# PostgreSQL Compatibility

yugabyteDB

Build

Meet

Learn

# PostgreSQL Compatibility

Read Committed Isolation

Build

Meet

Learn

# PostgreSQL Compatibility

Read Committed Isolation &
Pessimistic Locking

yugabyteDB

Build

Meet

Learn

# Why compatibility?

Why compatibility?

Scaling your existing app should be straightforward

**NO** to "our db scales flawlessly & is fast but we don't support this yet, can your app workaround this?"

# Levels of compatibility

**Wire** - does the db work with existing Postgres drivers, same byte format, serialization etc.

**Syntax** - would the same syntax as PostgreSQL work?

**Feature** - parity in terms of different functionalities: triggers, stored procedures, gin indexes, etc

**Runtime** - matching execution semantics. An app shouldn't be able to say whether the db point is PostgreSQL or something else (barring any theoretical performance differences that stem from distribution)

Compatibility at the transactional layer is even more important because……

Compatibility at the transactional layer is even more important because......

ACID transactions -> tricky

Compatibility at the transactional layer is even more important because……

ACID transactions -> tricky

Distributed ACID transactions -> trickier

Fault tolerance & always available

Clock skew          Sharding, automatic shard splitting

Load balancing          and more ….
                    Horizontal scalability

Compatibility at the transactional layer is even more important because……

ACID transactions -> tricky

Distributed ACID transactions -> trickier
=> Impossible for the app to workaround incompatibility at the transactional layer

Fault tolerance & always available

Clock skew          Sharding, automatic shard splitting

Load balancing          and more ….
Horizontal scalability

# Isolation Levels in PostgreSQL

1. Serializable (supported by YB since long)
2. Repeatable Read (supported by YB since long)
3. Read Committed (new addition to YB!)
4. Read Uncommitted (same as Read Committed)

Strictness of isolation increases bottom to top

Performance increases top to bottom due to lower conflicts

# Read Committed Isolation Level

# Read Committed Isolation Level

Key ideas that define "read committed" -

1. From Postgres -
    a. New snapshot per statement in the transaction
    b. Apps never face serialization errors (40001), so don't have to retry those
2. Read restart errors no more [stems from **clock-skew** due to YugabyteDB's distributed nature]

# Read Committed Isolation Level

Key ideas that define "read committed" -

1. From Postgres -
   a. New snapshot per statement in the transaction
   b. Apps never face serialization errors, so don't have to retry those

**yugabyteDB**

# [Insert live terminal demo for point (1 & 2) here]

# Read Committed Isolation Level

Key ideas that define "read committed" -

1. From Postgres -
   a. New snapshot per statement in the transaction
   b. Apps never face serialization errors, so don't have to retry those
2. Read restart errors no more [specific to Yugabyte due to its distributed nature]

a 2 min detour to declutter read restarts …

Person Y

> _

N4

t=2

k=1 v=1 @ t=0

N1

t=5

N2

N3

Raft group for some tablet

N5

t=5

Person X

> _

Assume max clock skew = 8 units

yugabyteDB

Person Y

t=2

N4

k=1 v=1 @ t=0

N1

t=5

N2

N3

Raft group for some tablet

t=5

Write

N5

k=1 v=5 @ t=5

Person X

Assume max clock skew = 8 units

# 2 times units later .....

Person Y

t=4

N4

k=1 v=1 @ t=0
**k=1 v=5 @ t=7**

N1

t=7

N2

N3

Raft group for some tablet

t=7

Commit

N5

Person X

Assume max clock skew = 8 units

# 1 time unit later …..

Person Y

t=5

N4

Phone call

Person X

N5

t=8

k=1 v=1 @ t=0
**k=1 v=5 @ t=7**

N1

t=8

N2

N3

Raft group for some tablet

Assume max clock skew = 8 units

yugabyteDB

# 1 time unit later .....

Person Y

```
>_
```

t=6

Read

N4

N5

t=9

Person X

```
>_
```

k=1 v=1 @ t=0
**k=1 v=5 @ t=7**

t=9

N1

N2

N3

Raft group for some tablet

Assume max clock skew = 8 units

# 1 time unit later …..

Person Y

`>_`

🕐 t=6

N4

Read →

k=1 v=1 @ t=0
**k=1 v=5 @ t=7**

😅

N1

🕐 t=9

N2

N3

N5

🕐 t=9

Person X

`>_`

Raft group for some tablet

Assume max clock skew = 8 units

# 1 time unit later …..

Person Y

k=1 v=1 @ t=0
**k=1 v=5 @ t=7**

t=6

N4

Read restart error

t=9

N1

N2

N3

t=9

N5

Raft group for some tablet

Person X

Assume max clock skew = 8 units

# [Insert live terminal demo for point (3) i.e., read restarts here]

# The Temenos High Water Benchmark In (Big) Numbers [SKIP slide]

**temenos**

| | | | |
|---|---|---|---|
| **3000** Global Banking Customers | **41/50** Of The Top Global Banks | **1.2 Bn** Global Bank Customers | Investing **20%** Revenue in R&D |

