

Benchmarking Cloud Native Database

Iqbal Farabi and Tara Baskara
KubeCon + CloudNativeCon Europe 2019



KubeCon



CloudNativeCon

Europe 2019

Hola!

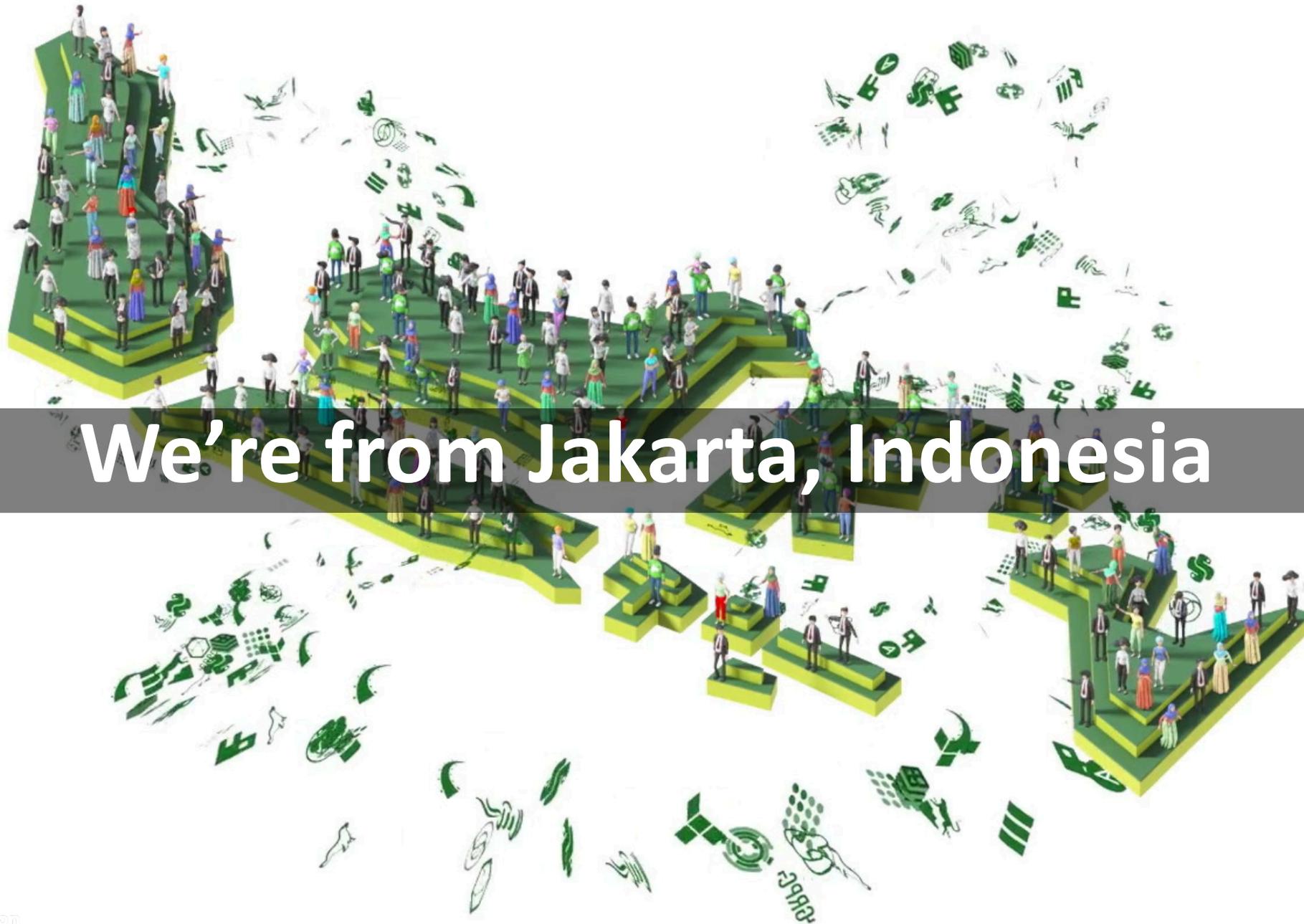


Iqbal Farabi
System Engineer
Go-Jek Indonesia
@iqbal_farabi



Tara Baskara
System Engineer
Go-Jek Indonesia
@iqbal_farabi





We're from Jakarta, Indonesia



KubeCon



Cloud Native Con

China 2018







JAKARTA



BANGALORE



SINGAPORE



THAILAND



VIETNAM

The Journey Continues...

GOJEK
INTERNATIONAL
EXPANSION

We officially announce our international expansion to four different markets: Vietnam, Thailand, Singapore and the Philippines!

#JUSTGOJEKIT



Outline



Today we will discuss about...

Cloud Native Database

- What is cloud native database? What databases do we pick to run the experiments on? What are their characteristics?

YCSB

- What is YCSB? What does it measure? What are the different kind of workloads defined?

Experiments and Results

- Setup, experiments, and brief explanation of the results.



Cloud Native Database



Definition

CNCF [defines](#) cloud native technologies as, "technologies that empower organizations to build and run scalable applications in modern and dynamic environments. Cloud native technologies enable loosely coupled systems that are resilient, manageable, and observable".



Landscape

 <p>CarbonData APACHE Apache CarbonData ★ 898 Apache Software Foundation</p>	 <p>Ignite apache Apache Ignite ★ 2,469 Apache Software Foundation</p>	 <p>ArangoDB ArangoDB ★ 7,905 Funding: \$16.9M</p>	 <p>BIGCHAINDB BigchainDB ★ 3,163 Funding: \$5.37M</p>	 <p>cassandra Cassandra ★ 5,110 Apache Software Foundation</p>	 <p>Cockroach LABS CockroachDB ★ 16,036 Cockroach Labs Funding: \$53.5M</p>	 <p>Couchbase Couchbase ★ 159 Funding: \$150M</p>	 <p>CRATE.IO Crate.io ★ 2,416 Funding: \$17.9M</p>	 <p>Dgraph Dgraph Labs ★ 9,433 Funding: \$2.95M</p>
 <p>druid Druid ★ 7,979 Apache Software Foundation</p>	 <p>FoundationDB FoundationDB ★ 9,242 MCap: \$965B</p>	 <p>hadoop Hadoop ★ 8,927 Apache Software Foundation</p>	 <p>hazelcast IMDG Hazelcast ★ 3,073 Funding: \$13.6M</p>	 <p>IBM DB2 IBM DB2 MCap: \$124B</p>	 <p>iguazio iguazio Funding: \$48M</p>	 <p>Infinispan Infinispan ★ 732 MCap: \$32.4B</p>	 <p>InterSystems IRIS Data Platform InterSystems IRIS Data Platform InterSystems</p>	 <p>MariaDB MariaDB ★ 2,633 Funding: \$98.2M</p>
 <p>MEMSQL MemSQL Funding: \$108M</p>	 <p>Microsoft SQL Server Microsoft SQL Server MCap: \$994B</p>	 <p>mongoDB MongoDB ★ 15,921 MCap: \$7.56B</p>	 <p>MySQL MySQL ★ 3,724 MCap: \$188B</p>	 <p>neo4j Neo4j ★ 6,372 Funding: \$160M</p>	 <p>noms NomsDB ★ 6,849 MCap: \$128B</p>	 <p>ORACLE DATABASE Oracle Database MCap: \$188B</p>	 <p>OrientDB OrientDB ★ 3,828 MCap: \$155B</p>	 <p>PERCONA Percona Server for MySQL ★ 554</p>
 <p>piloosa Piloosa ★ 1,648 Funding: \$3.67M</p>	 <p>PostgreSQL PostgreSQL ★ 5,165</p>	 <p>presto Presto ★ 362 Presto Software Foundation</p>	 <p>Qubole Qubole Funding: \$75M</p>	 <p>redis Redis ★ 36,054 Funding: \$147M</p>	 <p>RethinkDB RethinkDB ★ 22,072 Linux Foundation</p>	 <p>SCYLLA Scylla ★ 5,043 Funding: \$35M</p>	 <p>snowflake Snowflake Funding: \$929M</p>	 <p>software AG Software AG MCap: \$2.7B</p>
 <p>STOLON Stolon ★ 2,068 Sorint.Lab</p>	 <p>TiDB TiDB ★ 18,411 Funding: \$71.6M</p>	 <p>TiKV TiKV ★ 5,019 Cloud Native Computing Foundation (CNCF)</p>	 <p>VERTICA Vertica Funding: \$30.5M Vertica Systems</p>	 <p>Vitess Vitess ★ 7,893 Cloud Native Computing Foundation (CNCF)</p>	 <p>YugaByte YugaByte DB ★ 1,063 Funding: \$24M</p>			



