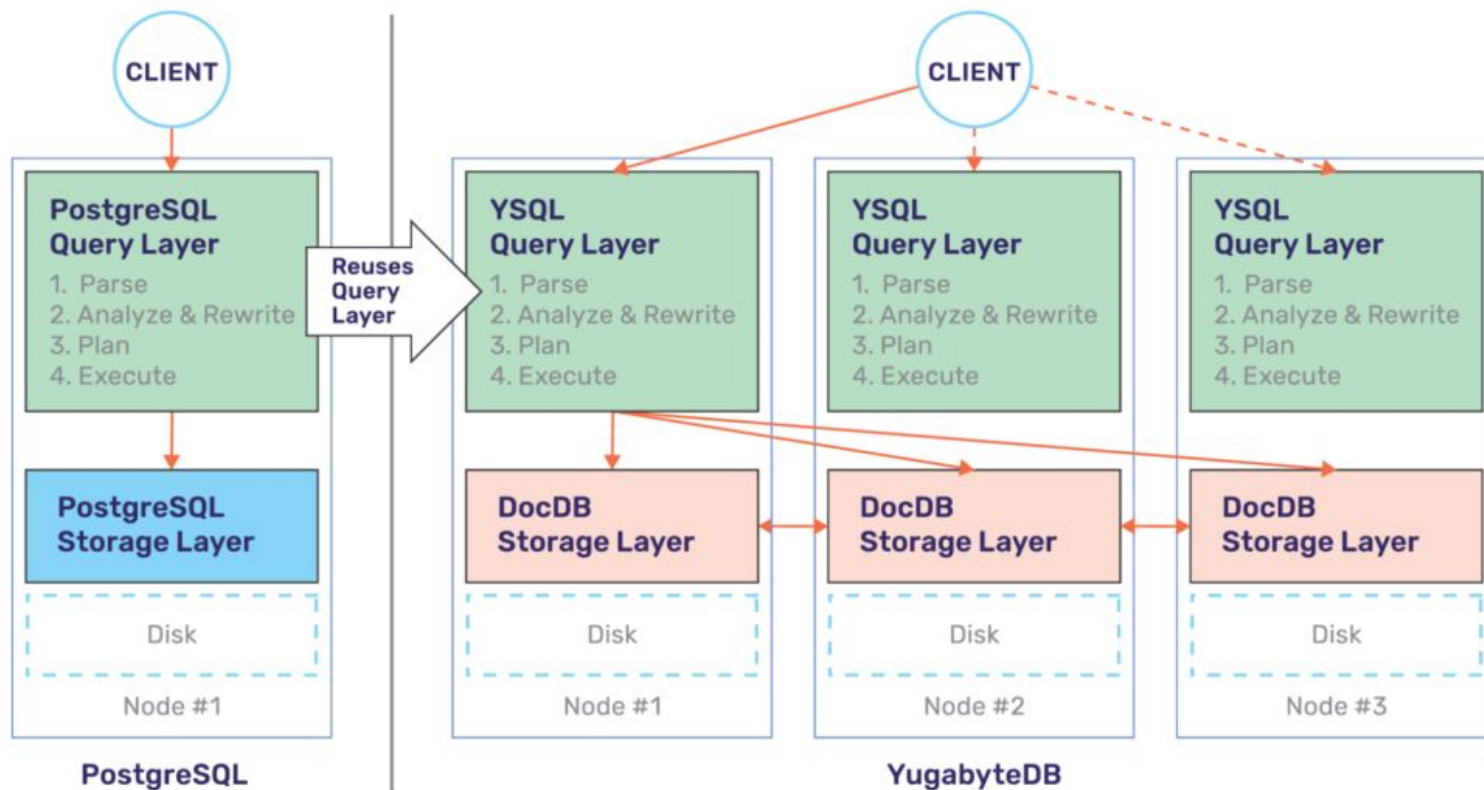
# Design Goals for YugabyteDB

**Transactional**, distributed SQL database designed for **resilience** and **scale**

- 100% open source
- PostgreSQL compatible
- Enterprise-grade RDBMS
  - Day 2 operational simplicity
  - Secure deployments
- Public, private, hybrid clouds
- High performance

	PostgreSQL	Google Spanner	YugabyteDB
SQL Ecosystem	✓ Massively adopted	✗ New SQL flavor	✓ Reuse PostgreSQL
RDBMS Features	✓ Advanced Complex	✗ Basic cloud-native	✓ Advanced Complex and cloud-native
Highly Available	✗	✓	✓
Horizontal Scale	✗	✓	✓
Distributed Txns	✗	✓	✓
Data Replication	Async	Sync	Sync + Async

# YugabyteDB Reuses PostgreSQL Query Layer



Transactions are fundamental to SQL...

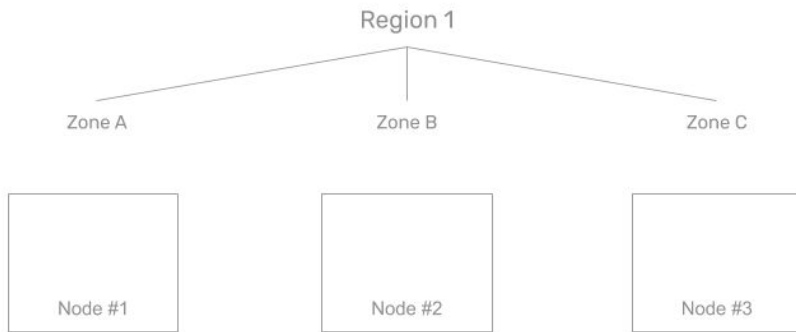
But they require time synchronization between nodes.

**Why?**

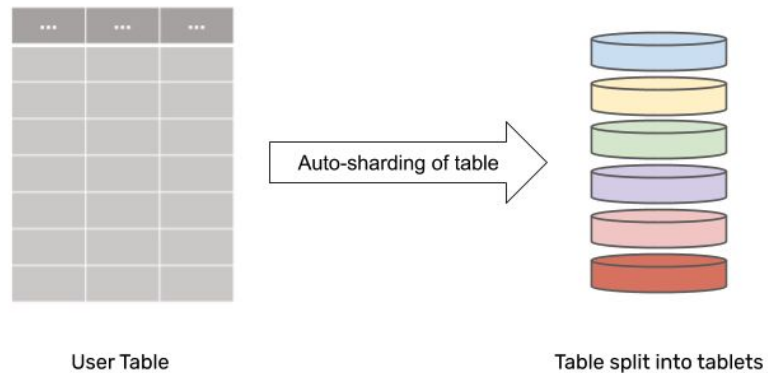
Let's look at single-row transactions before answering this

# Single-Row Transactions: Raft Consensus

# Distributing Data For Horizontal Scalability



- Assume 3-nodes across zones
- How to distribute data across nodes?



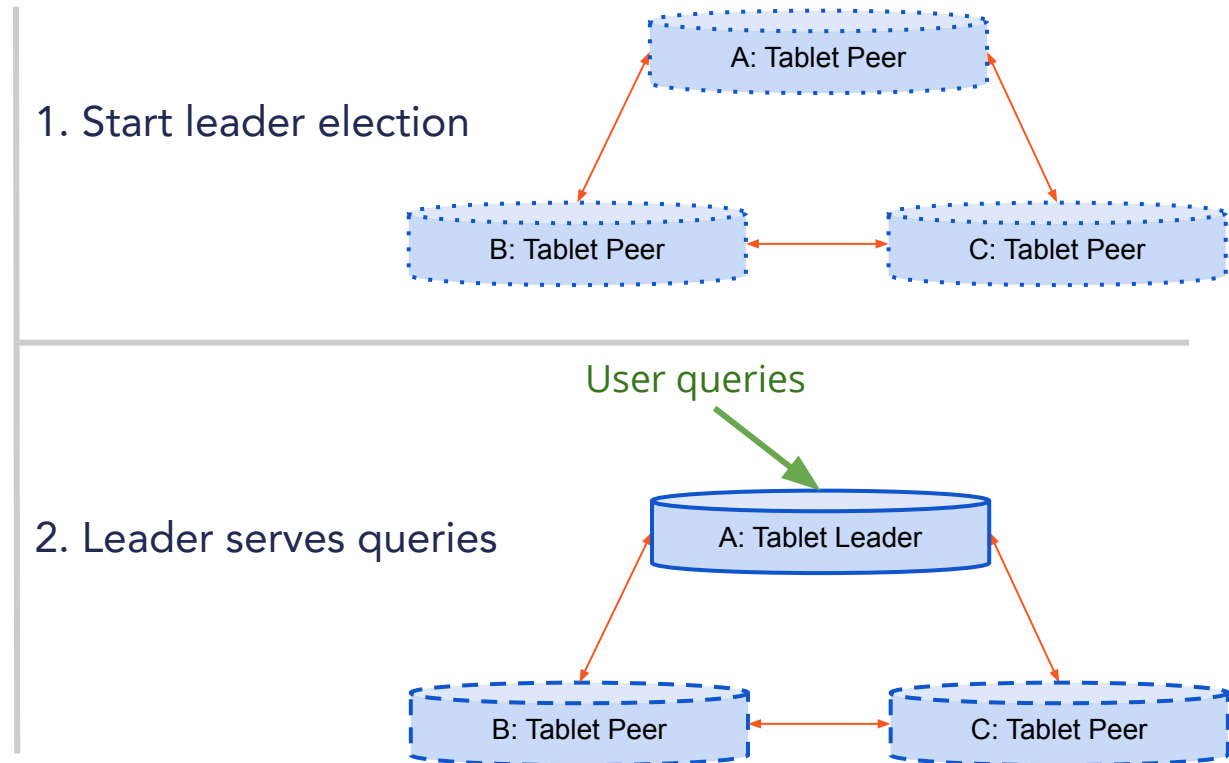
- User tables sharded into tablets
- Tablet = group of rows
- Sharding is transparent to user



# Tablets Use Raft-Based Replication



*Raft Algorithm* for replicating data: per-row linearizability



# Replication in a 3 Node Cluster

- Assume  $rf = 3$
- Survives 1 node or zone failure
- Tablets replicated across 3 nodes
- Follower (replica) tablets balanced across nodes in cluster

