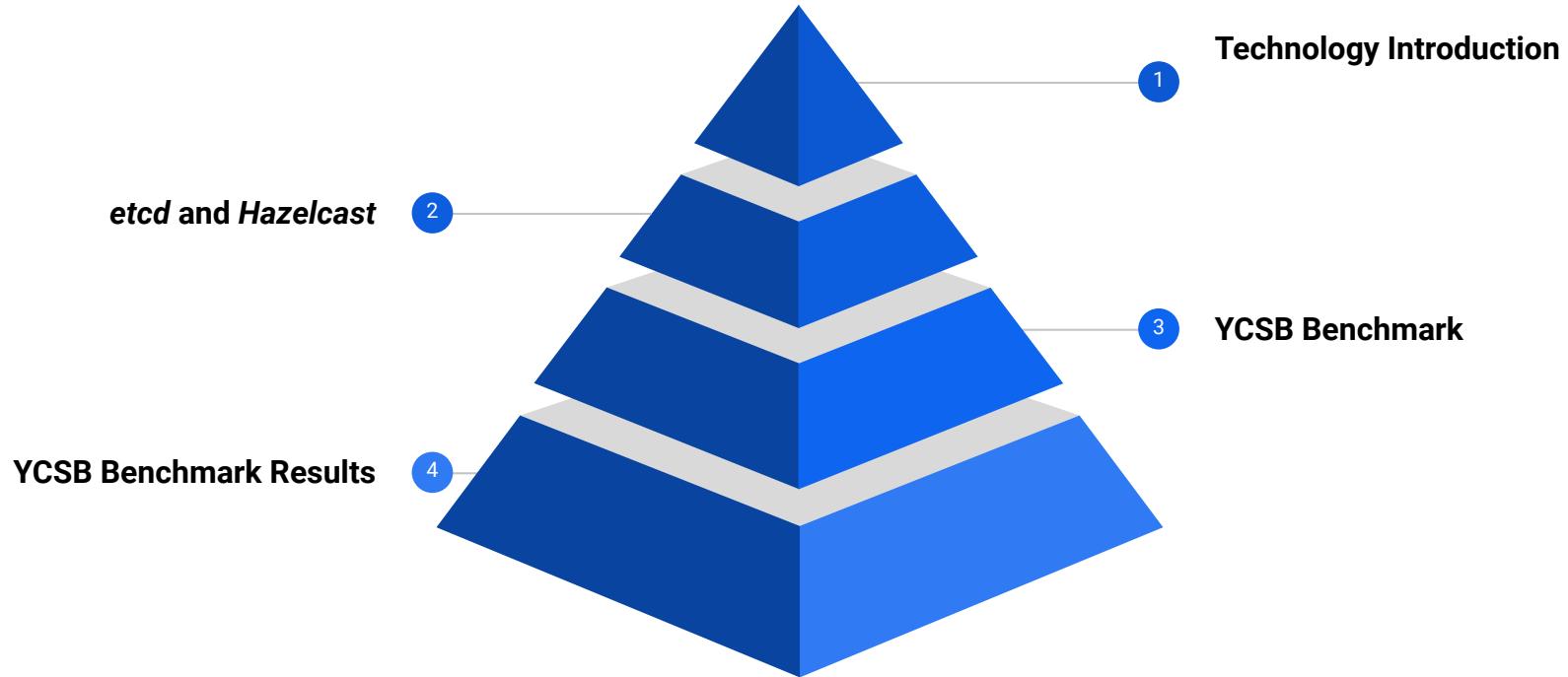


Key-value databases with *etcd* and *Hazelcast*

INFO-H415 Advance Databases
Fall 2022

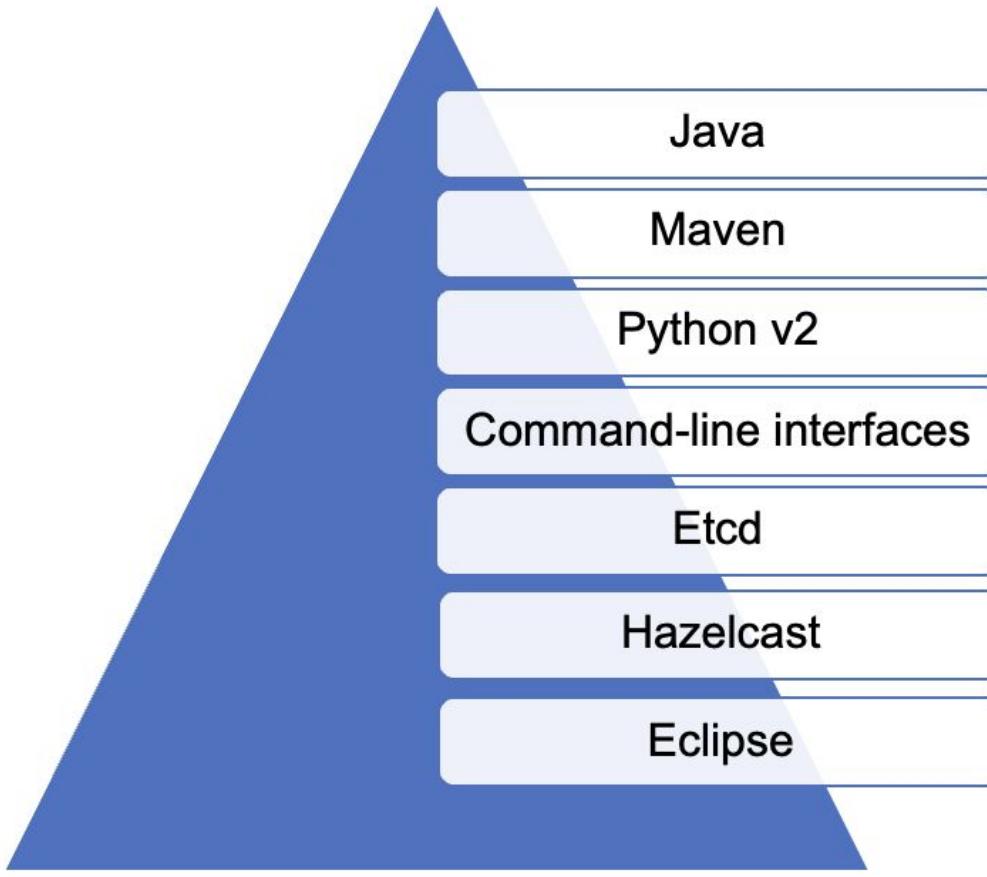
Aliakberova, Liliia,
Gepalova, Arina,
Lorencio Abril, Jose Antonio
Mayorga Llano, Mariana