Filtering

- Open source
- Operational database
- ACID compliance
- Provides SQL-like API



The List

CockroachDB

- <https://www.cockroachlabs.com/>

TiDB

- <https://pingcap.com/en/>

YugaByte DB

- <https://www.yugabyte.com/>



CockroachDB

CockroachDB is a distributed SQL database built on a transactional and strongly-consistent key-value store.

It aims to:

- Scales horizontally
- Provides fault resiliency
- Supports strongly-consistent ACID transactions
- Provides a familiar SQL API with PostgreSQL-like syntax



CockroachDB Architecture

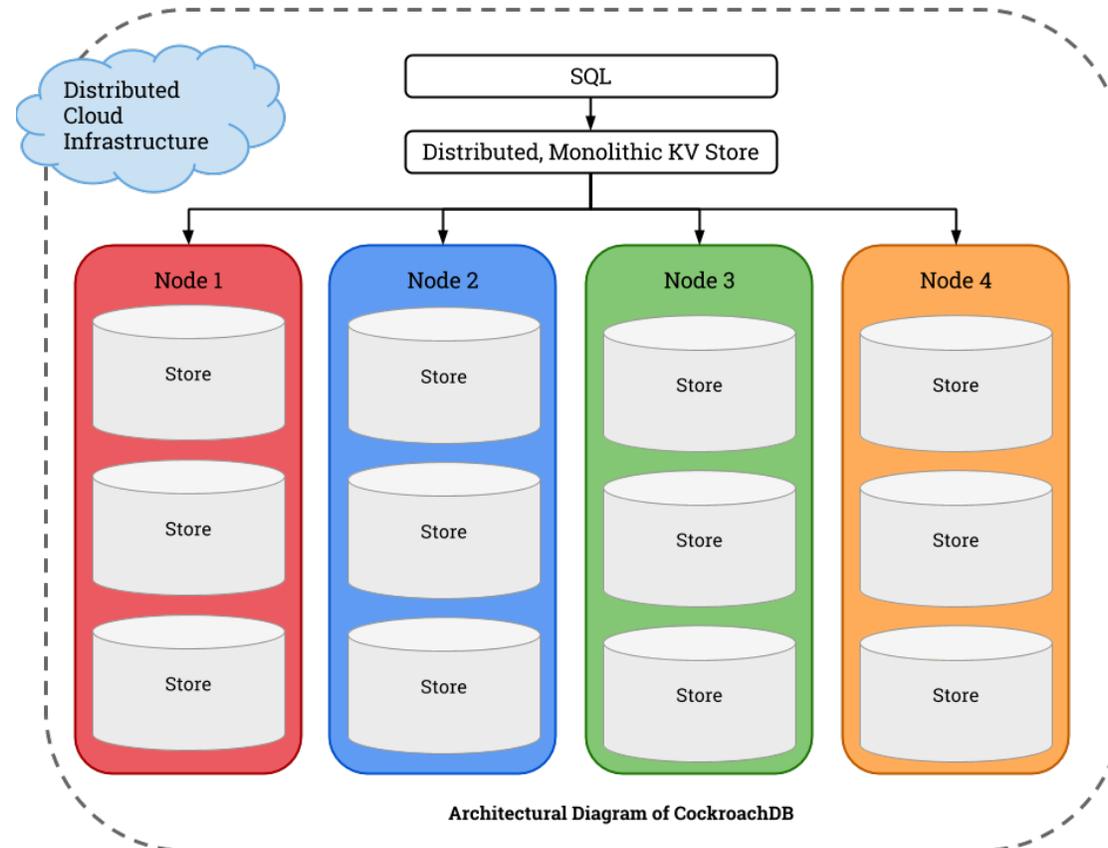


image source: <https://thenewstack.io/cockroachdb-unkillable-distributed-sql-database/>

TiDB

TiDB is an open-source NewSQL database that supports Hybrid Transactional and Analytical Processing (HTAP) workloads.

It aims to:

- Scales horizontally
- Provides fault resiliency
- Supports distributed transactions with strong consistency
- Provides a familiar SQL API with MySQL-like syntax