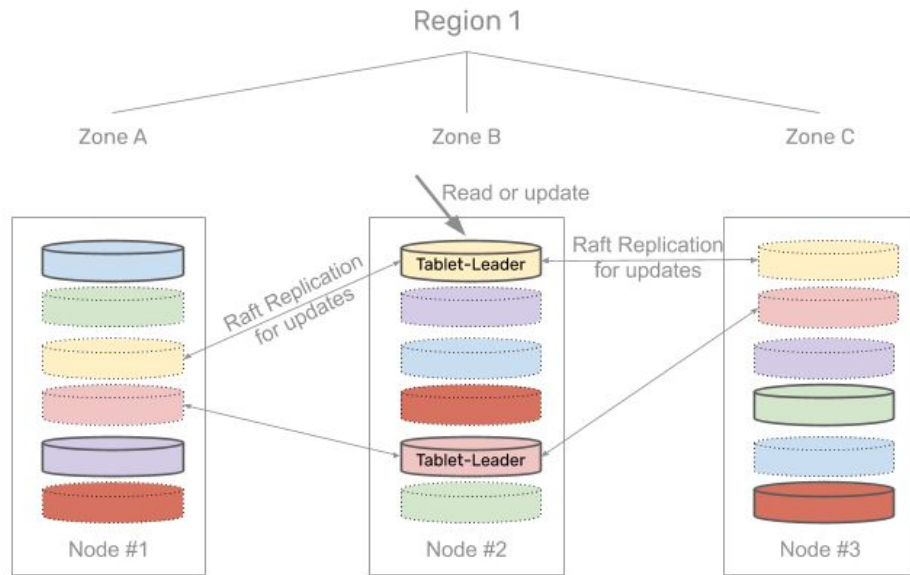


Diagram with replication factor = 3

# No Time Synchronization Needed So Far

YugabyteDB Raft replication based on:

`leader leases + time intervals + CLOCK_MONOTONIC`

From Jepsen Testing Report:

Whatever the case, this is a good thing for operators: nobody wants to worry about clock safety unless they have to, and YugaByte DB appears to be mostly robust to clock skew. Keep in mind that we cannot (rigorously) test YugaByte DB's use of CLOCK\_MONOTONIC\_RAW for leader leases, but we suspect skew there is less of an issue than CLOCK\_REALTIME synchronization.

# Need for Time Synchronization: Distributed Transactions

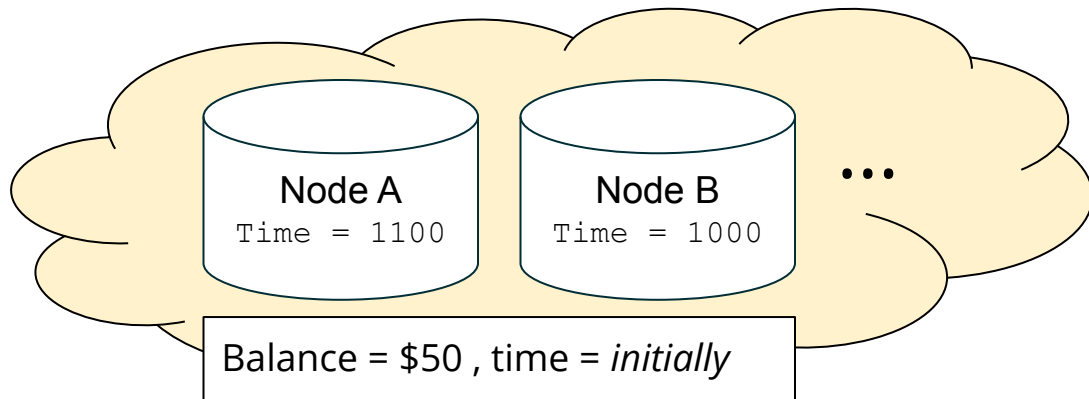
# Let's Take a Simple Example Scenario

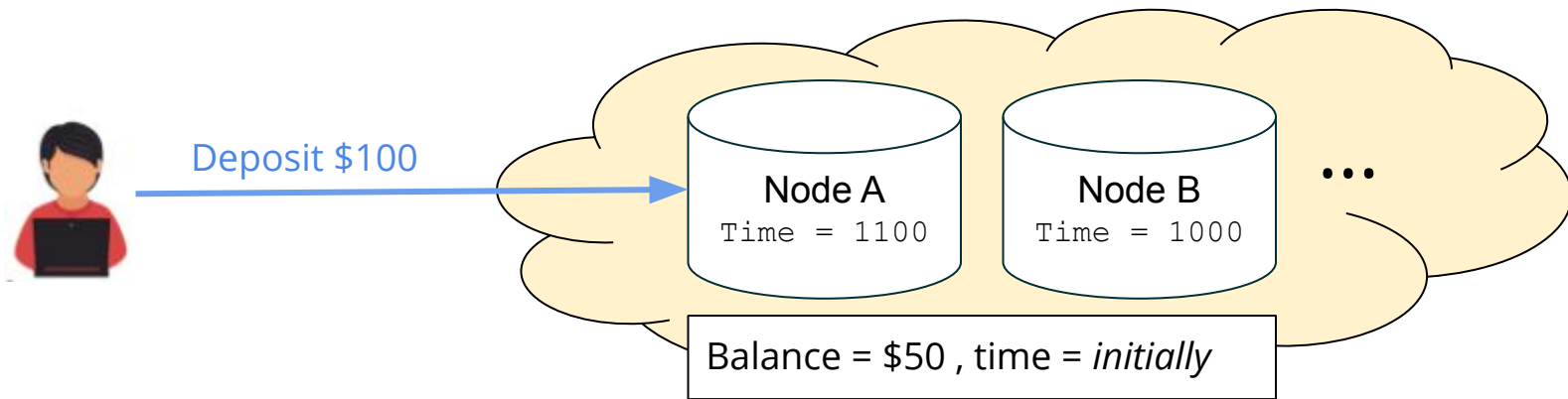
## User:

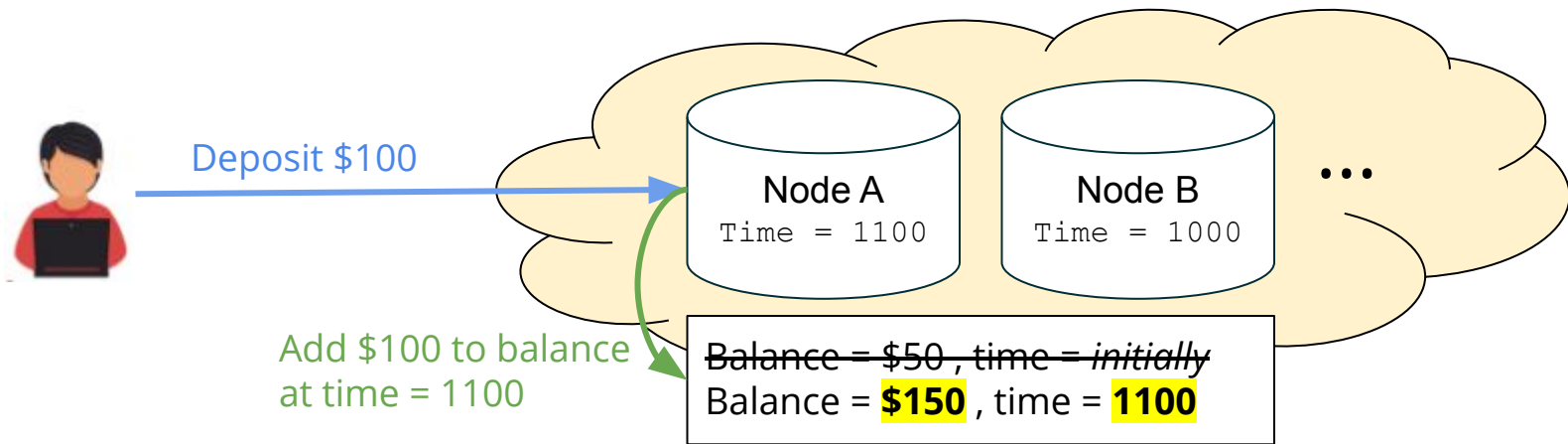
- *Deposit followed by withdrawal*
- A user has \$50 in bank account
- Deposit \$100, new balance = \$150
- Withdraw \$70, should be ok always

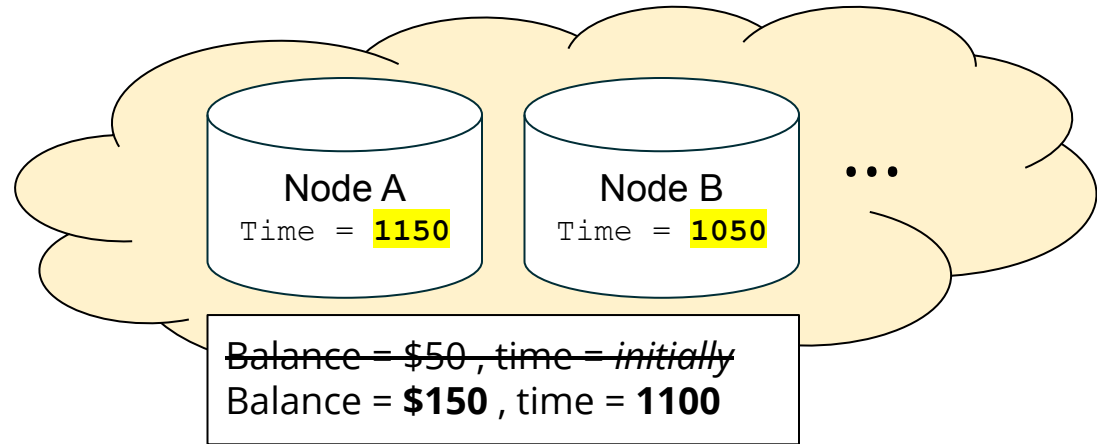
## Cluster:

- *Nodes have clock skew*
- Node A is 100ms ahead of node B
- Deposit on Node A
- Withdraw from B, 50ms after deposit