**High Water Benchmark**

| | | | |
|---|---|---|---|
| **102K** Business Transactions Per Second | **100M** Customers **200M** Accounts | **4.1x** More Efficient For A Smaller $CO_2$ Footprint | **+40%** Better Performance |

**yugabyteDB**

| | | | |
|---|---|---|---|
| **350K** Database Reads Per Second | **80K** Database Writes Per Second | Inserts **3 ms** Selects **1 ms** Deletes **1 ms** | **39/3** DB Nodes / AWS AZ |

# Roadmap

Track further enhancements and their progress at
https://github.com/yugabyte/yugabyte-db/issues/13557

# View the docs



https://docs.yugabyte.com/preview/architecture/transactions/read-committed/

# Pessimistic locking

yugabyteDB

# Optimistic locking

**On conflict, roll-back** one of the conflicting txns based on priority

**Deadlocks are avoided** since transactions never wait

**Aborts transactions unnecessarily** in contentious workloads

**Behavior is incompatible** with PostgreSQL

# Pessimistic locking

**On conflict, wait** for the blocking transaction to commit or rollback

**Deadlocks are detected** among waiting transactions

**Aborts transactions minimally** in contentious workloads

**Behavior is compatible** with PostgreSQL

Along with pessimistic locking, comes the problem of detecting distributed deadlocks.

# Deadlock Detection

**Detection happens quickly** with detector triggering 1s after conflict

yugabyteDB

# Deadlock Detection

**Detection happens quickly** with detector triggering 1s after conflict

**Detection is guaranteed** once a deadlock is created

yugabyteDB

# Deadlock Detection

**Edge Chasing Algorithm**

Chandy and Misra, 1982 - *A Distributed Algorithm for Detecting Resource Deadlocks in Distributed Systems*

https://www.cs.utexas.edu/users/misra/scannedPdf.dir/ResourceDeadlock.pdf

**Wait-for graph** is formed by considering transactions as nodes and waiting-status as a directed edge

**Probes are sent** between transaction coordinators to detect cycles in this graph
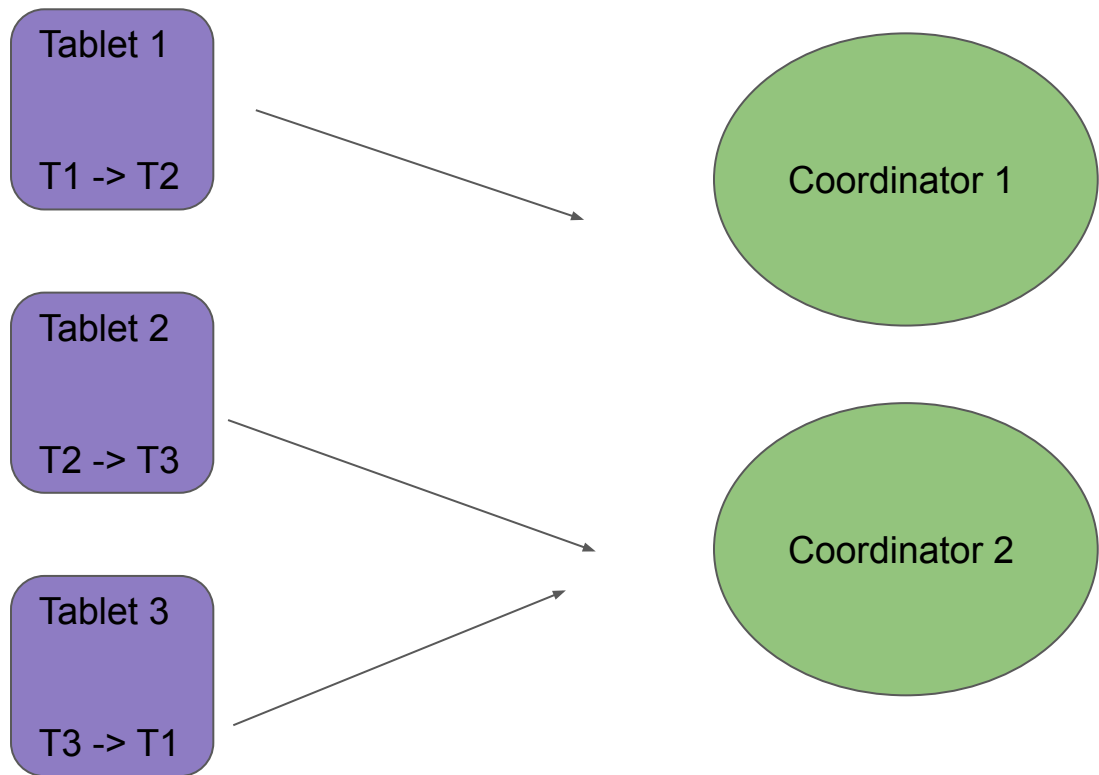
yugabyteDB