TiDB Architecture

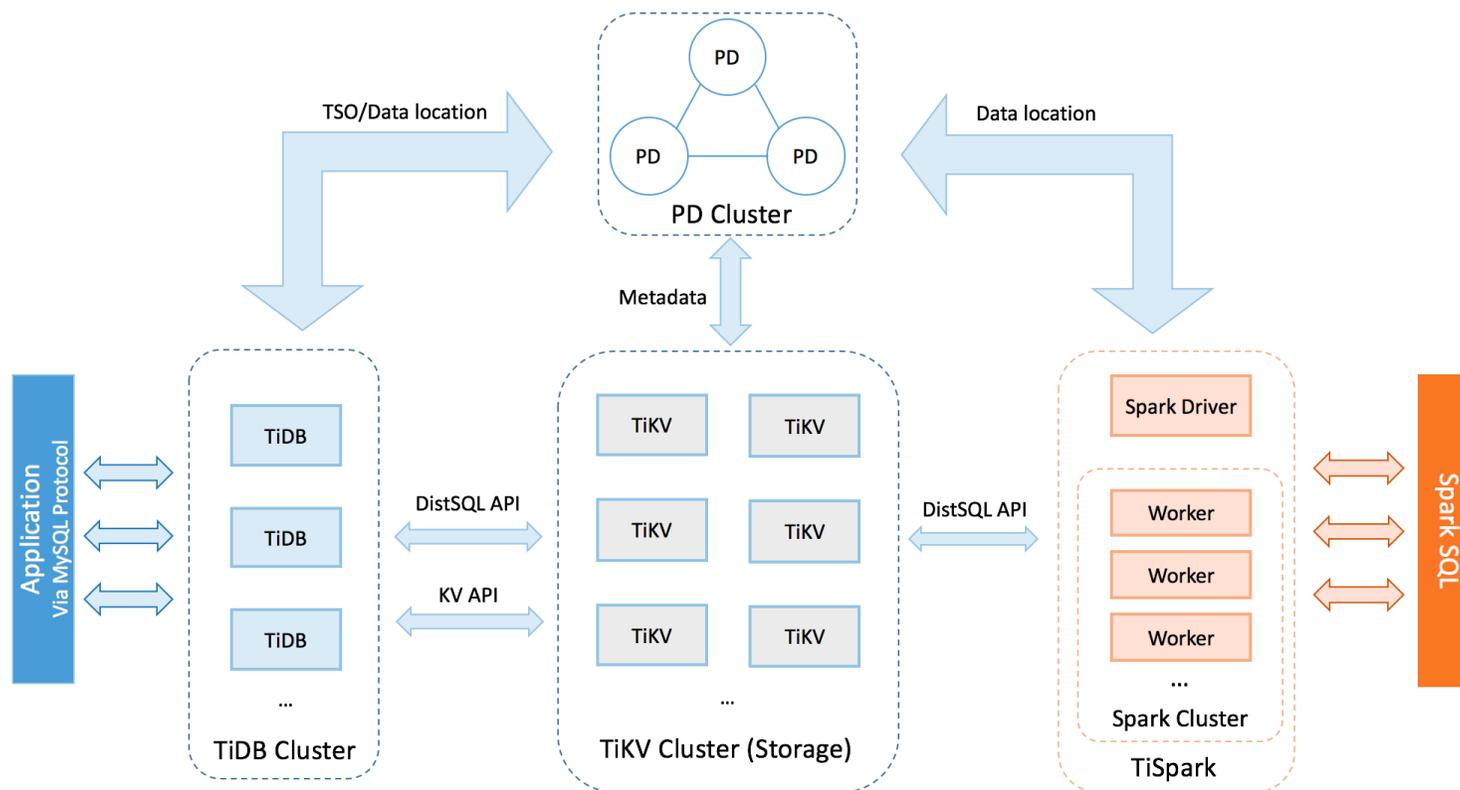


image source: <https://pingcap.com/docs/architecture/>

YugaByte DB

CockroachDB is a distributed SQL database built on a transactional and strongly-consistent key-value store.

It aims to:

- Scales with autosharding
- Provides fault resiliency
- Supports multi-shard ACID transactions
- Provides YugaByte Structured Query Language (YSQL) and YugaByte Cloud Query Language (YCQL) APIs.



YugaByte DB Architecture

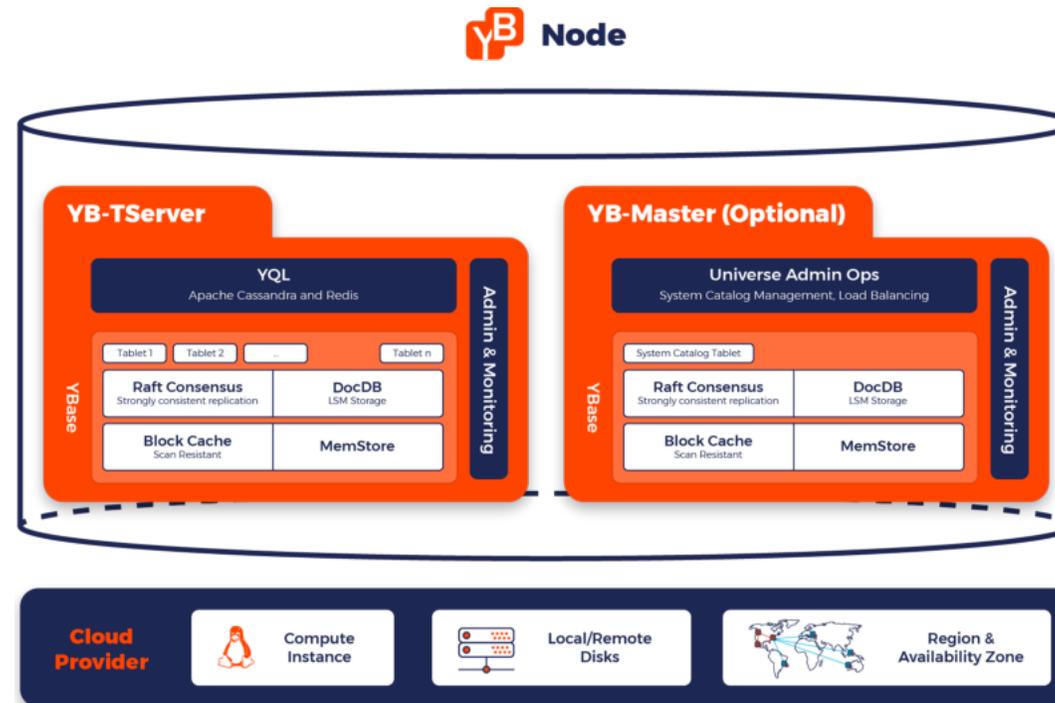


image source: <https://blog.yugabyte.com/yugabyte-db-architecture-diverse-workloads-with-operational-simplicity/>

Classifications (1)

According to Brian F. Cooper et al., in their [YCSB paper](#), there are four main tradeoffs faced by cloud serving systems:

- read vs write performance
- latency vs durability
- synchronous vs asynchronous replication
- data partitioning



Classifications (2)

Latency vs Durability

Write operations synced to disk before returning success will increase durability in the case of system failures. However it might lower throughput and increase latency.

Write operations stored in memory and synced later will increase throughput and decrease latency, but might increase the risk of data loss in case of system failures.



Classifications (3)

Synchronous vs Asynchronous Replication

Synchronous replication ensures consistency among all nodes but might increase latency.

Asynchronous replication decreases latency but might cause data loss if failure happens to nodes with data that not replicated yet.



Classifications (4)

Data Partitioning

Strictly row-based partitioning allows efficient access to an entire record of data.

Column-based partitioning allows efficient access for a subset columns when retrieving multiple records.



Classifications (5)

Database	Latency/durability	Sync/async replication	Row/column partitioning
CockroachDB	Durability	Synchronous	Row
TiDB	Durability	Synchronous	Row
YugaByte DB	Latency	Asynchronous	Row



YCSB



Yahoo! Cloud Serving Benchmark

Created by Brian F. Cooper et al. to create a **standard benchmark** and benchmarking framework to assist in the evaluation of different **cloud systems**.

Focus on **serving systems**, which serve read and write workloads, over batch or analytical systems.



Workloads Data

YCSB workloads data look like the following:

- 1 table named “usertable”
- 10 string fields, 1 primary key with content like “user123456”, 9 fields with content a random string of ASCII characters with 100 bytes length
- 1,000,000 records
- 1,000,000 operations with 1,000 threads for Workload A, B, C, and D
- 1,000,000 operations with 100 threads for Workload E



Operations

Operations performed by YCSB are:

- Insert: insert new record.
- Update: update a record by replacing the value of one field.
- Read: read a record, either one randomly chosen field or all fields.
- Scan: scan records in order, starting at randomly chosen record key with randomly chosen number of records.