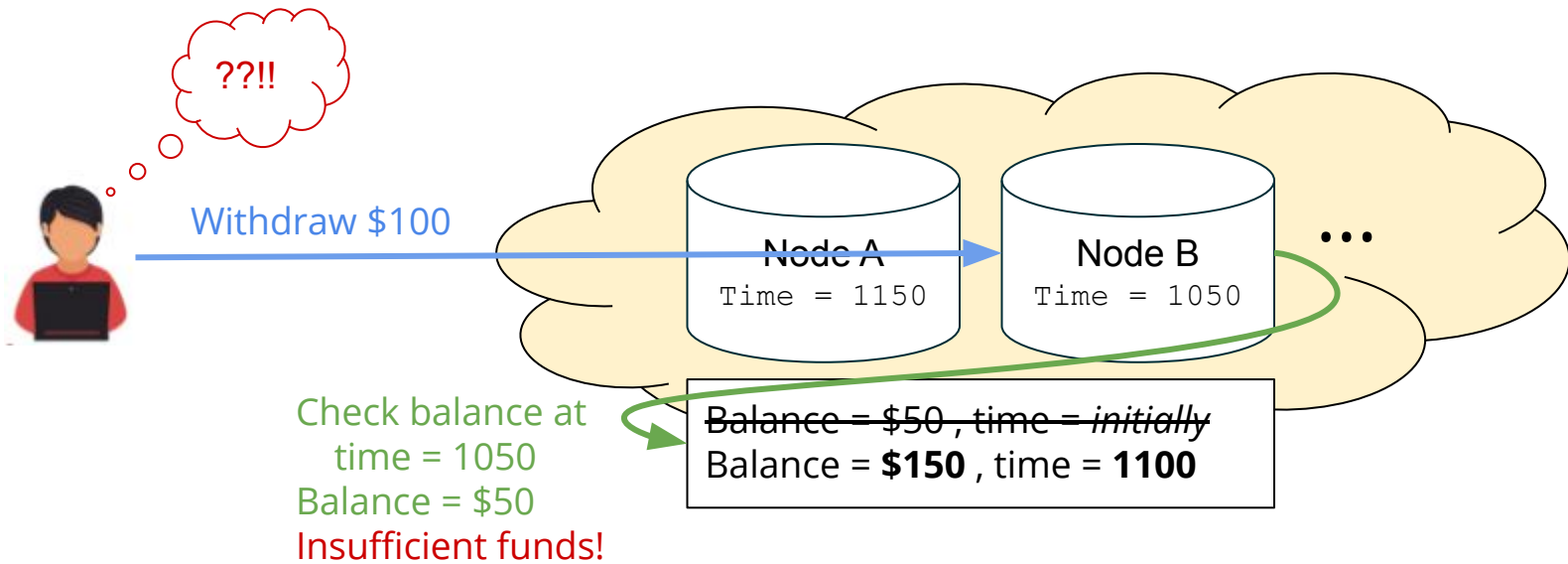












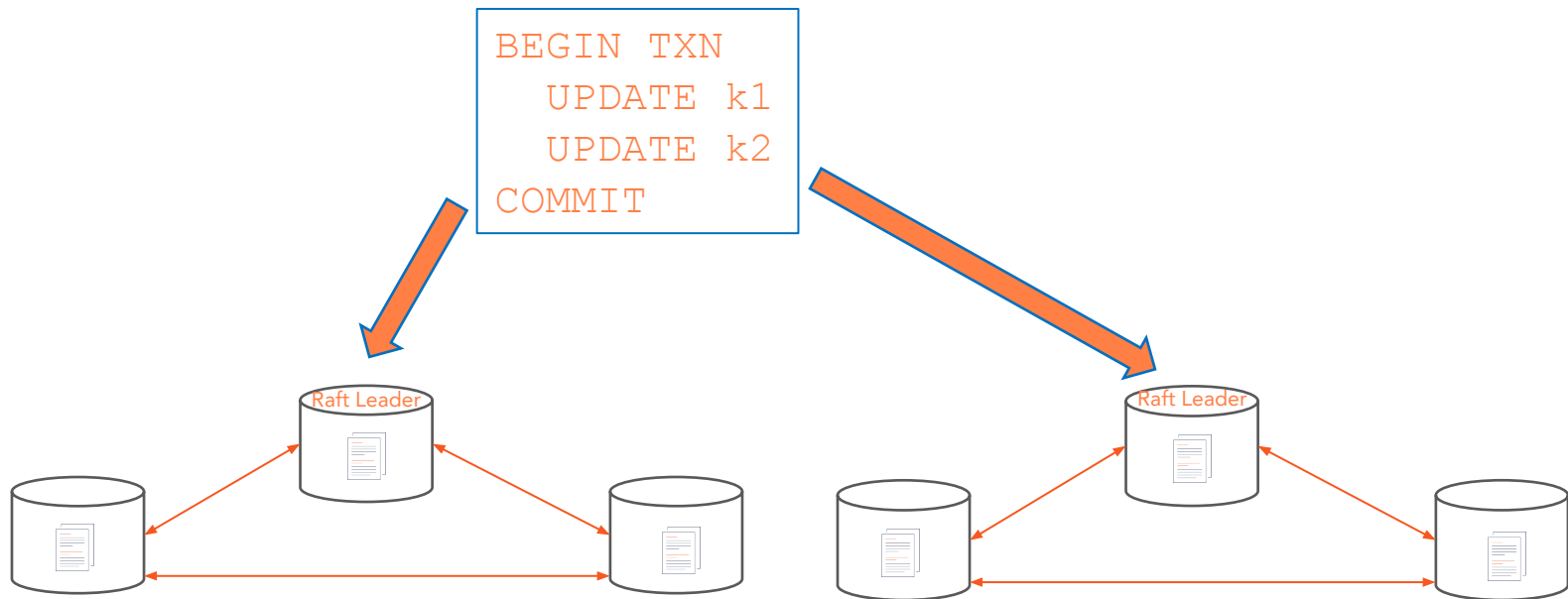
# Time Sync Needed For Distributed Txns Also

```
BEGIN TXN
  UPDATE k1
  UPDATE k2
COMMIT
```

k1 and k2 may belong to **different shards**

Belong to **different Raft groups** on completely **different nodes**

# What Do Distributed Transactions Need?



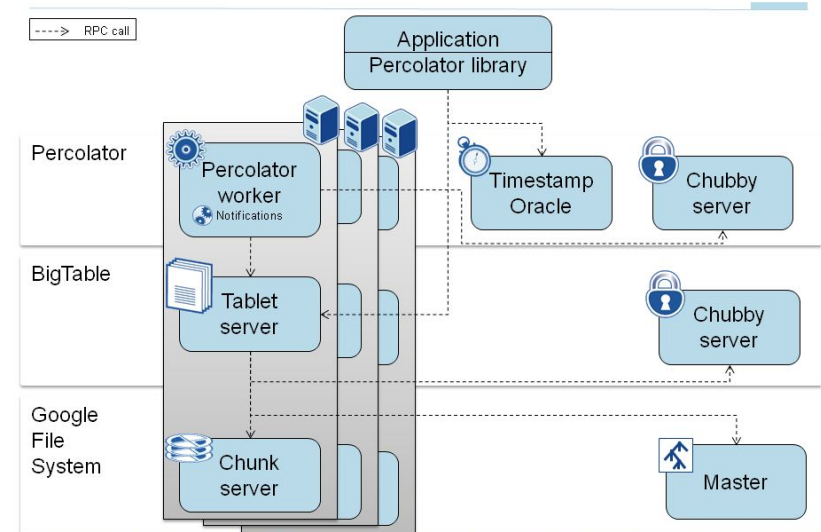
Updates should get written at the **same time**

But how will **nodes agree on time?**

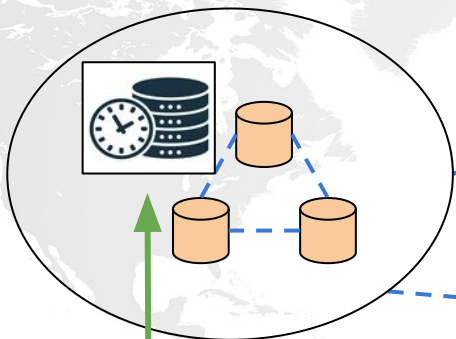
# Time Synchronization: Timestamp Oracle vs Distributed Time Sync

# Timestamp Oracle - Google Percolator, Apache Omid

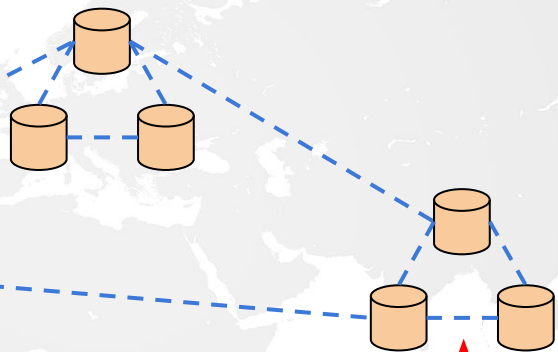
- Not scalable
  - Bottleneck in system
- Poor multi-region deployments
  - High latency
  - Low availability
- Prediction
  - Clock sync will improve esp in public clouds over time



Sources: <http://research.google.com/pubs/pub36726.html>, <http://labs.google.com/papers/bigtable.html>, <http://labs.google.com/papers/gfs.html>



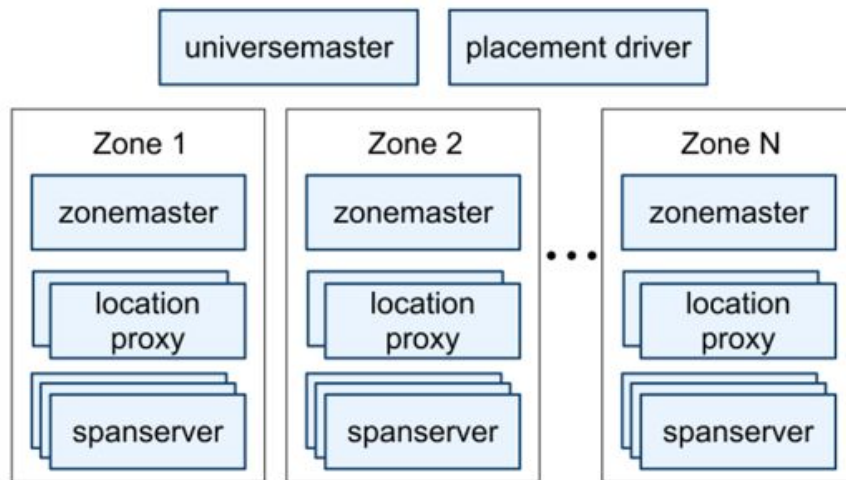
**1. Timestamp Oracle** gets partitioned away from rest of the cluster



**2. Remote transactions** fail without connectivity to the Timestamp Oracle

# Distributed Time Sync - Google Spanner

- Scalable
  - Only nodes involved in a txn need to coordinate
- Multi-region deployments
  - Distributed, region-local time synchronization
- Based on 2-phase commit
- Uses hybrid logical clocks



# Distributed Time Sync Using GPS/Atomic Clock Service



# Atomic Clock Service



Atomic Clock Based Time Service:  
highly available, globally synchronized clocks, tight error bounds

Not a commodity service

Most of physical clocks are **very poorly synchronized**:  
`ntp` has clock skew of **100ms - 250ms**