Professor: Zimányi, Esteban



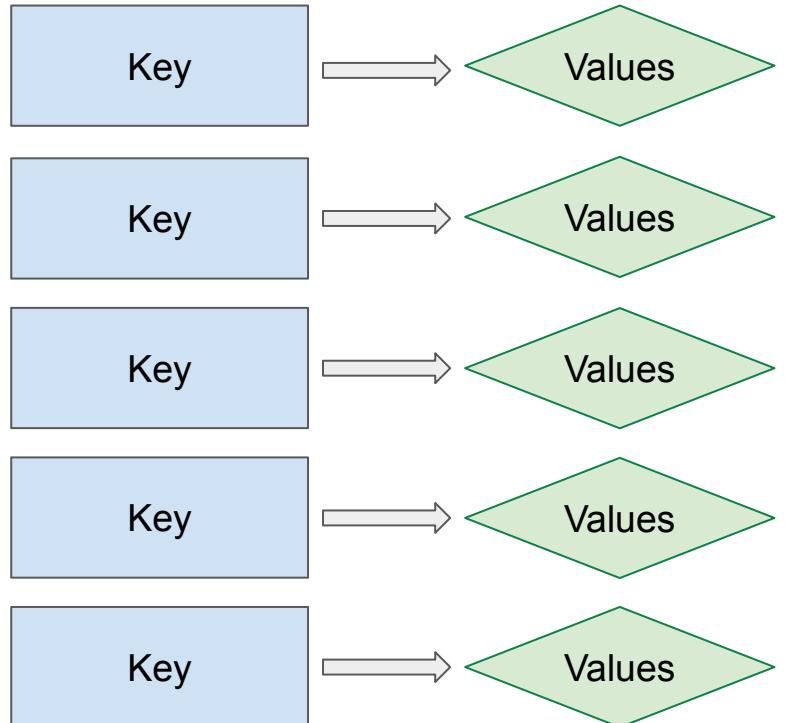
Technology Introduction

Used Tools



Features

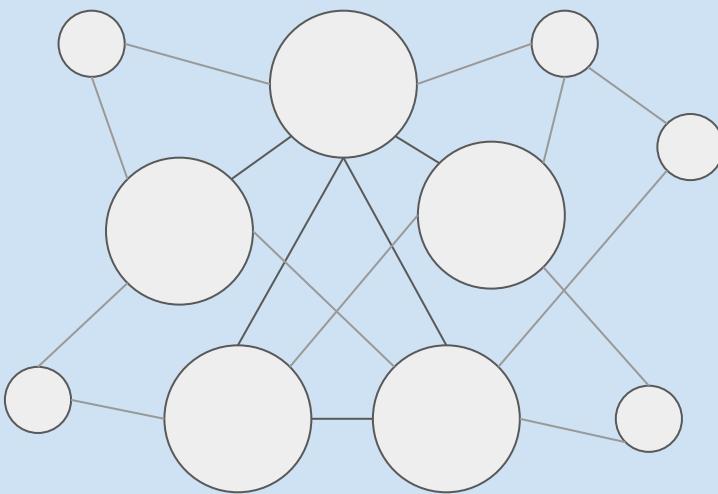
- NoSQL database
- Schemaless data model
- Flexibility
- High performance
- Etcd and Hazelcast:
distributed systems



etcd and Hazelcast



Distributed system



Independent components connected

Open-source, key-value database

Distributed systems challenge:

- **Consistent** single source of information
- **Coordination** amongst clusters

Designed to

- Store infrequently updated data
- Provide reliable watch queries

Have you heard about **RAFT** algorithm?