Distributions

To choose which operations (insert, update, read, or scan) to perform on which records and how many records, YCSB has several built-in distributions:

- Uniform: choose an item uniformly at random.
- Zipfian: some item will be extremely popular, most records will be unpopular.
- Latest: like Zipfian with preference of latest inserted records as popular distribution.
- Multinomial: probabilities of each item can be specified.



Workloads

Workload	Operations	Distribution	Application Example
A – Update Heavy	Read: 50% Update: 50%	Zipfian	User session: session store recording recent actions
B – Read Heavy	Read: 95% Update: 5%	Zipfian	Photo tag: add tag is update, but most ops are reading tags
C – Read Only	Read: 100%	Zipfian	User profile cache
D – Read Latest	Read: 95% Insert: 5%	Latest	User status updates; people want to read the latest
E – Short Ranges	Scan: 95% Insert: 5%	Zipfian/Uniform	Threaded conversations



Experiments and Results



Cluster Setup

3 nodes GKE cluster with following specifications:

- n1-standard-16 machine type
- 1000 GB Local SSD
- 60 GB RAM



Statefulset Setup

Resource:

- 14 vCPU request, 16 vCPU limit
- 30 GB RAM request, 60 GB RAM limit
- 500 GB SSD local persistent volume



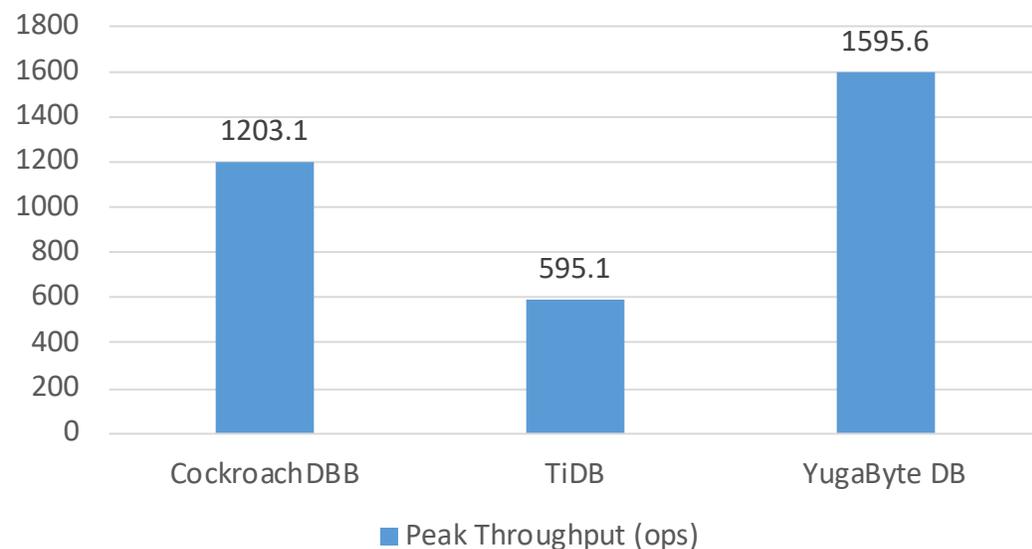
go-ycsb

All experiments in this presentation all done using a Go port of YCSB called [go-ycsb](#) created by engineers at PingCap, the company that creates TiDB.