# Deadlock Detection

Tablet 1

T1 -> T2

Tablet 2

T2 -> T3

Tablet 3

T3 -> T1

# Deadlock Detection

Tablet 1

T1 -> T2

Tablet 2

T2 -> T3

Tablet 3

T3 -> T1

Coordinator 1

Coordinator 2

# Deadlock Detection

Tablet 1

T1 -> T2

Tablet 2

T2 -> T3

Tablet 3

T3 -> T1

Coordinator 1

T1 -> T2

Coordinator 2

T2 -> T3
T3 -> T1

# Deadlock Detection

Tablet 1

T1 -> T2

Tablet 2

T2 -> T3

Tablet 3

T3 -> T1

Coordinator 1

T1 -> T2

probe_id
T1 -> T2
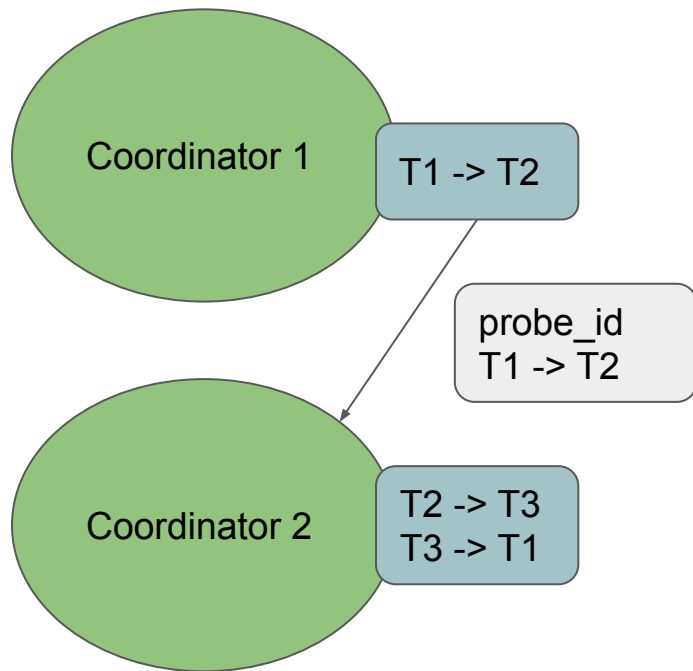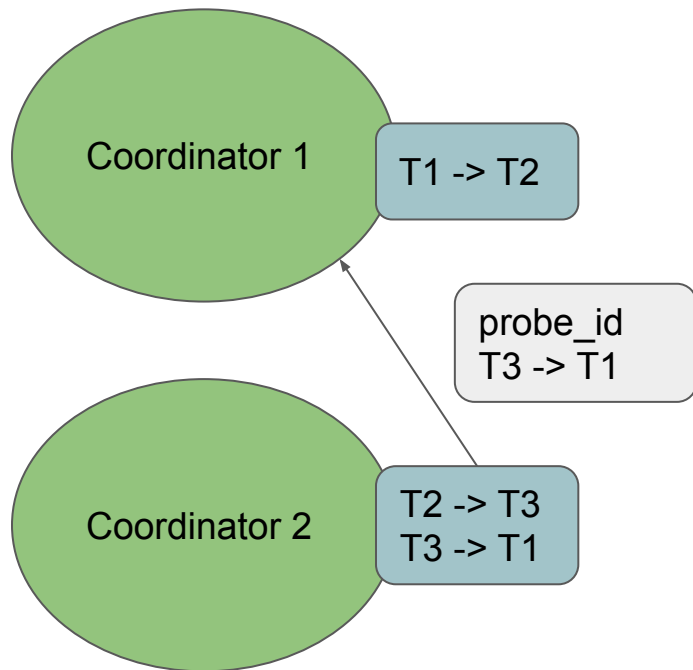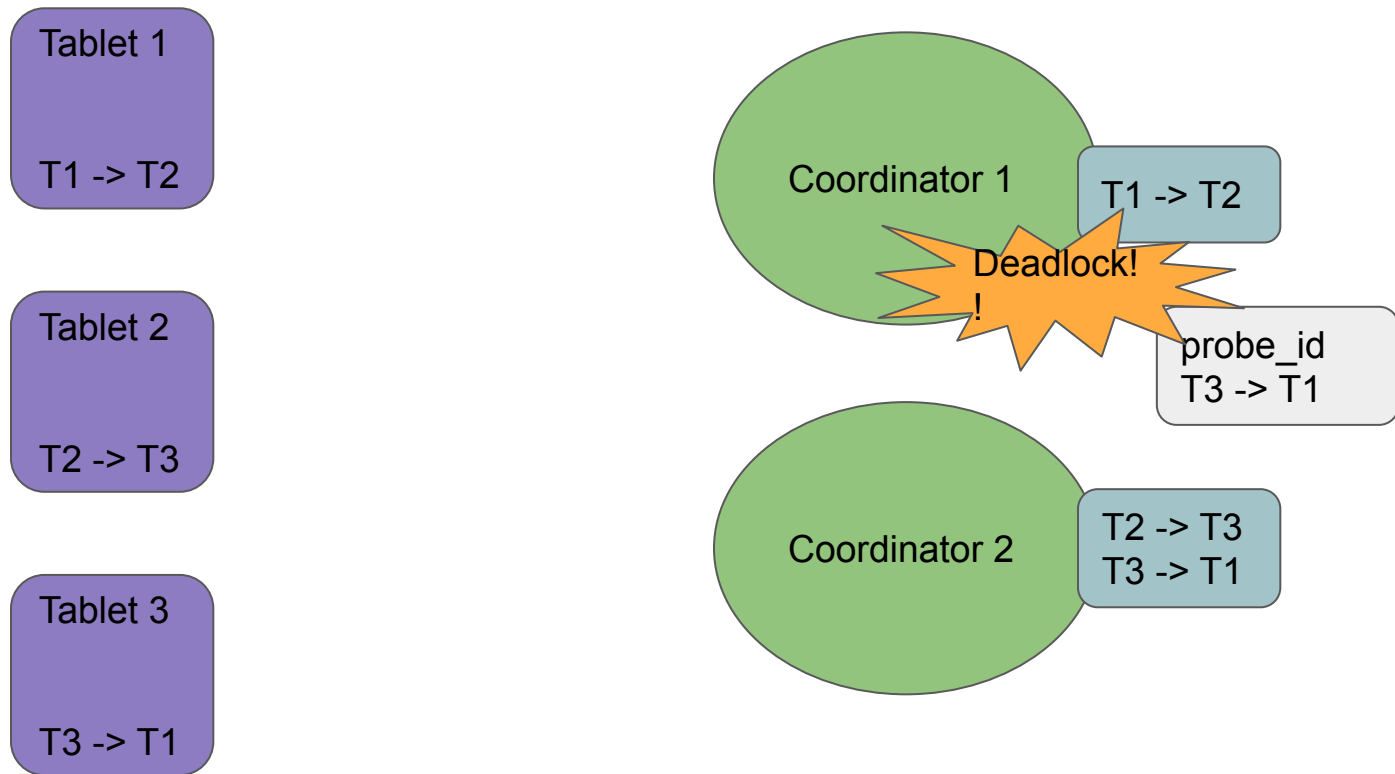
Coordinator 2

T2 -> T3
T3 -> T1

# Deadlock Detection

Tablet 1

T1 -> T2

Tablet 2

T2 -> T3

Tablet 3

T3 -> T1

Coordinator 1

T1 -> T2

probe_id
T3 -> T1

Coordinator 2

T2 -> T3
T3 -> T1

# Deadlock Detection

# Deadlock Detection

**Detection happens quickly** with detector triggering 1s after conflict

**Detection is guaranteed** once a deadlock is created

**Overhead is optimal** requiring constant additional memory and only one probe per conflicting transaction. Probes could trigger many RPCs depending on wait-for graph

**yugabyteDB**

# Additional Resources

View our docs page for [Read Committed isolation level](#)

View our spec document for [pessimistic locking](#)

See our [roadmap](#) for this area, more exciting Items coming your way!

# Thank You

Join us on Slack: yugabyte.com/slack

Star us on Github:
github.com/yugabyte/yugabyte-db