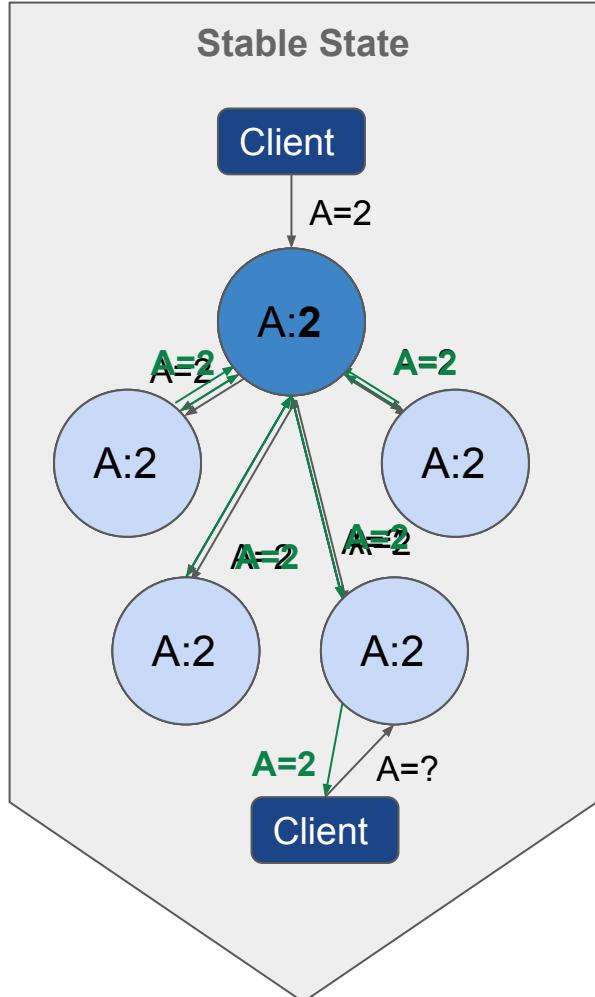
RAFT Algorithm

Roles:

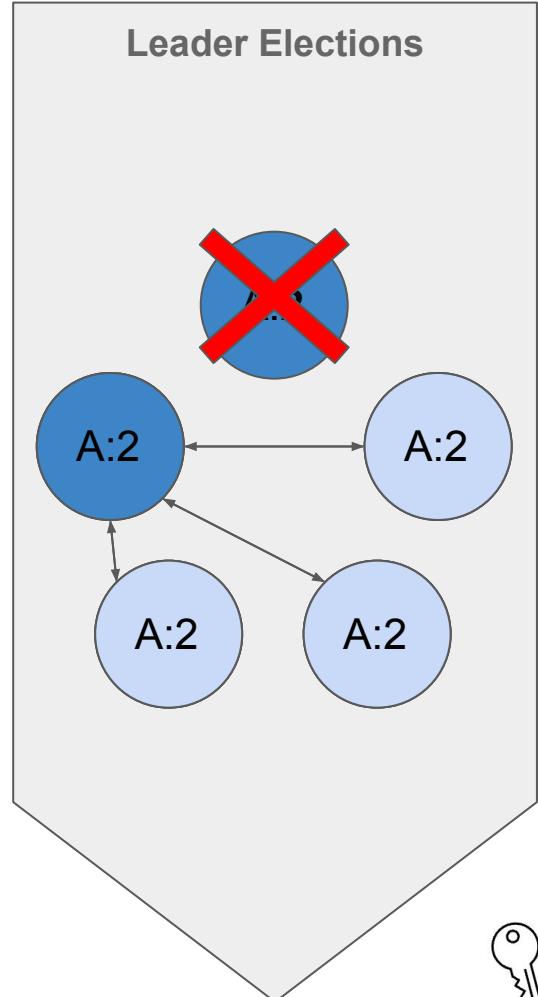
- Leader
- Follower
- Candidate
- Learner

Advantages:

- Replicated
- Consistent
- Highly available
- Fault Tolerant
- Performance



Leader Elections





Hazelcast IMDG:

- in-memory object store

Data partitioning:

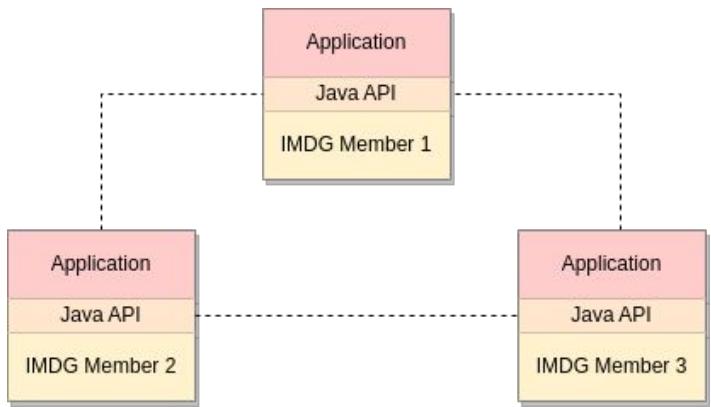
- segments

Features and advantages:

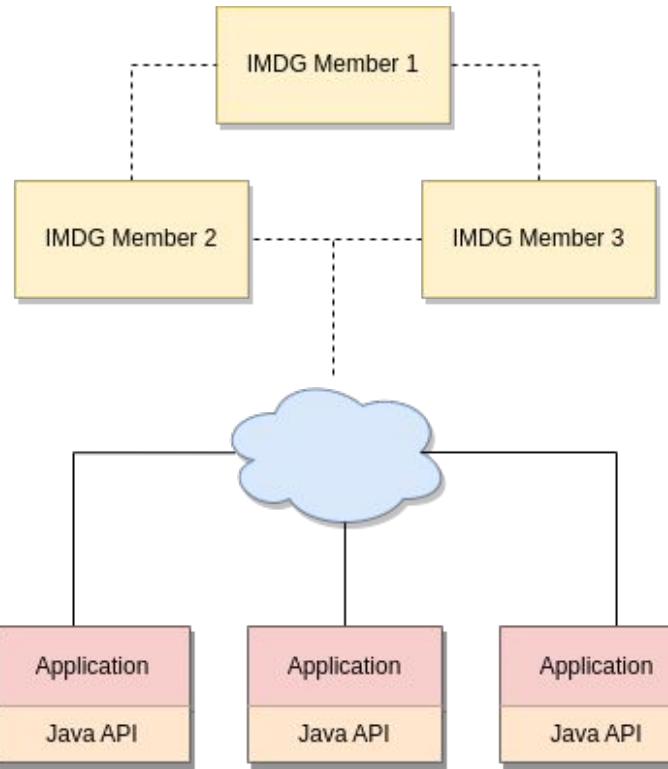
- open source
- used with a simple installation of Java
- uses in-memory storage
- peer-to-peer
- scalability can be performed dynamically
- resilient to failure
- hash-functions for keys