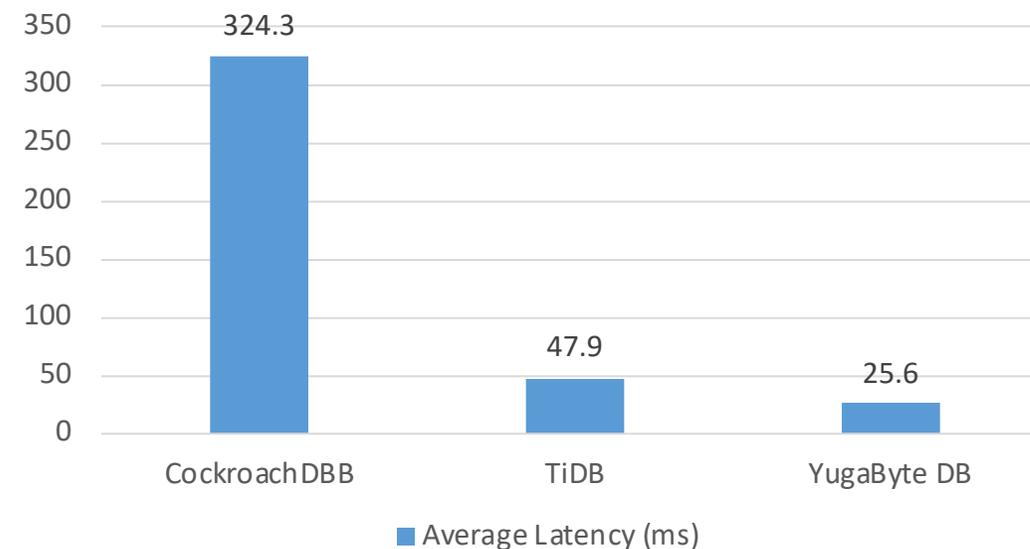


Workload A - Read

Throughput

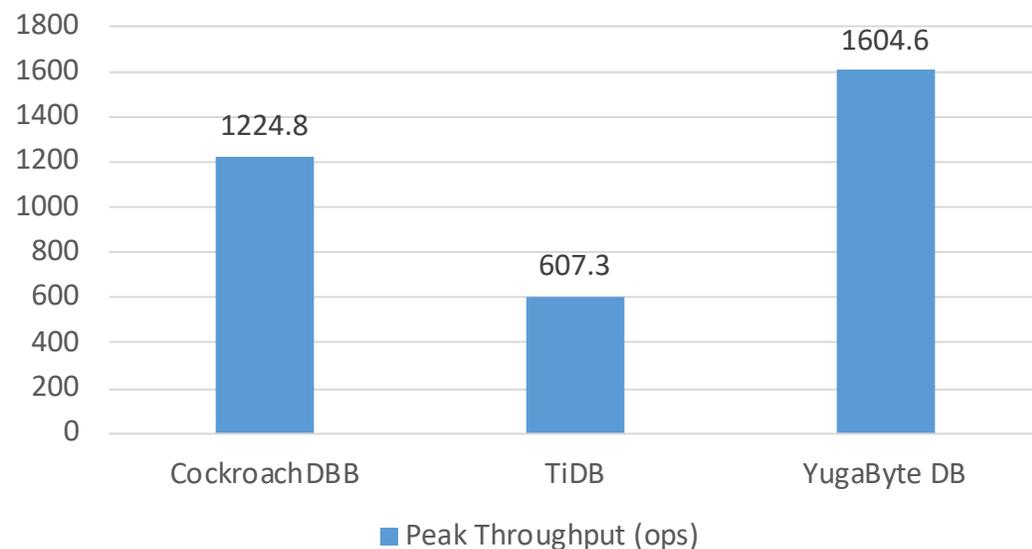


Latency

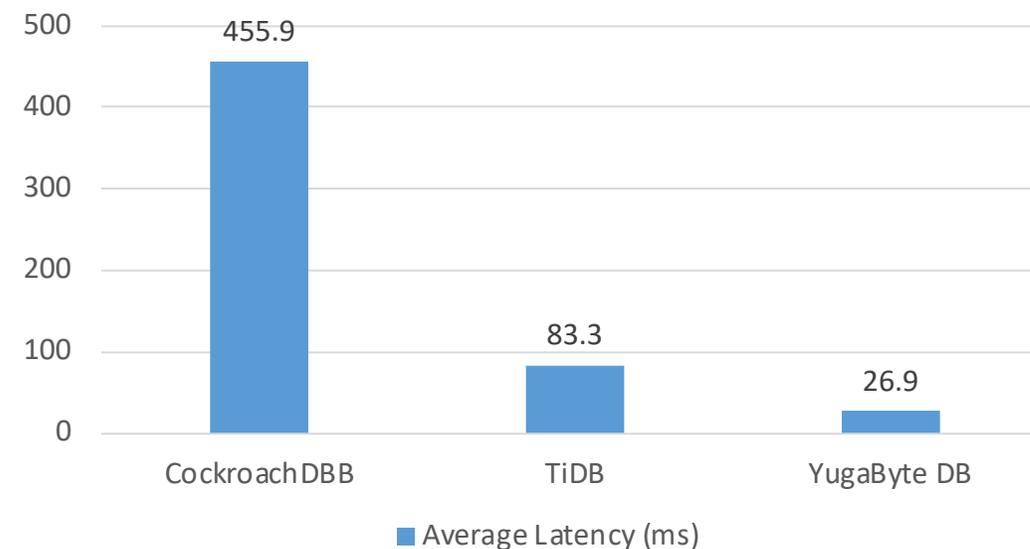


Workload A - Update

Throughput

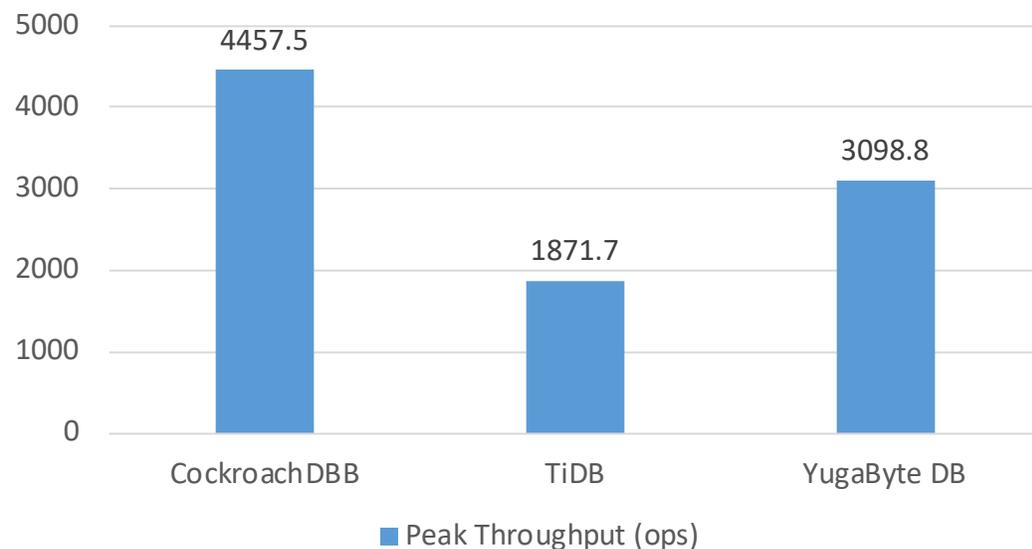


Latency

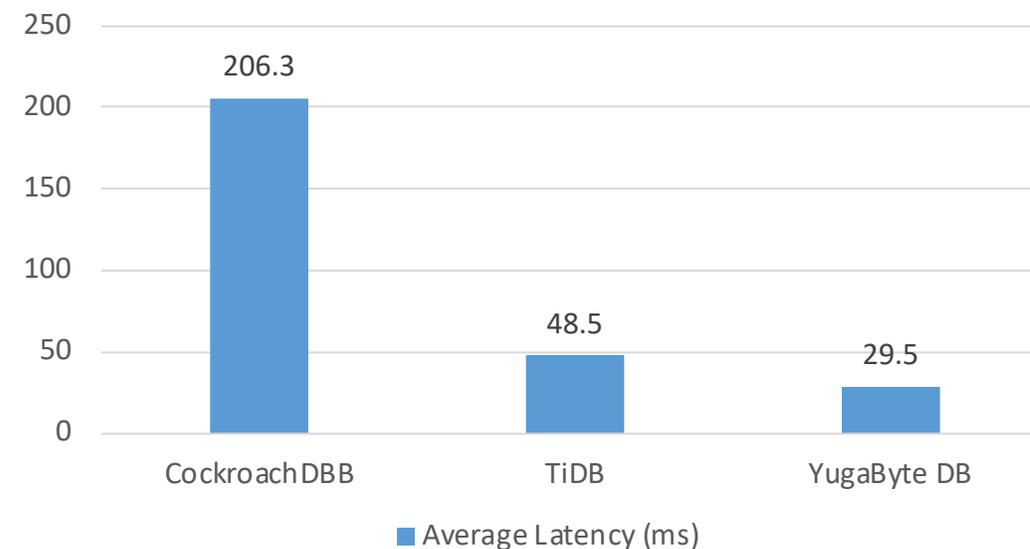


Workload B - Read

Throughput

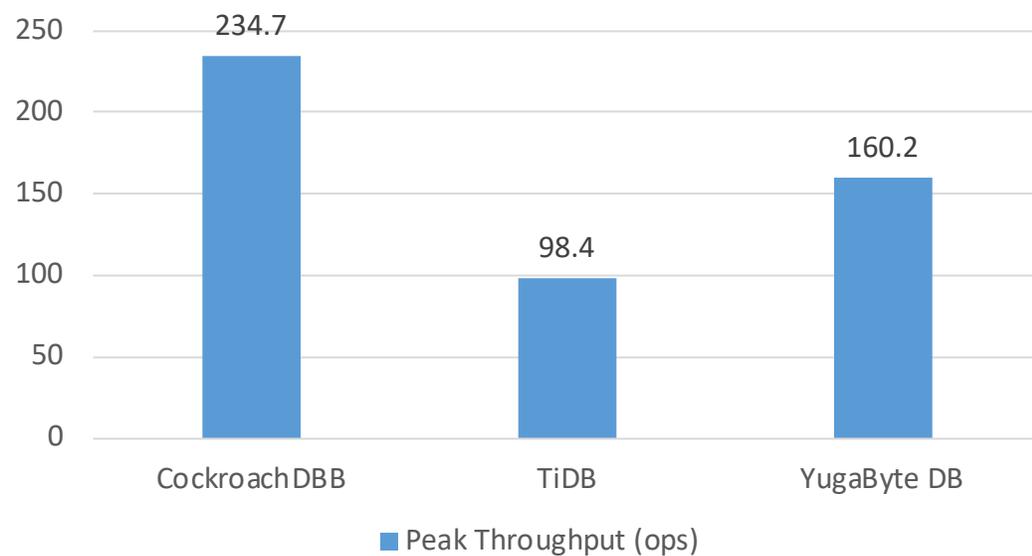


Latency

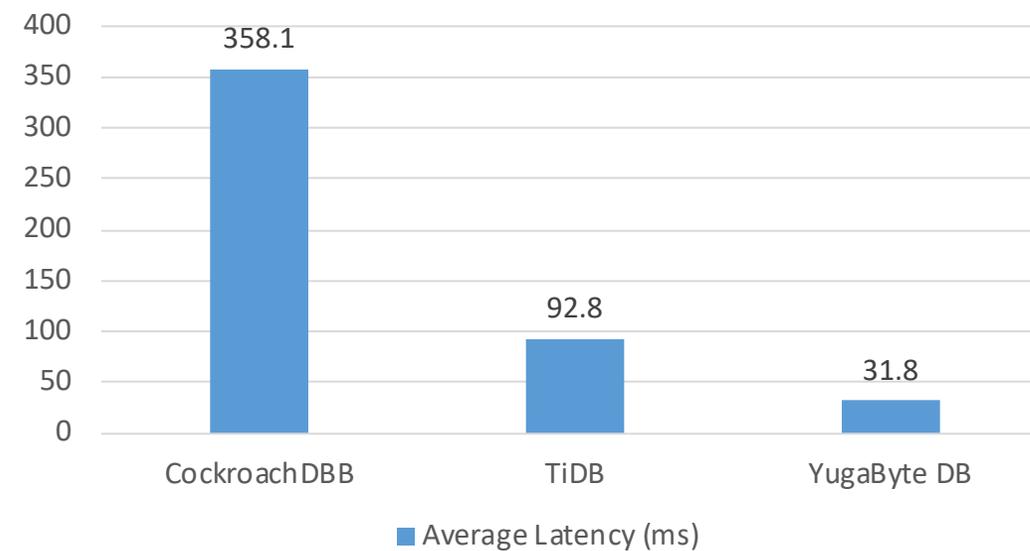


Workload B - Update

Throughput

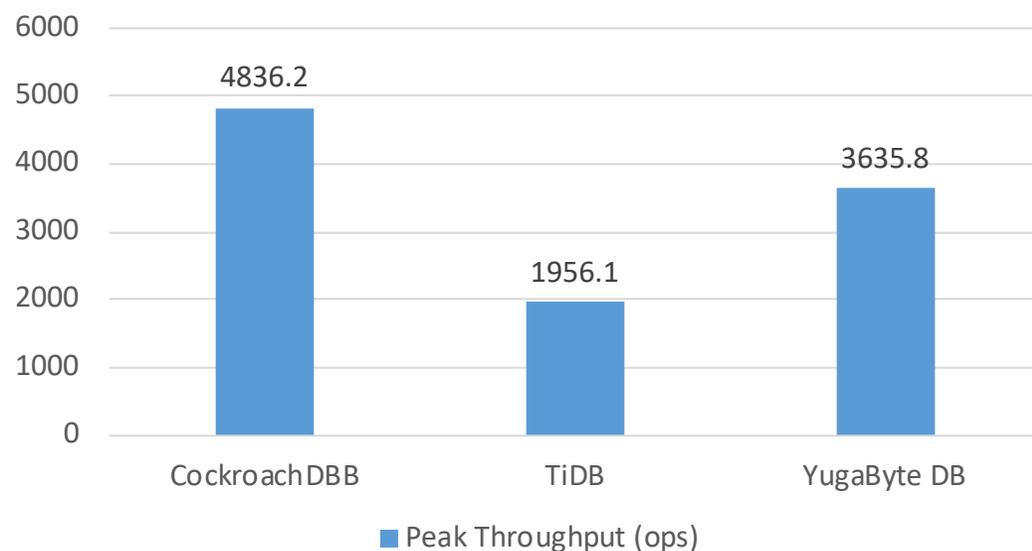


Latency