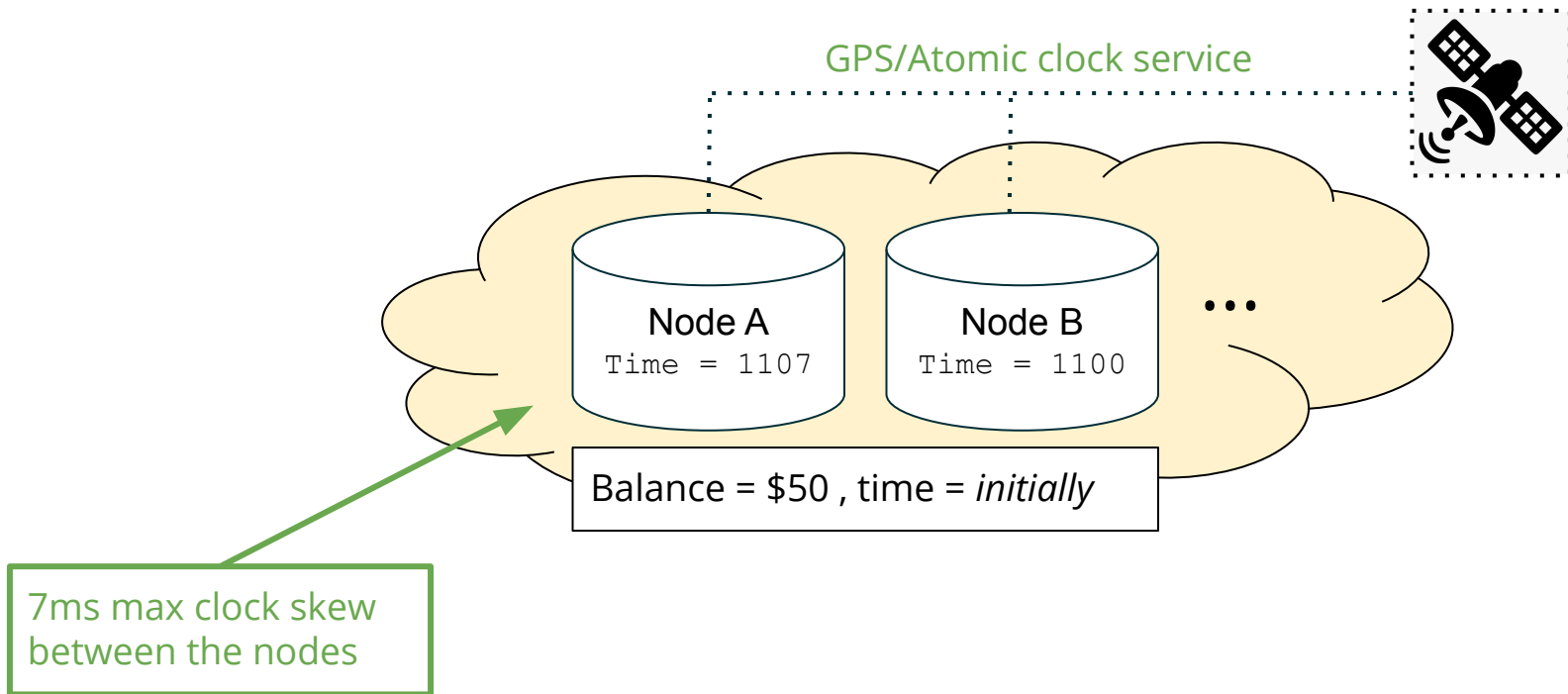
# How Does an Atomic Clock Service Help?

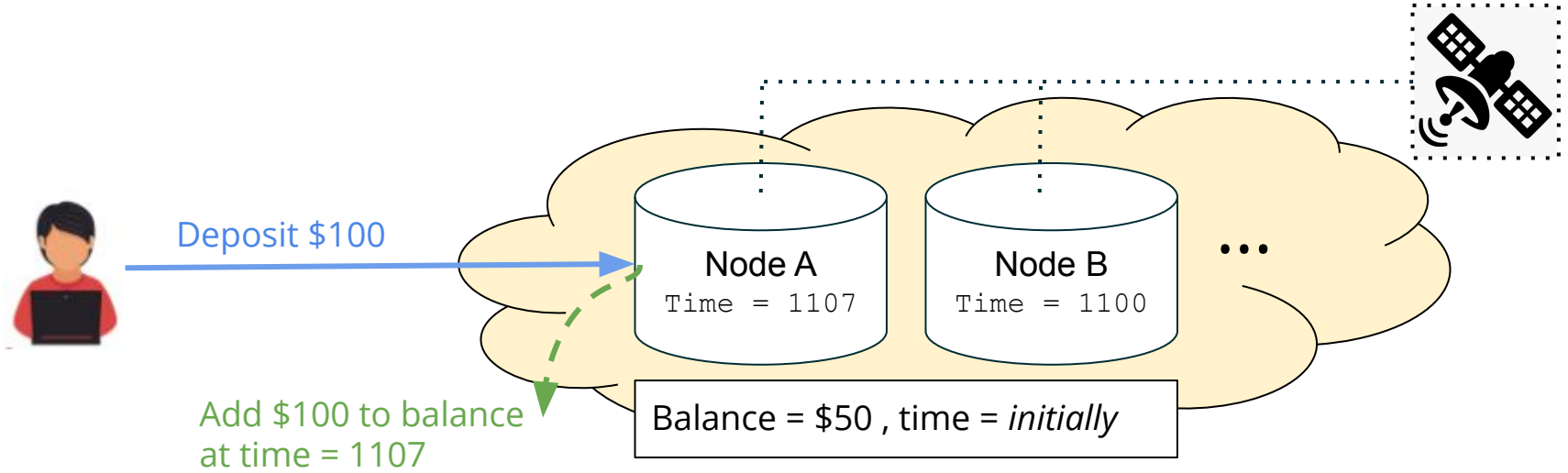
Guarantees an upper bound on clock skew between nodes.

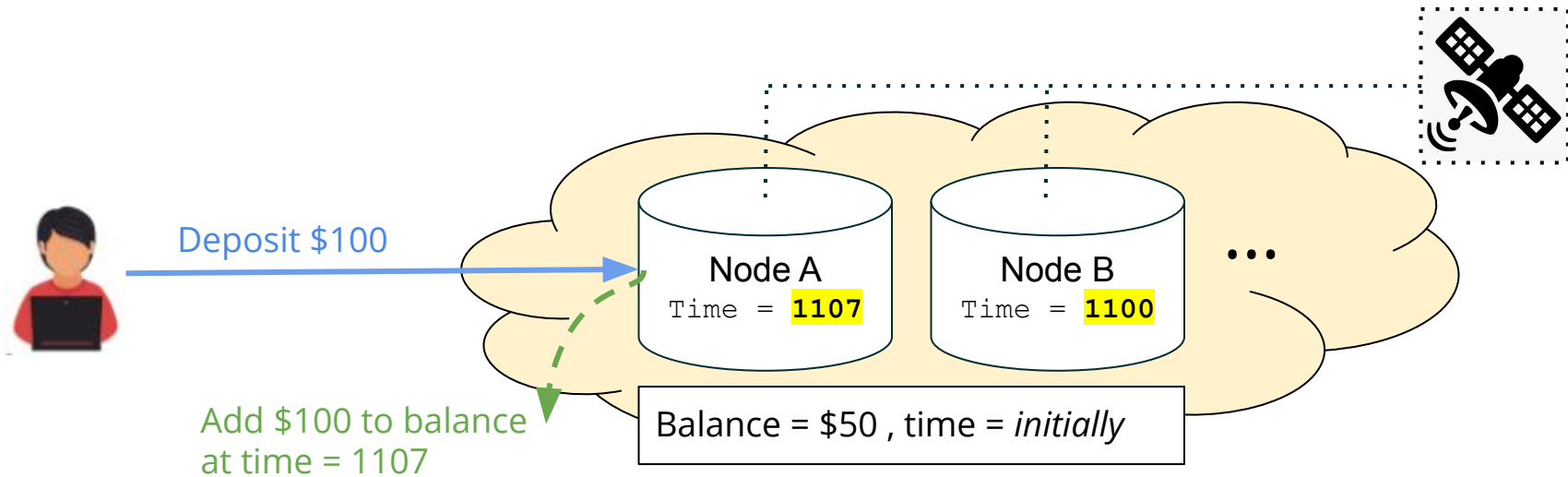
TrueTime service used by Google Spanner: 7ms max skew