Hazelcast Topology



Embedded Topology



Client/Server Topology





Benchmark



Developed and published:

- By **Yahoo! Research** (2010).
- As an **open source benchmark** that could be **extended** to different technologies.

Objective:

- Standard benchmark
- Evaluation of different systems
- Focusing on read/write operations.
- Extensibility

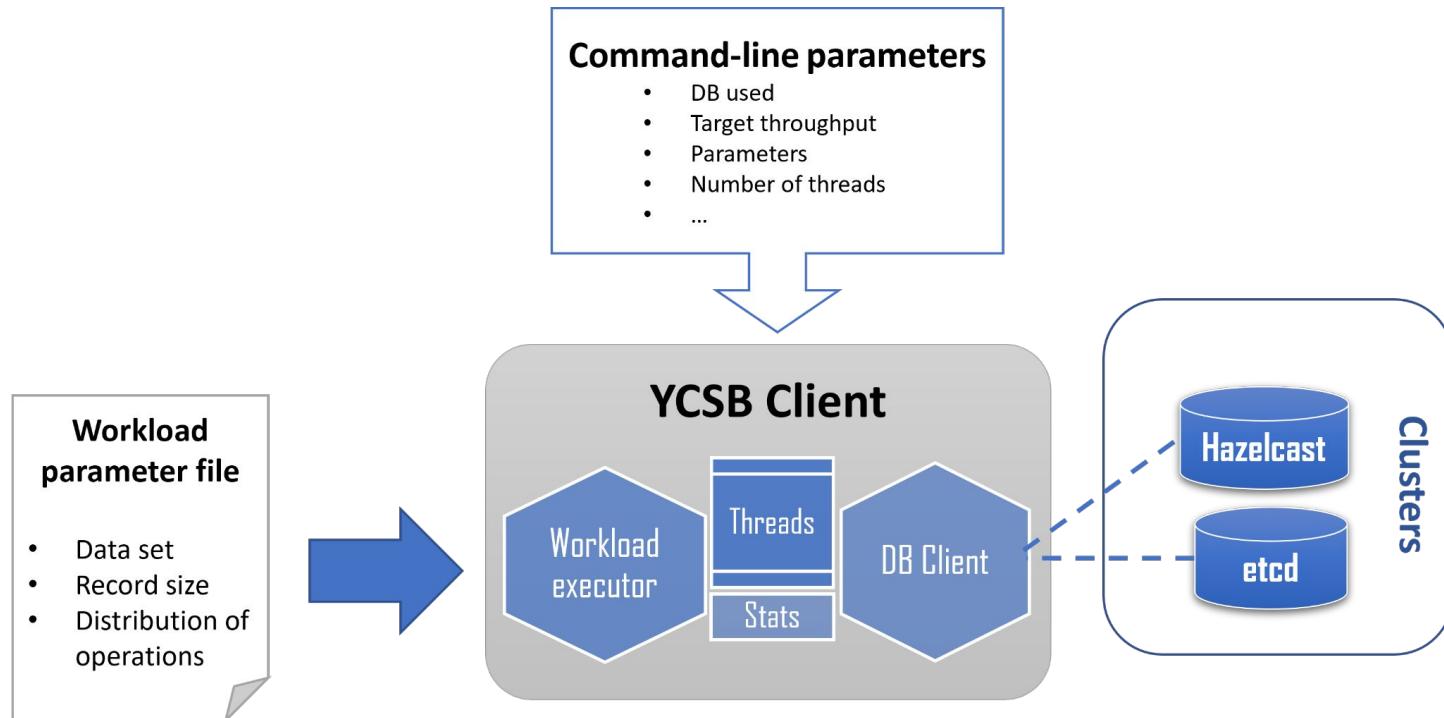
Benchmarking Cloud Serving Systems with YCSB

Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, Russell Sears

Yahoo! Research
Santa Clara, CA, USA
{cooperb,silberst,etam,ramakris,sears}@yahoo-inc.com



Architecture





Benchmark

OPERATIONS

- **Insert** a new record.
- **Update** a record replacing its value.
- **Read** either a randomly chosen field or all of them.
- **Scan** records in order, starting at a randomly chosen record key. The number of records to scan is randomly chosen.
- **Delete** a single record

Workloads

Workload A	Update heavy workload	50 % read 50% write (not necessarily the same records)
Workload B	Read mostly workload	95 % read 5 % write (not necessarily the same records)
Workload C	Read only	100% Read
Workload D	Read latest workload	5% insert 95 % read
Workload E	Short ranges	95 % scan 5 % insert
Workload F	Read-modify-write	Forces read of the records to be modified and updated

How about a **custom workload** ?



Implementation of YCSB in *etcd* and *Hazelcast*

Downloading YCSB

From the official repository.

Installation of tools

Install databases to benchmark.

Java 8.

Maven.

Python 2.

Modifications in pom files

Add pom files to define the new DB Class extensions.