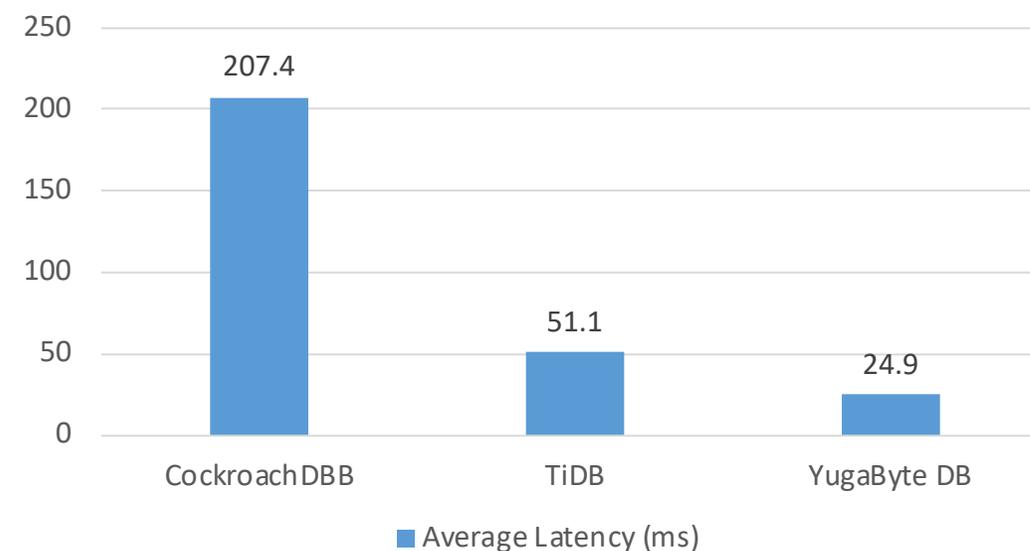


Workload C - Read

Throughput

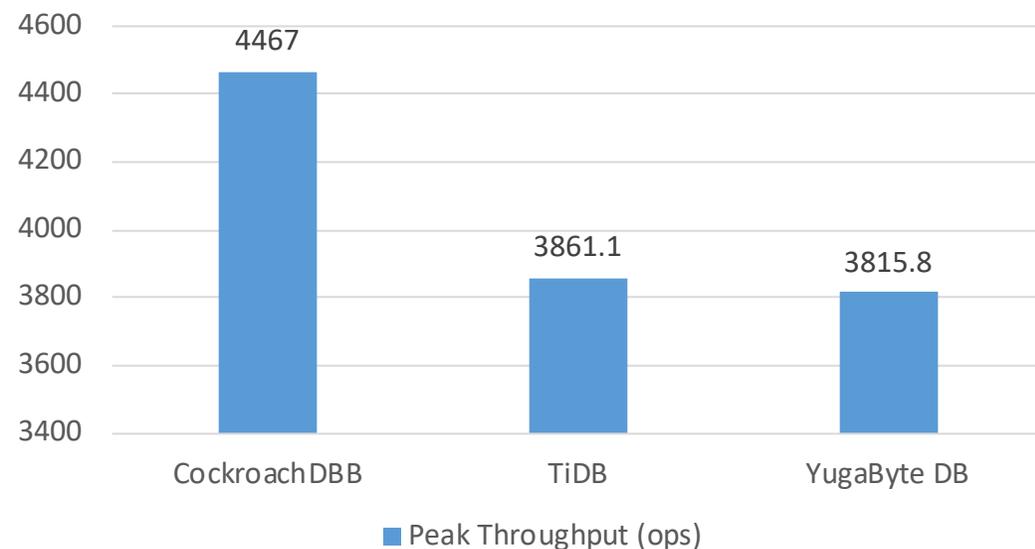


Latency

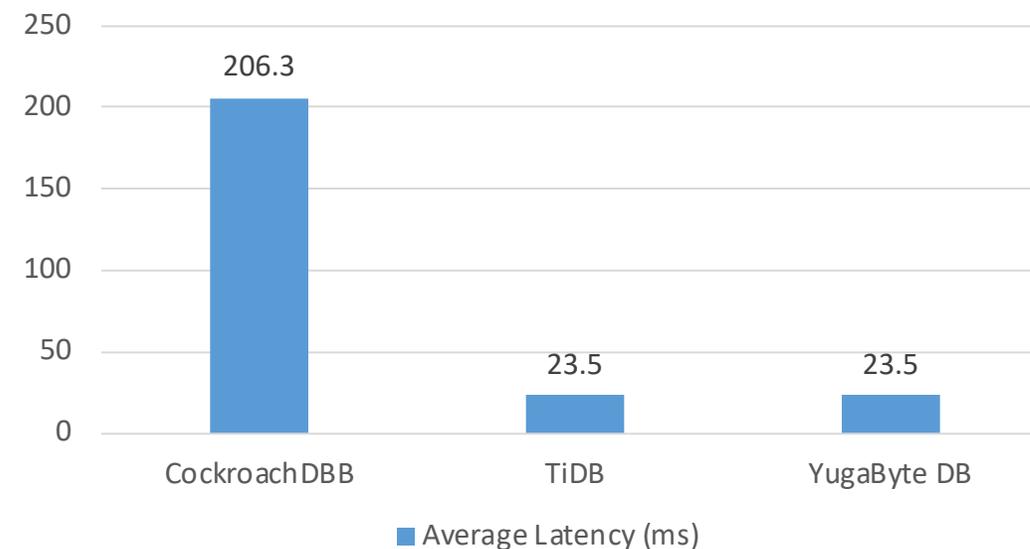


Workload D - Read

Throughput

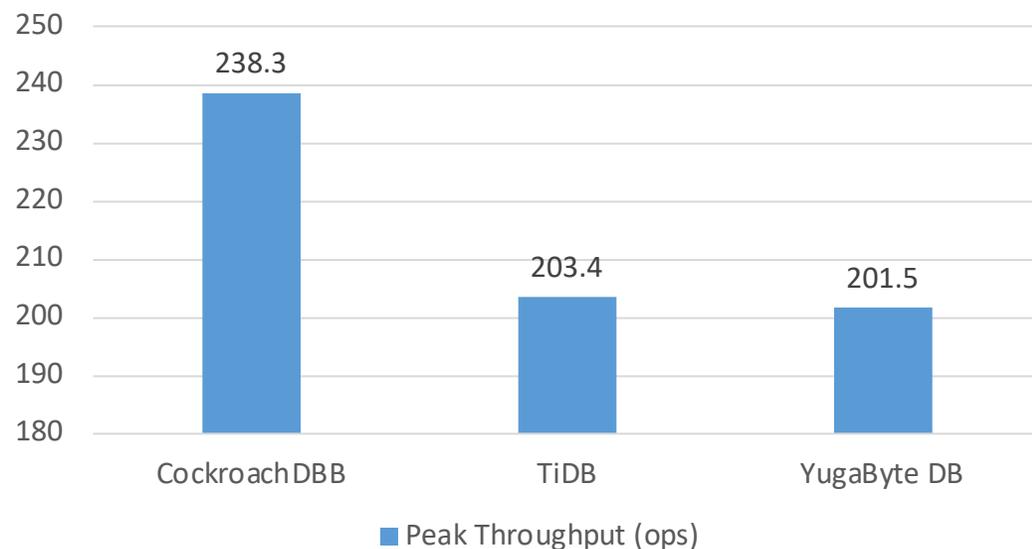


Latency

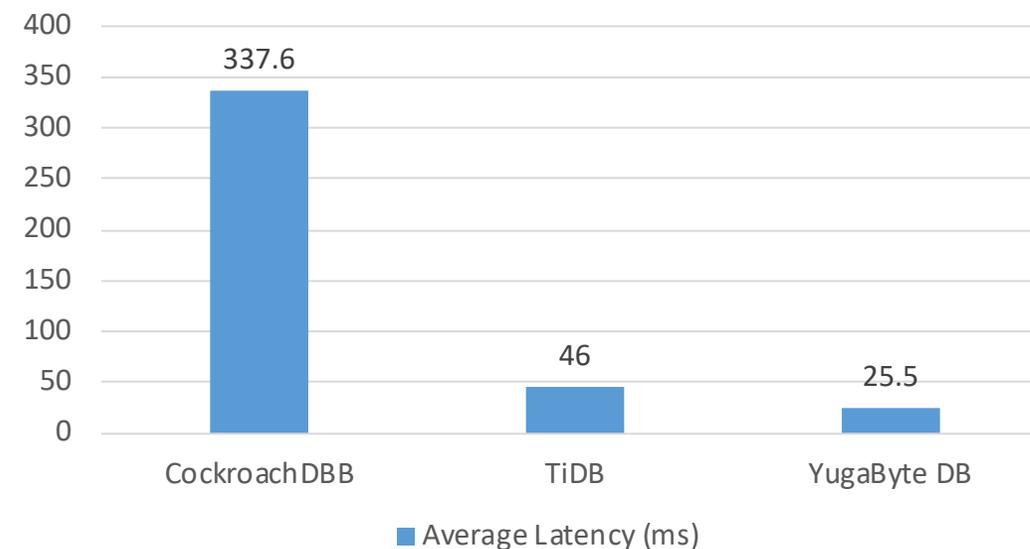


Workload D - Insert

Throughput

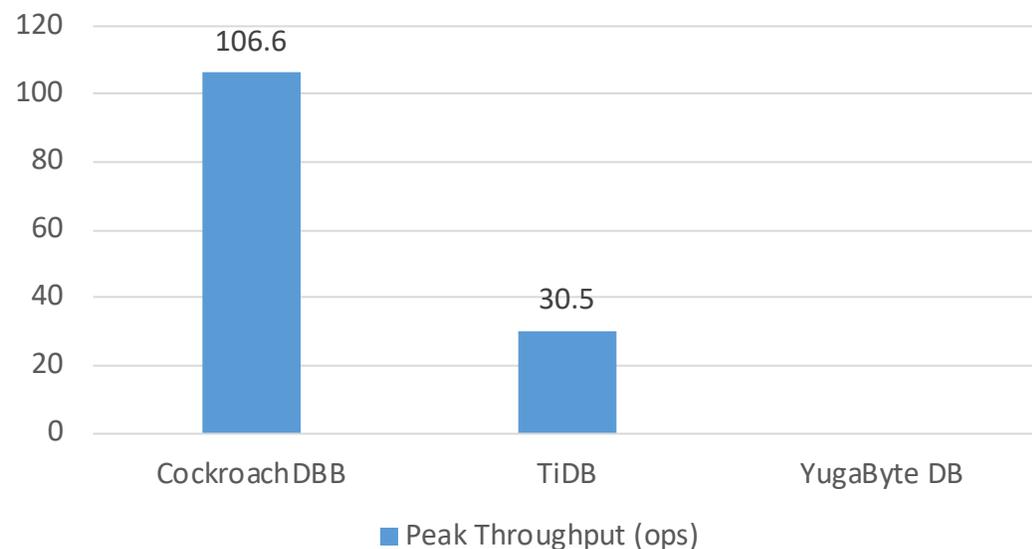


Latency

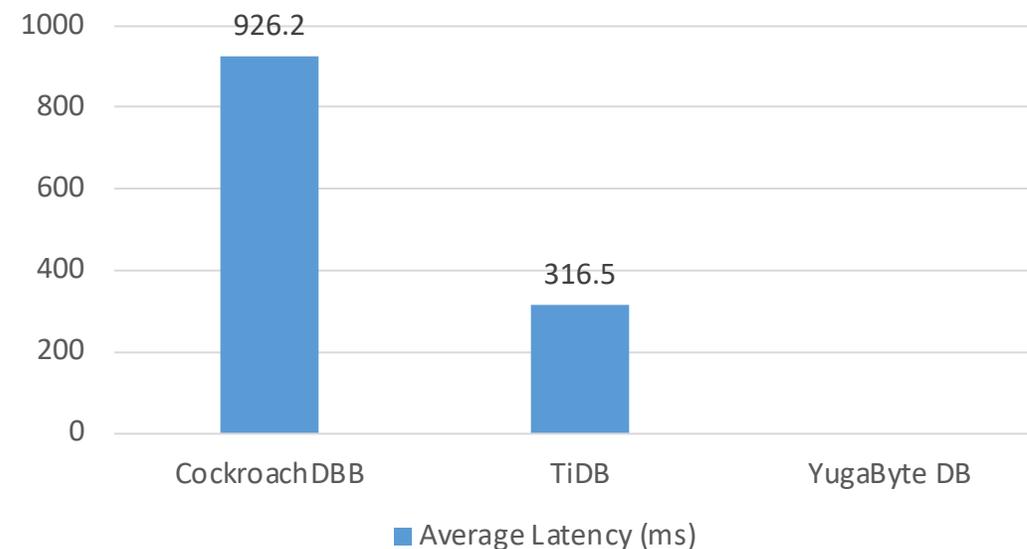


Workload E - Scan

Throughput

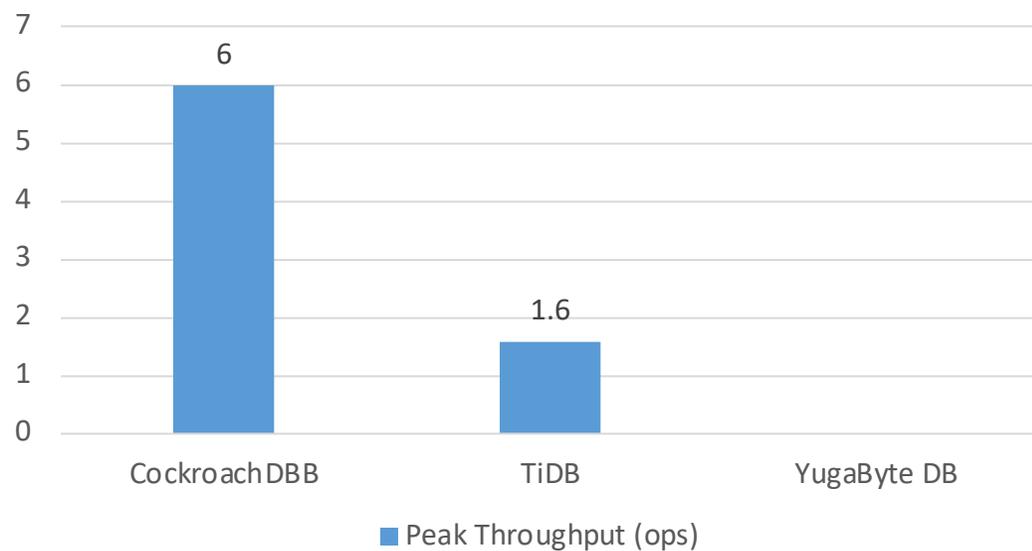


Latency

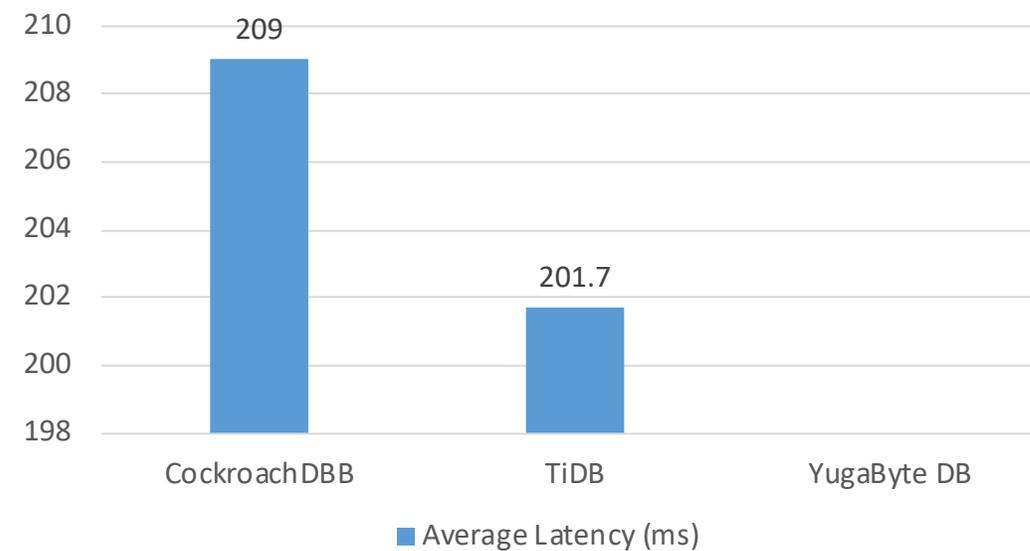


Workload E - Insert

Throughput



Latency



Discussions (1)

- All databases do not perform as well as they would if they run on a dedicated VM cluster with the same spec.
- All databases perform very well on read operations, especially in read-heavy workload (Workload B and Workload C).
- All databases perform fairly well on update operations in update-heavy workload (Workload A) but not so much in read-heavy workload (Workload B).



Discussions (2)

- All databases perform fairly well on insert operations in read-heavy workload (Workload D), but perform poorly in scan-heavy workload (Workload E).
- All databases does not perform well in scan operations in scan-heavy workload (Workload E).



Further Study

For further study, we recommend the following:

- Figure out the bottleneck that prevent databases to perform as well as they do in dedicated VM cluster instead of on top of Kubernetes cluster
- Communicate closely with engineers from respective database to gain more insights on how to fine tune each database



Q & A





Gracias!