Let's try that scenario again with an atomic clock service

# Let's Take an Example Scenario



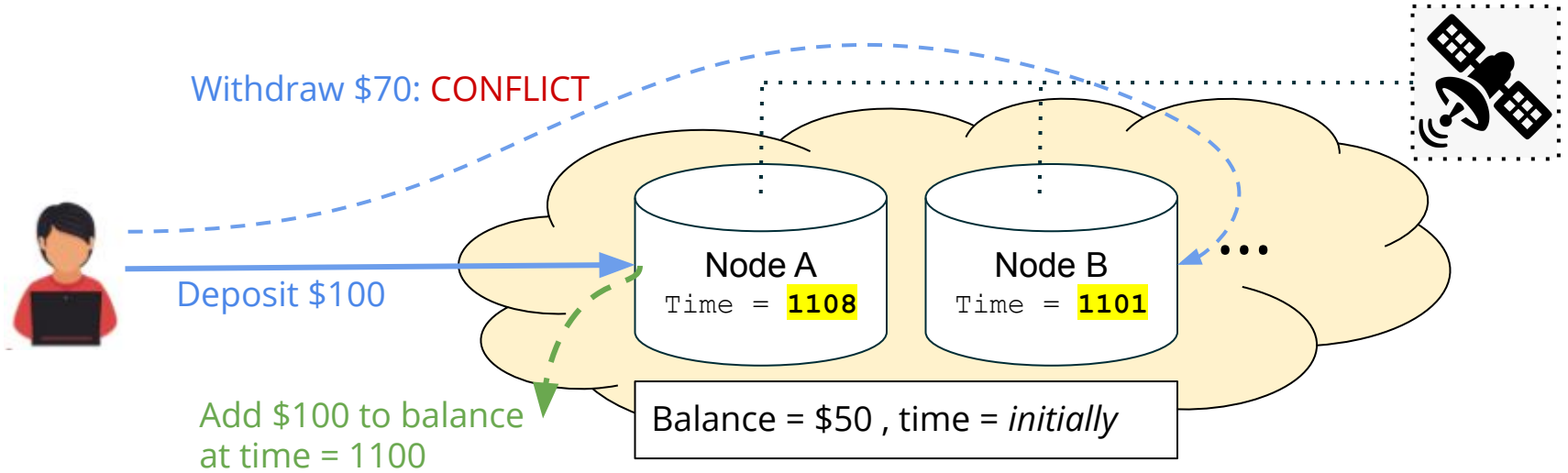


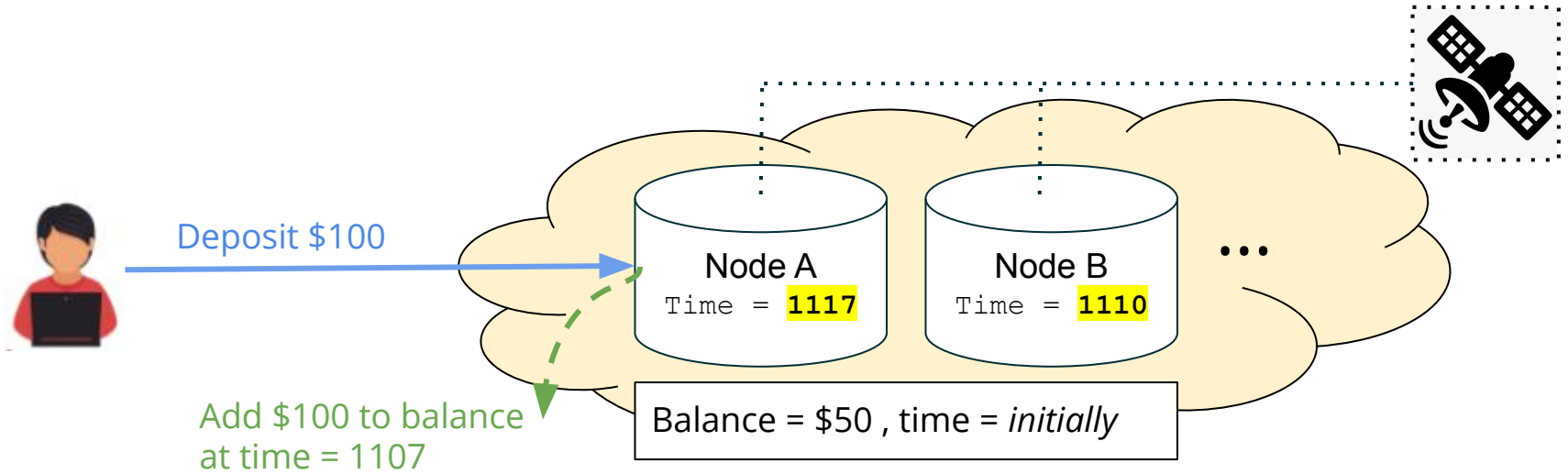


Introduce  $> 7\text{ms}$  `commit_wait` delay - allows all clocks to catch up

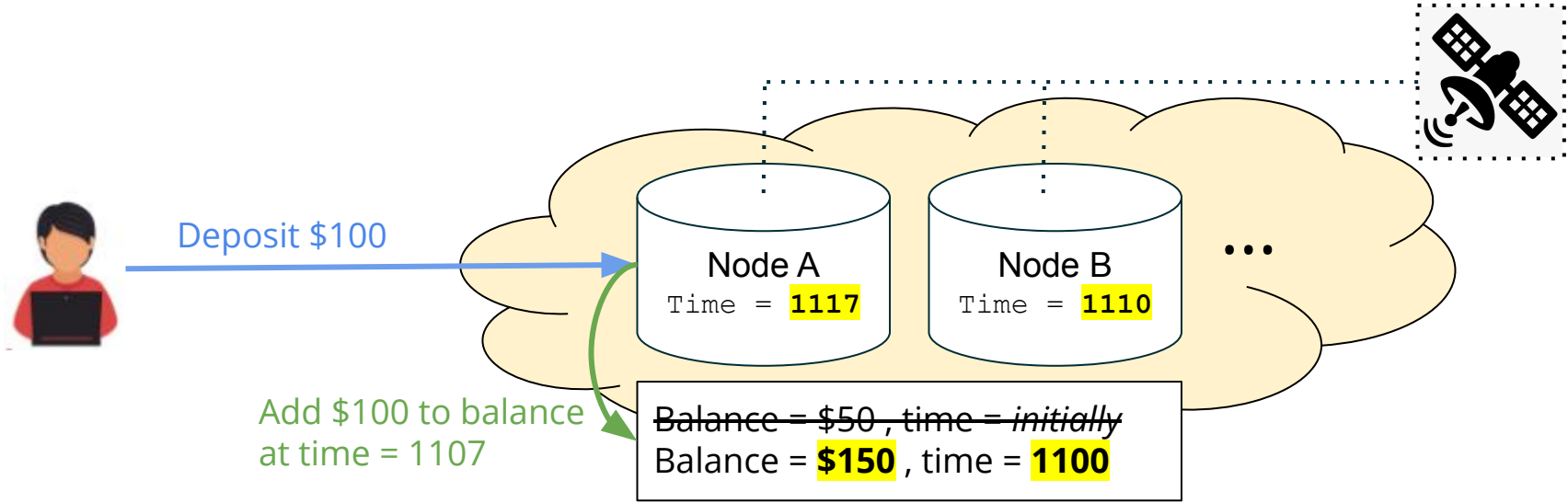


## Transaction conflicts detected and handled appropriately





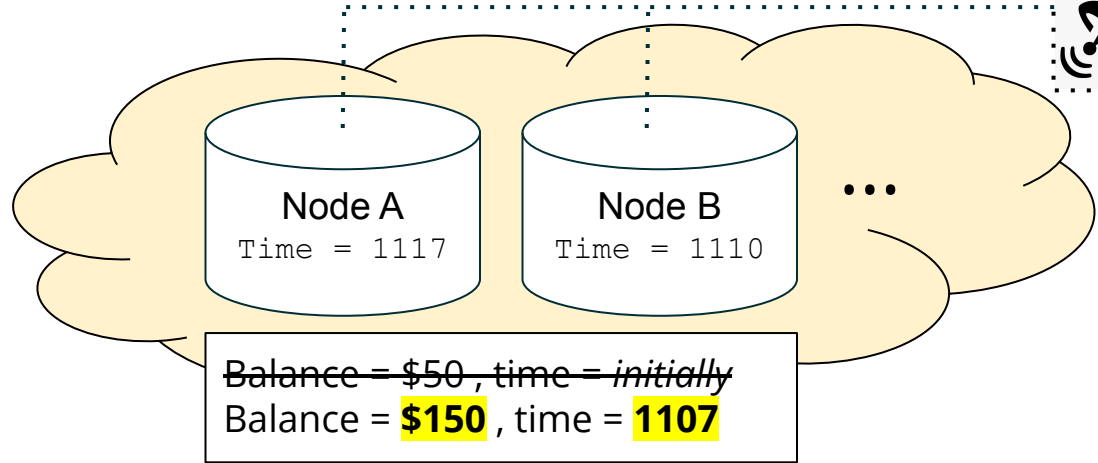
All clocks are guaranteed to have caught up to commit time 1107







All transactions will now occur after time **1110** and hence are safe



# Doing This Without Atomic Clocks: Hybrid Logical Clocks

# Hybrid Logical Clock (HLC)

Combine coarsely-synchronized **physical clocks** with **Lamport Clocks** to track causal relationships

Hybrid Logical Clock = (physical component, logical component)

synchronized using NTP

a monotonic counter

Nodes update HLC on each **Raft exchange** for things like **heartbeats**, **leader election** and **data replication**

← 64-bit HLC value →

(physical component , logical component)

Physical timestamp

- 52-bit value
- Microsecond precision
- Synchronized using `ntp`

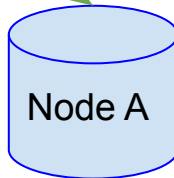
Logical timestamp

- 12-bit value
- 4096 values max
- Vector clock component

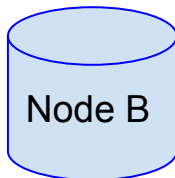


UTC Time: T0

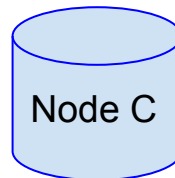
1. Update request



```
Wall clock: 1200  
HLC time  : (1200, 10)
```



```
Wall clock: 1000  
HLC time  : (1000, 30)  
⚠️ Skew : 200ms
```

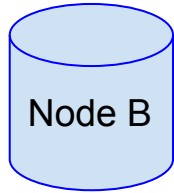
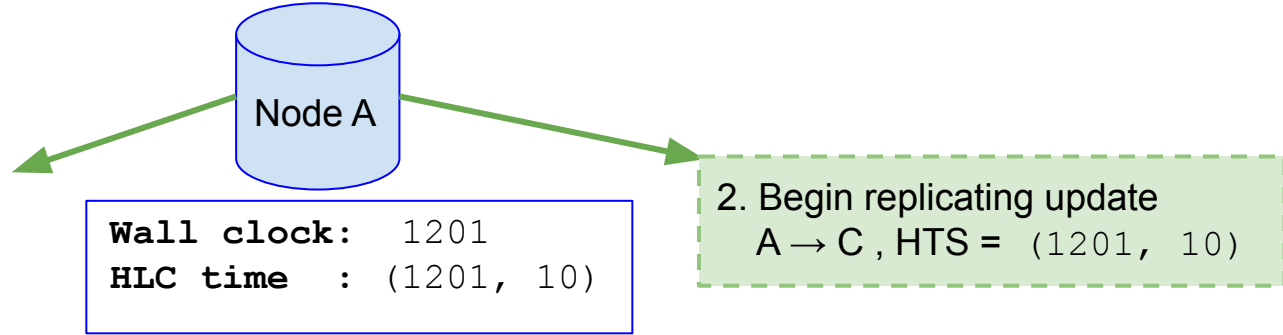


```
Wall clock: 1100  
HLC time  : (1100, 20)  
⚠️ Skew : 100ms
```

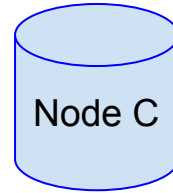


# UTC Time: T0 + 1

Time has progressed on all nodes by 1 second



Wall clock: 1001  
HLC time : (1001, 30)  
⚠️🕒 Skew : 200ms



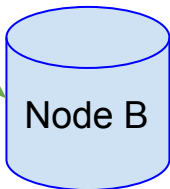
Wall clock: 1101  
HLC time : (1101, 20)  
⚠️🕒 Skew : 100ms