Required for Maven to be able to detect the newly implemented databases and compilation.

Extending ycsb.DB class

Implement the five main functions of the DB Class interface:

- Read
- Insert
- Delete
- Update
- Scan

Adaptation of the YCSB core

We have to add some information about the new databases in the main python script of ycsb, which orchestrates the benchmarking and needs to know how to find the DB.



Use Case

Simple Process Orchestrator:

- Receives processes to execute from a **producer**.
- Have to distribute the processes between **servers**.
- Keeps track of all processes and servers.
- For this, it makes use of the **database**.

All in all:

- **42% read**
- **42% update**
- **16% insert**



```
orchestrator(KVstore db, Servers S, ProcessGen pGen, int coresPerServer):
    init(db, S, coresPerServer)
    while(pGen.generateProcess proc):
        search db to find a server with freecores>0 -> s
        core = db.read(s, freecores) # obtain core in which execute the process
        db.update(s, freecores = core-1) # reduce free cores
        search db to find first free processID -> pid
        db.insert(proc.id, localpid=pid) # associate real id to local pid
        db.update(pid, server=s, core=core, realpid=proc.id) # associate server s and core
        to process pid

        if(process_ended(Server s, Process pid, int core)):
            result = db.read(s, core)
            real_pid = db.read(pid, realpid)
            db.insert(real_pid, result=result) # Keep track of past processes outcomes
            db.update(pid, 'free') # Free localpid in db
            db.update(s, freecores = freecores + 1) # Free one core in S
            db.update(s, core, null) # Set result to null
            return (pid, result) to pGen

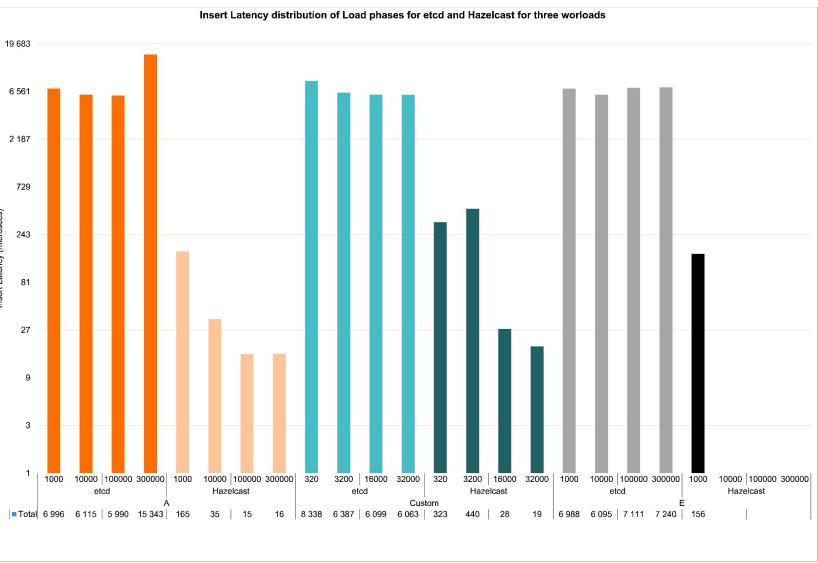
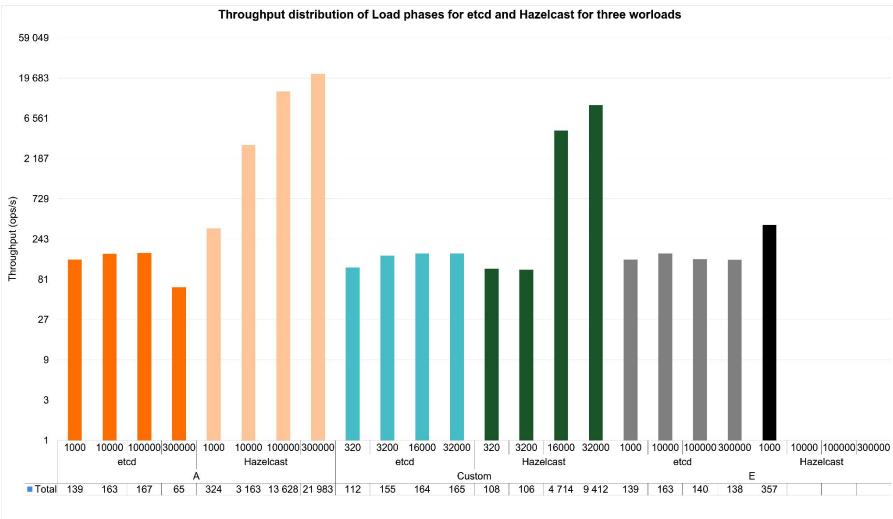
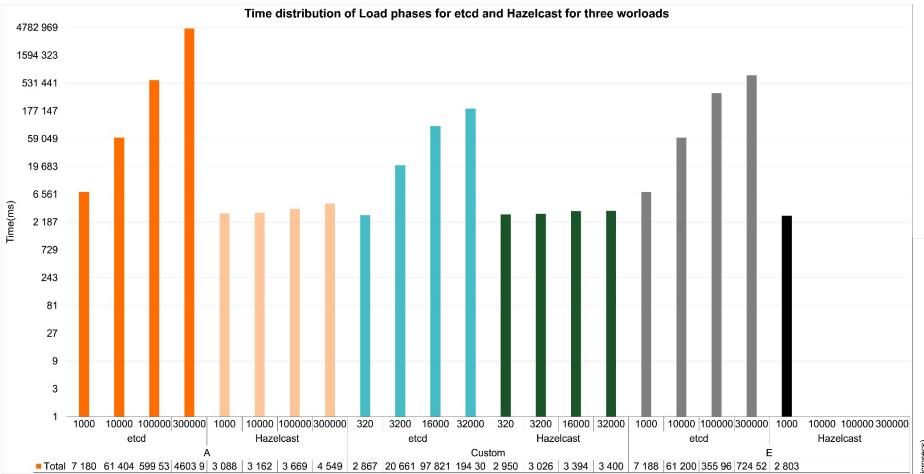
init(KVstore db, Servers S, int coresPerServer):
    for(s in S):
        db.put(s, freecores = coresPerServer)
        for(core in range(1, coresPerServer):
            db.put(s, core, null) # In this record we can introduce the result

    for(processID in range(1,coresPerServer*S):
        db.put(processID, 'free')
```