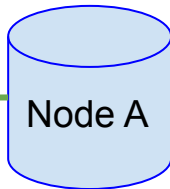
# UTC Time: T0 + 2

Time has progressed on all nodes by 2 second

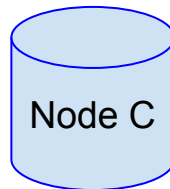
3. Finish replication A → B after 1ms, HTS = (1201, 10)



Wall clock: 1002  
HLC time : ~~(1002, 30)~~ (1201, 31)  
⚠️ Skew : ~~200ms~~ 1ms



Wall clock: 1202  
HLC time : (1202, 10)



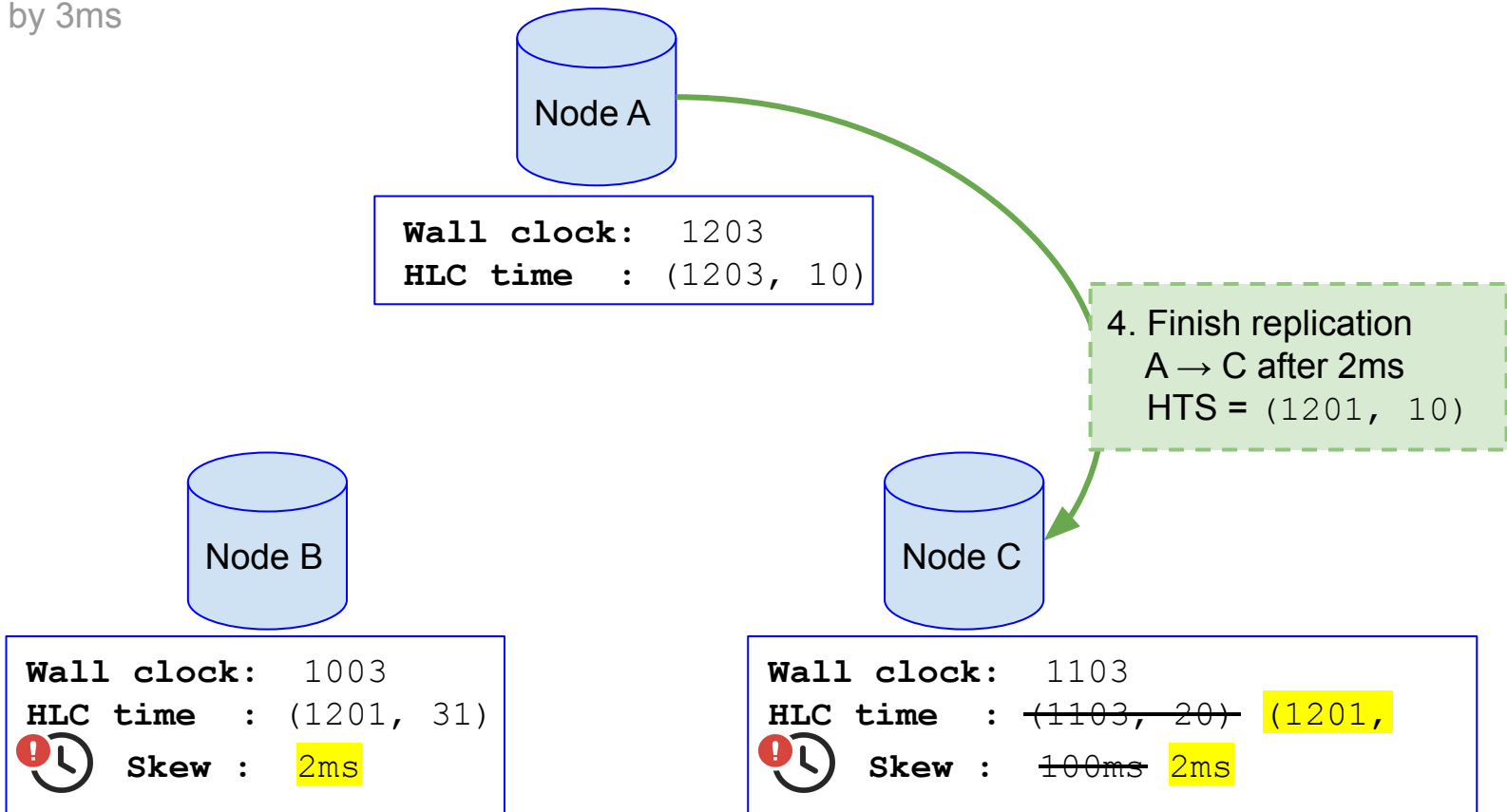
3. Replication still in flight A → C after 1ms, HTS = (1201, 10)

Wall clock: 1102  
HLC time : (1101, 20)  
⚠️ Skew : 100ms



# UTC Time: T0 + 3

Time has progressed on all nodes by 3ms



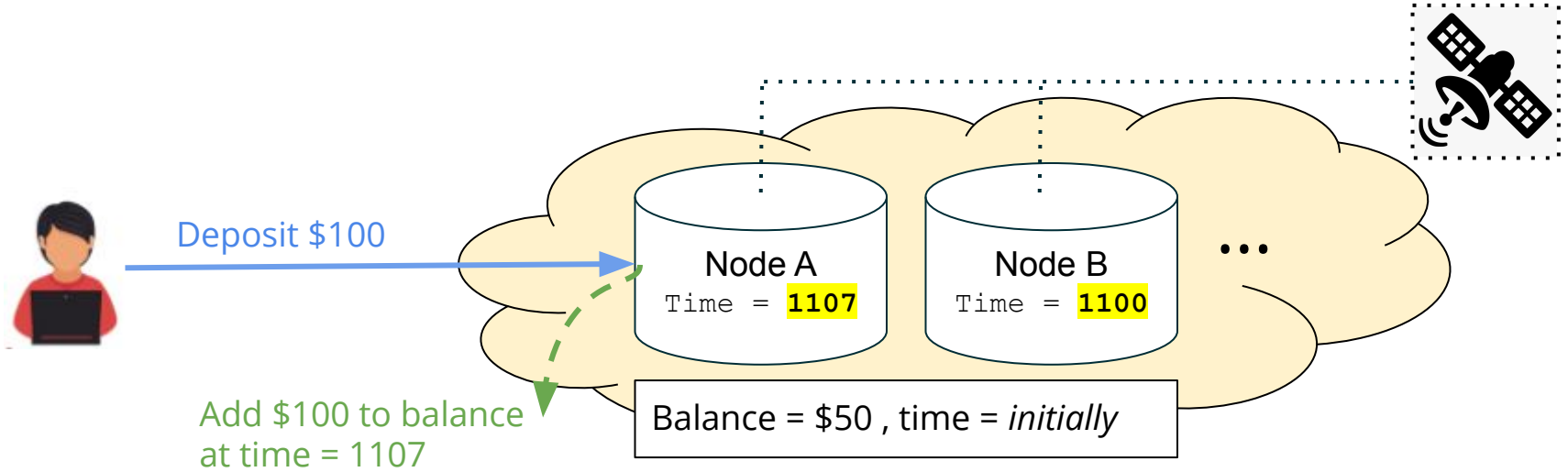


# Uncertainty Time-Windows and Max Clock Skew

Typically, **frequent RPCs between nodes** forces quick time synchronization in a cluster

Sometimes, this may not be the case. In such cases, DB **transparently detects conflicting transactions** and **retries them if possible** to do so.

# Recall This Scenario We Discussed Before

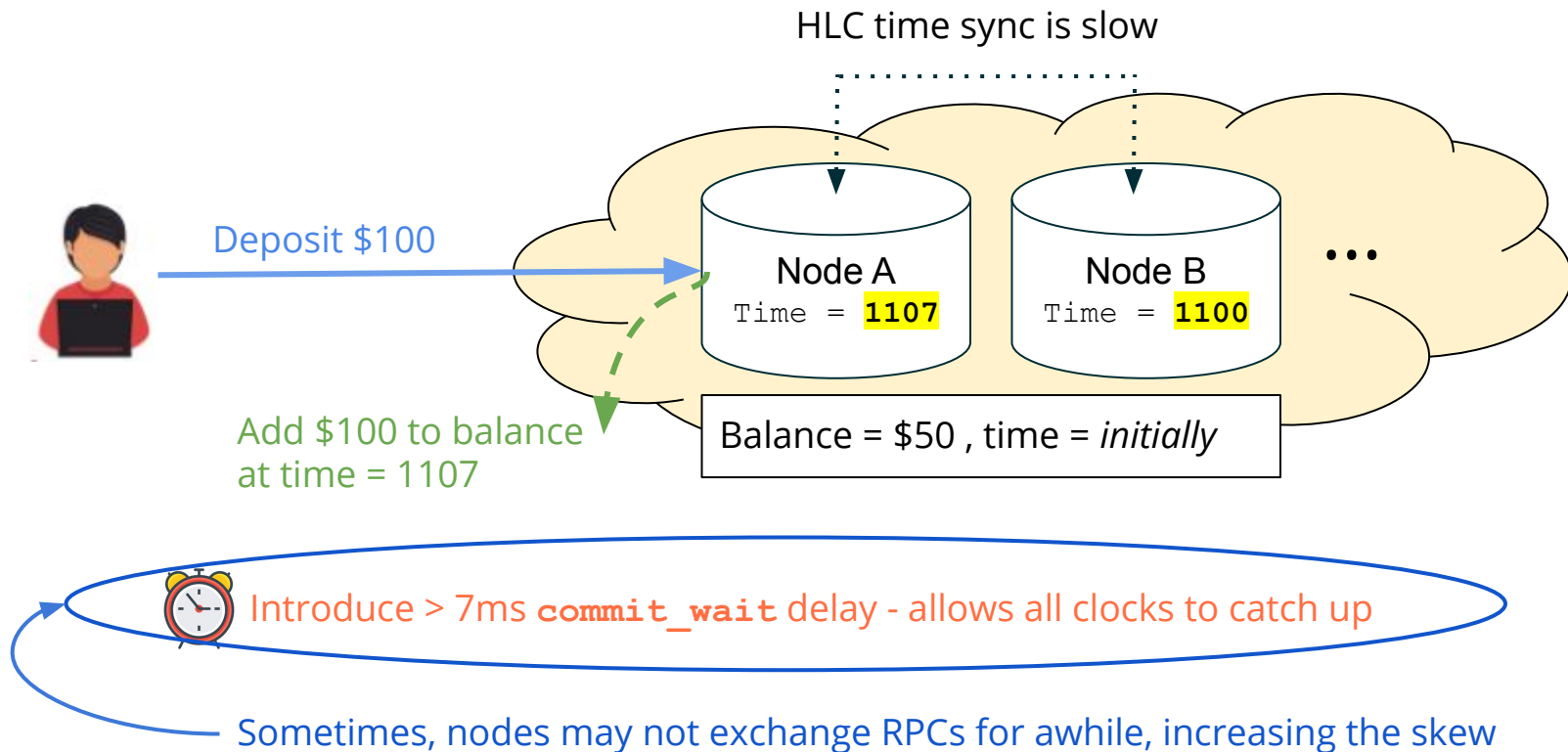


Add \$100 to balance  
at time = 1107



Introduce  $> 7\text{ms}$  `commit_wait` delay - allows all clocks to catch up

# Knowing the Max Skew Is Important with HLCs



Rare scenarios, but safety first!