YCSB Benchmark Results

Benchmark Results

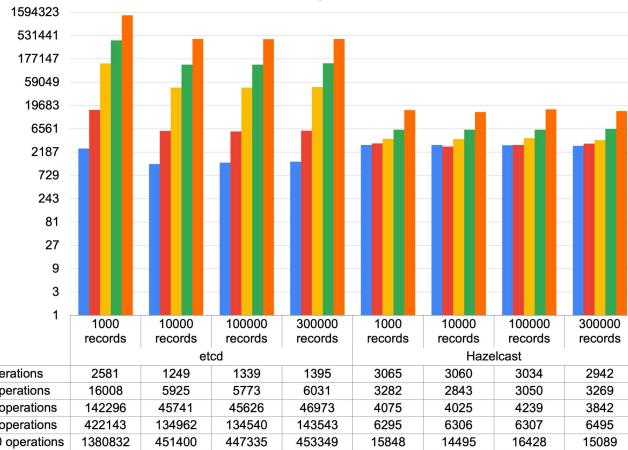


Load phase for 3 workloads

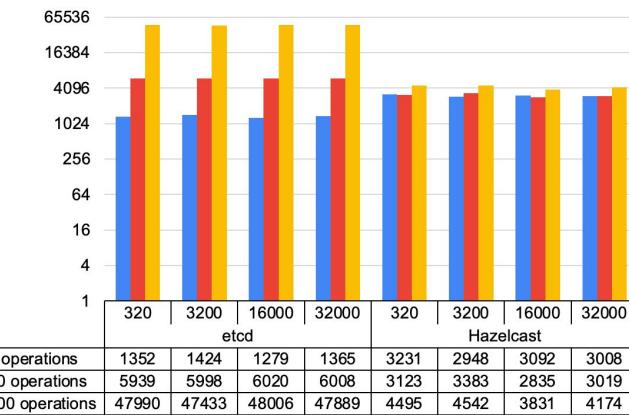


Time (ms)

Time distribution of the Workload A Run phase for etcd and Hazelcast with different number of operations



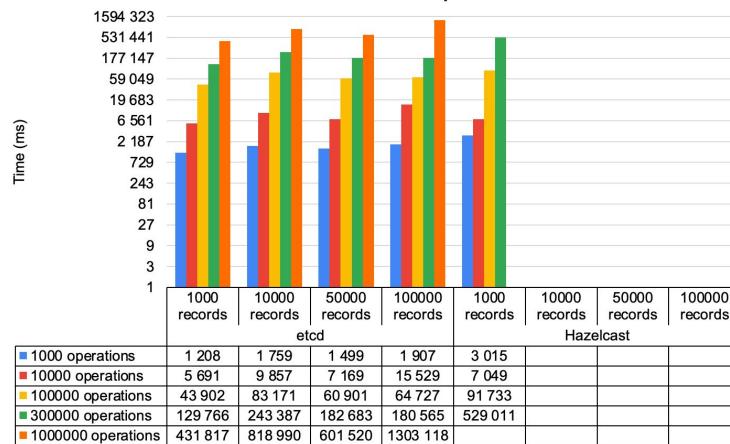
Time distribution of the Custom Workload Run phase for etcd and Hazelcast with different number of operations



Benchmark Results

Run time for 3 workloads

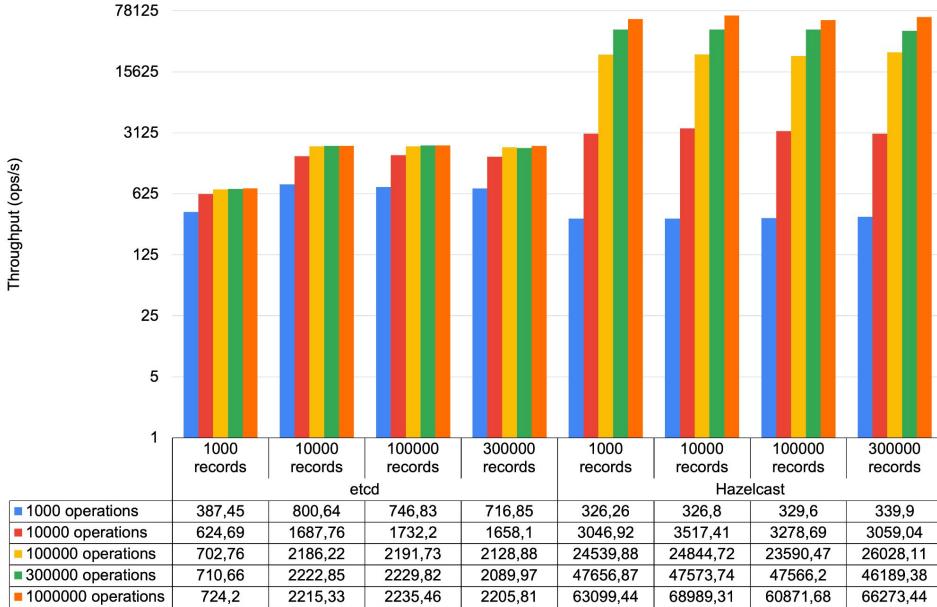
Time distribution of the Workload E Run phase for etcd and Hazelcast with different number of operations



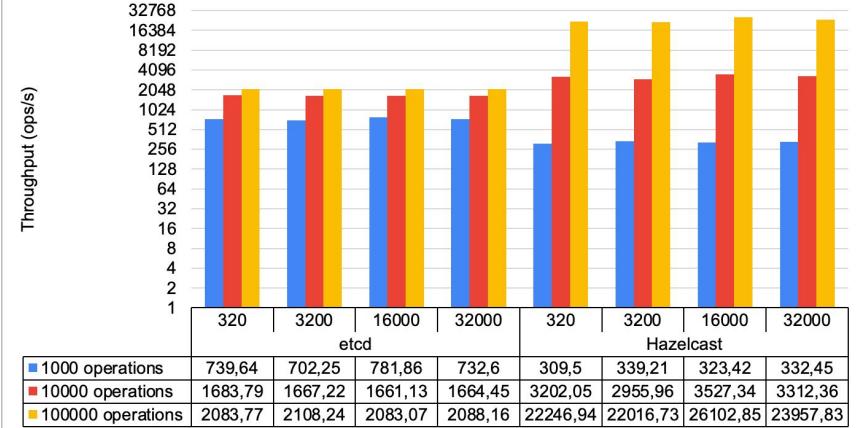
Benchmark Results

Throughput distribution of the Workload A Run phase for etcd and Hazelcast with different number of operations

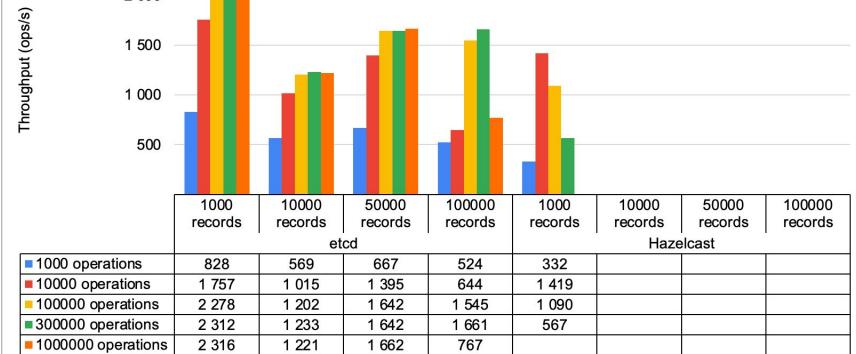
Throughput (ops/s)



Throughput distribution of the Custom Workload Run phase for etcd and Hazelcast with different number of operations



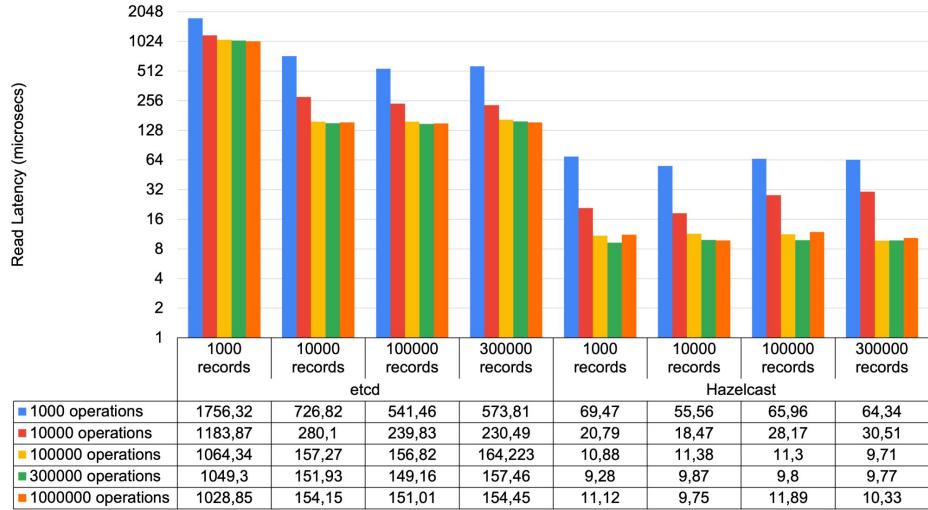
Throughput distribution of the Workload E Run phase for etcd and Hazelcast with different number of operations