Recommendation: set `max_skew` to large enough value

```
max_skew = 500ms # suggested default
```

# Integrating Raft with Hybrid Logical Clocks

# Raft vs Hybrid Logical Clock (HLC)

## Raft Consensus:

- Per tablet
- Issues Raft Sequence ID, which is a purely logical sequence number
- Single row transactions
- Raft sequence id is monotonically increasing

## Cluster:

- Per node
- Issues HLC timestamp, which can be compared across nodes
- Distributed transactions
- HLC timestamp is monotonically increasing

# Raft vs Hybrid Logical Clock (HLC)

## Raft Consensus:

- Per tablet
- Issues Raft Sequence ID, which is a purely logical sequence number
- Single row transactions
- Raft sequence id is monotonically increasing

## Cluster:

- Per node
- Issues HLC timestamp, which can be compared across nodes
- Distributed transactions
- HLC timestamp is monotonically increasing

Monotonicity allows these two ids to be correlated





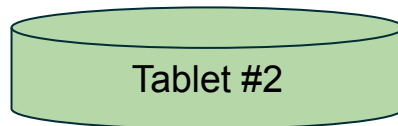
*Raft Log for tablet #1*

Txn1: Raft-Id=100,  
HLC=(1000, 32),  
<txn 1 details>

Txn5: Raft-Id=101,  
HLC=(1005, 21),  
<txn 5 details>

Txn7: Raft-Id=102,  
HLC=(1100, 75),  
<txn 7 details>

...



*Raft Log for tablet #2*

Txn2: Raft-Id=1200,  
HLC=(1004, 14),  
<txn 2 details>

Txn7: Raft-Id=1201,  
HLC=(1100, 75),  
<txn 7 details>

Txn8: Raft-Id=1202,  
HLC=(1105, 91),  
<txn 8 details>

...

Tablet #1

Tablet #2

**Raft Ids** are issued per-tablet and are not correlated across tablets.

**HLCs** are correlated across tablets because of distributed time sync

*Raft Log for tablet #1*

Txn1: Raft-Id=100,  
HLC=(1000, 32),  
<txn 1 details>

Txn5: Raft-Id=101,  
HLC=(1005, 21),  
<txn 5 details>

Txn7: **Raft-Id=102,**  
**HLC=(1100, 75),**  
<txn 7 details>

...

*Raft Log for tablet #2*

Txn2: Raft-Id=1200,  
HLC=(1004, 14),  
<txn 2 details>

Txn7: **Raft-Id=1201,**  
**HLC=(1100, 75),**  
<txn 7 details>

Txn8: Raft-Id=1202,  
HLC=(1100, 32),  
<txn 8 details>

...

# DocDB: Distributed Transactions

# Fully Decentralized Architecture

- **No single point of failure or bottleneck**
  - Any node can act as a Transaction Manager
- **Transaction status table distributed across multiple nodes**
  - Tracks state of active transactions
- **Transactions have 3 states**
  - Pending
  - Committed
  - Aborted
- **Reads served only for Committed Transactions**
  - Clients never see inconsistent data

# Isolation Levels

- **Serializable Isolation**

- Read-write conflicts get auto-detected
- Both reads and writes in read-write txns need provisional records
- Maps to SERIALIZABLE in PostgreSQL

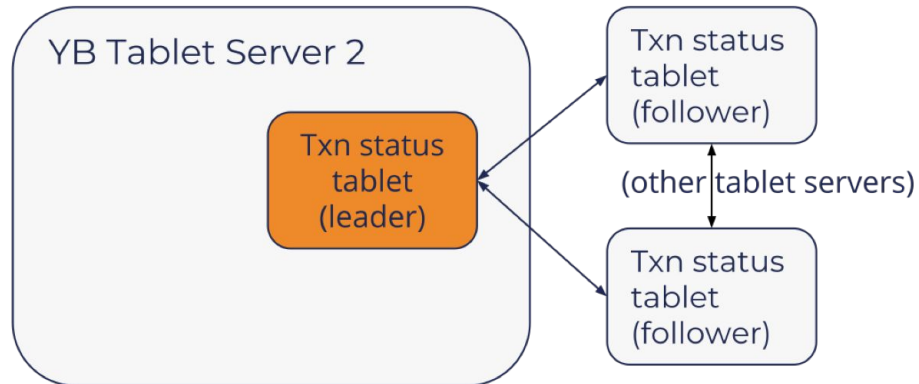
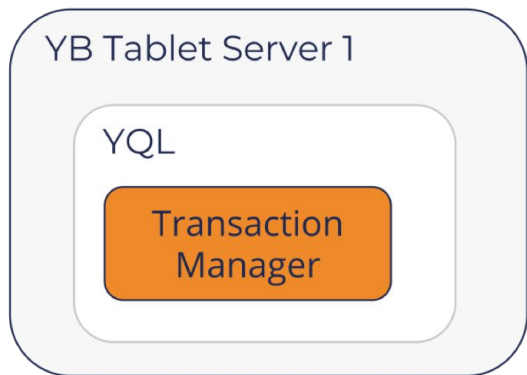
- **Snapshot Isolation**

- Write-write conflicts get auto-detected
- Only writes in read-write txns need provisional records
- Maps to REPEATABLE READ, READ COMMITTED & READ UNCOMMITTED in PostgreSQL

- **Read-only Transactions**

- Lock free

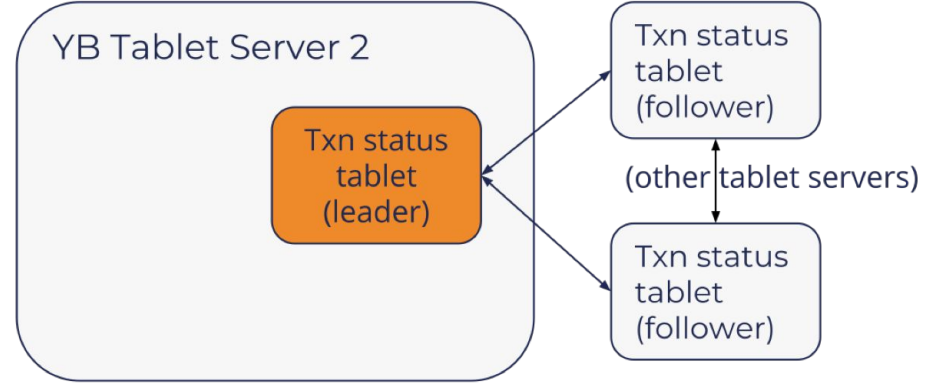
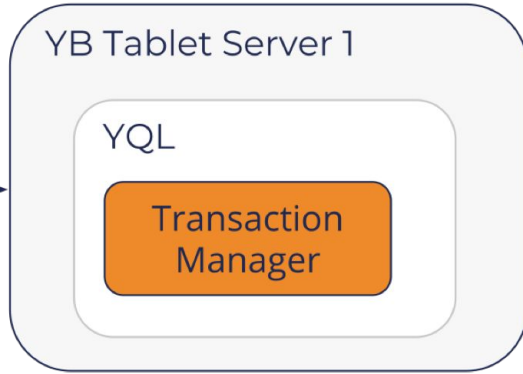
# Distributed Transactions - Write Path



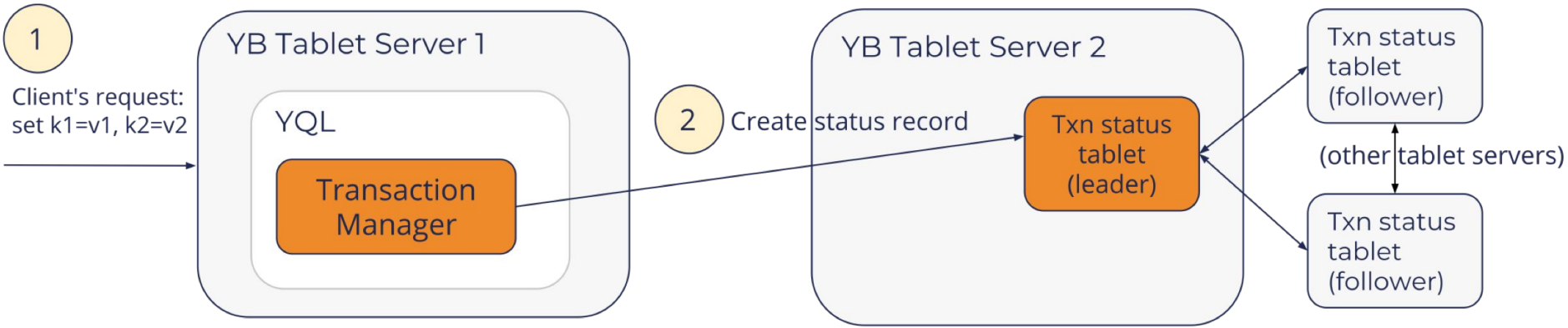
# Distributed Transactions - Write Path

1

Client's request:  
set k1=v1, k2=v2

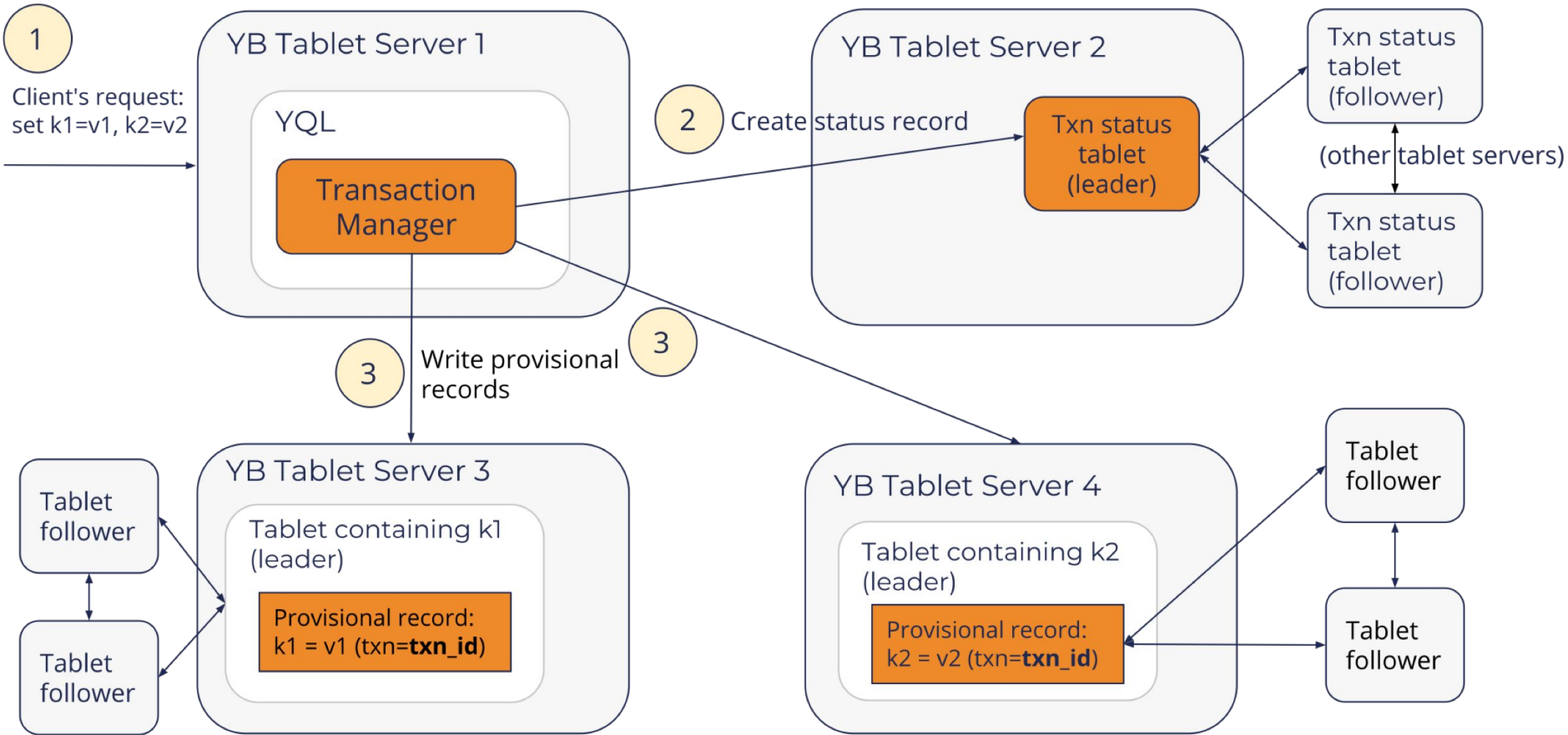


# Distributed Transactions - Write Path

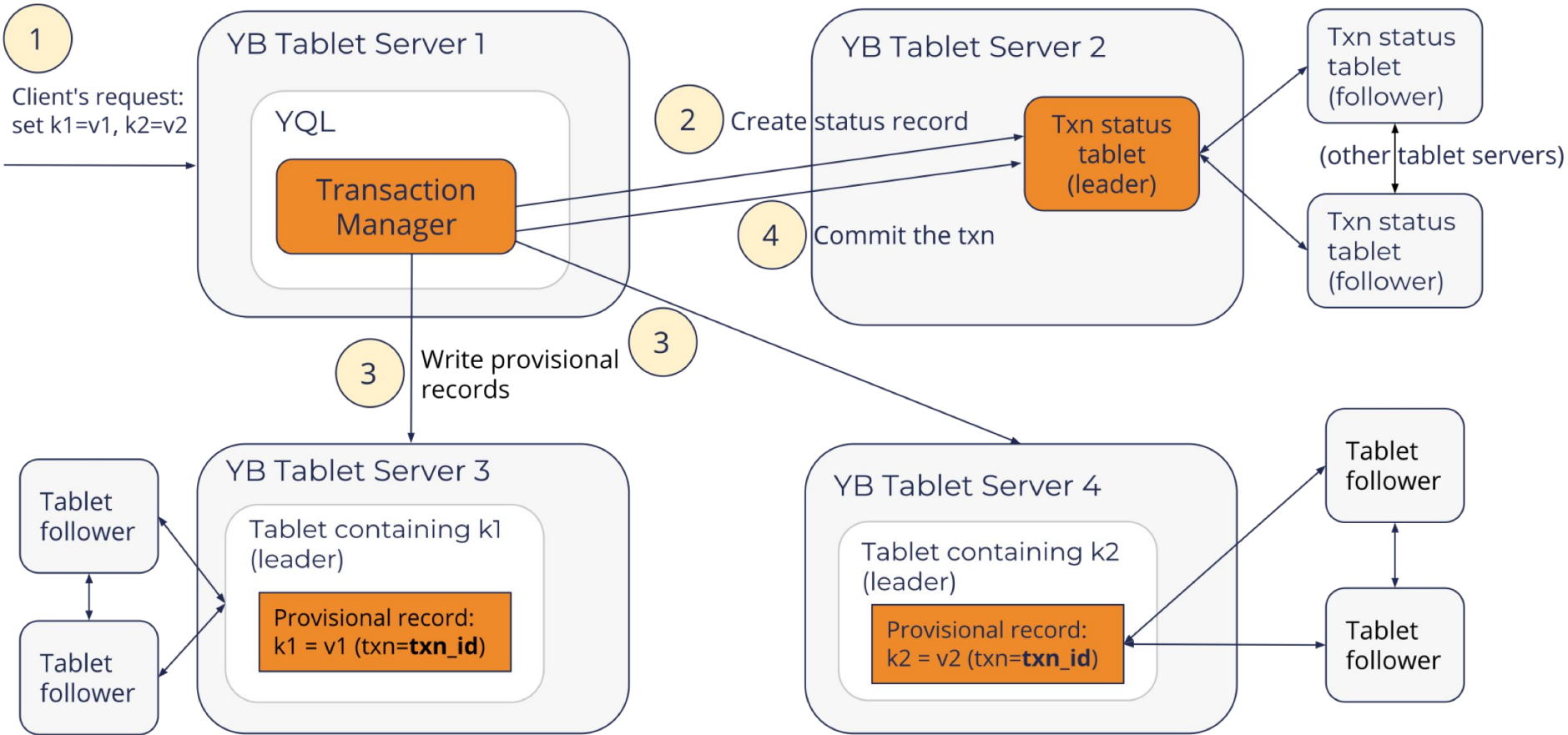




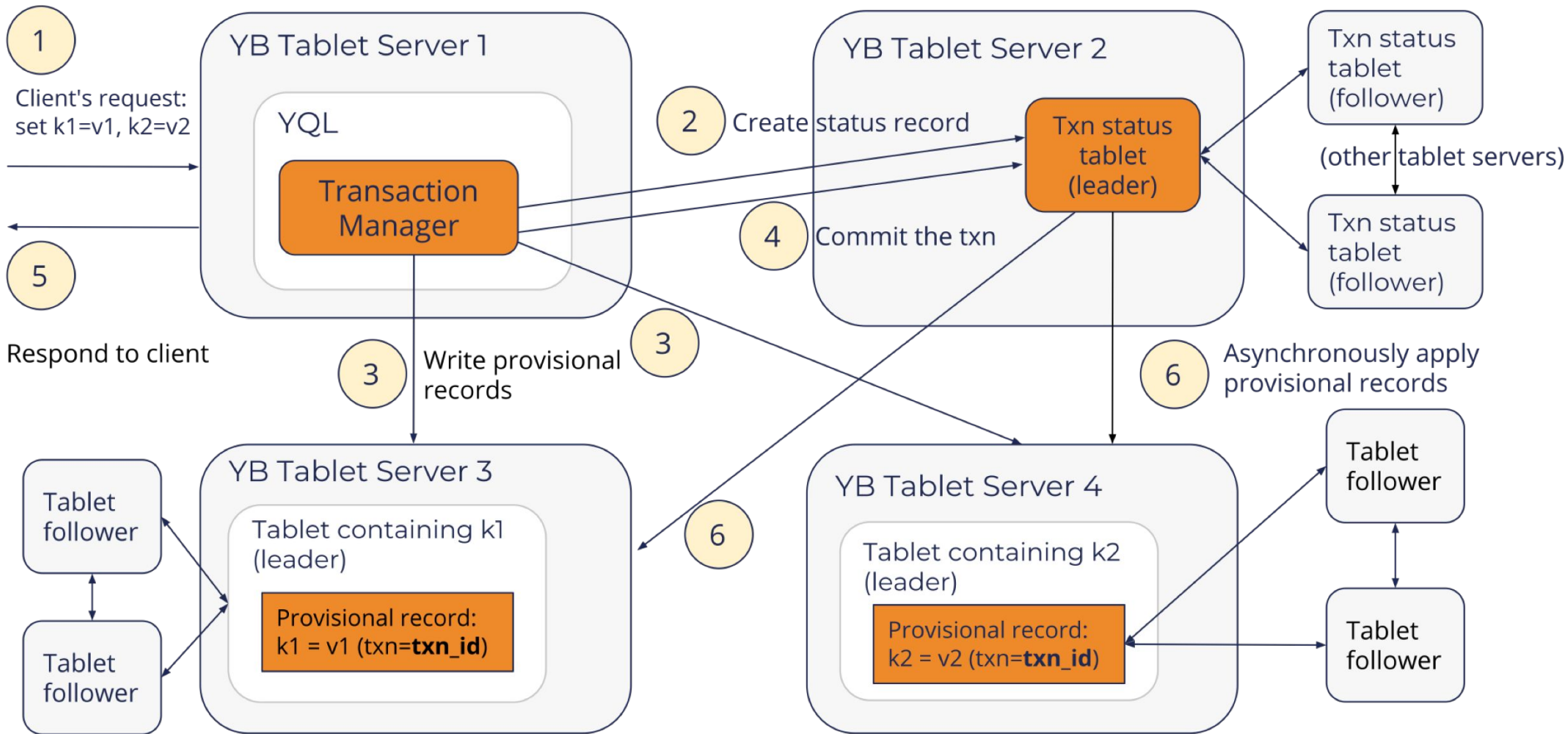
# Distributed Transactions - Write Path



# Distributed Transactions - Write Path



# Distributed Transactions - Write Path



# Was that interesting?

If so, join us on Slack for more:

[yugabyte.com/slack](https://yugabyte.com/slack)

# Thanks!