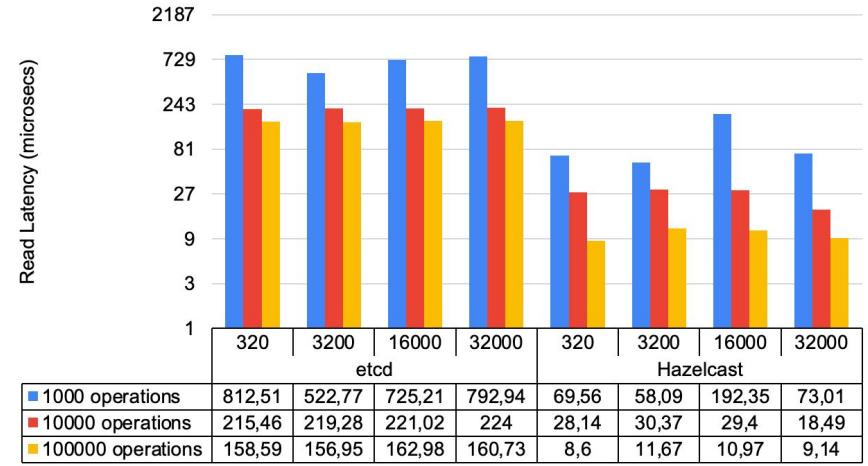
Throughput for 3 workloads

Benchmark Results

Read Latency distribution of the Workload A Run phase for etcd and Hazelcast with different number of operations



Read Latency distribution of the Custom Workload Run phase for etcd and Hazelcast with different number of operations

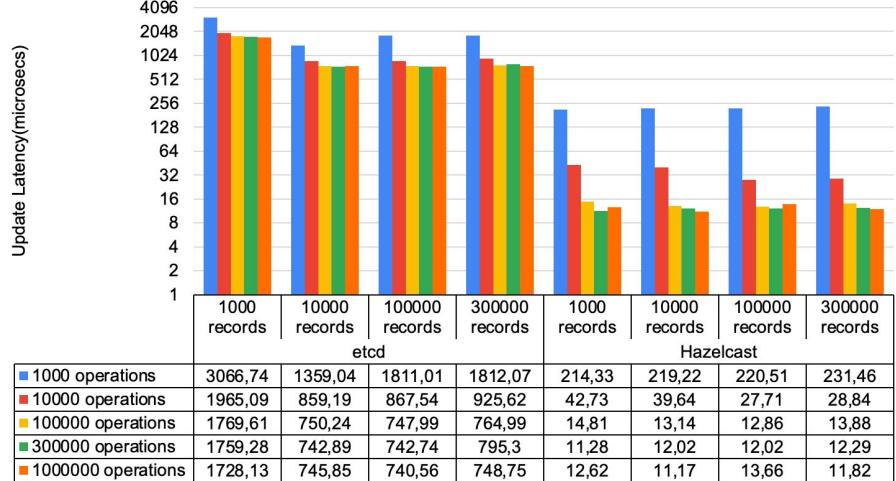


Read latency for 3 workloads

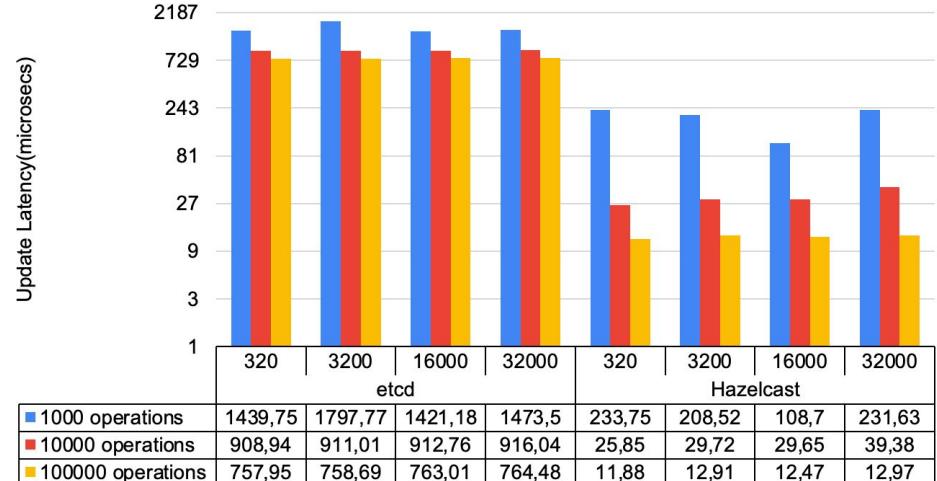


Benchmark Results

Update Latency distribution of the Workload A Run phase for etcd and Hazelcast with different number of operations



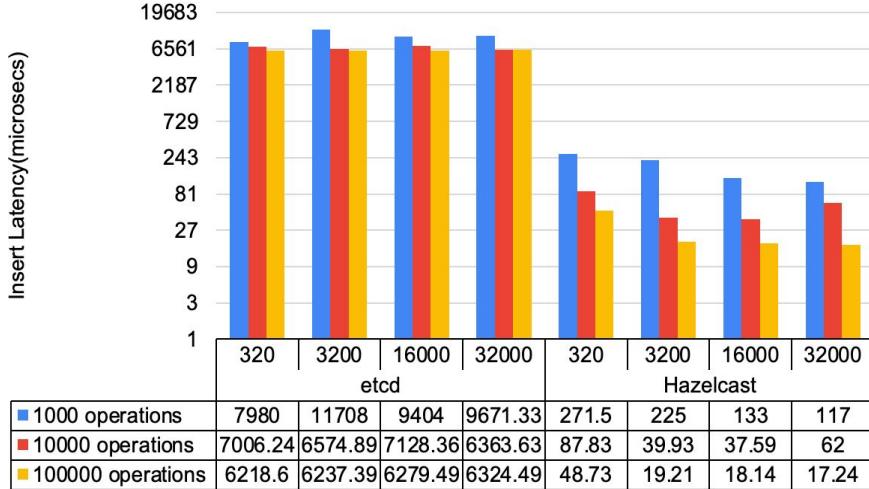
Update Latency distribution of the Custom Workload Run phase for etcd and Hazelcast with different number of operations



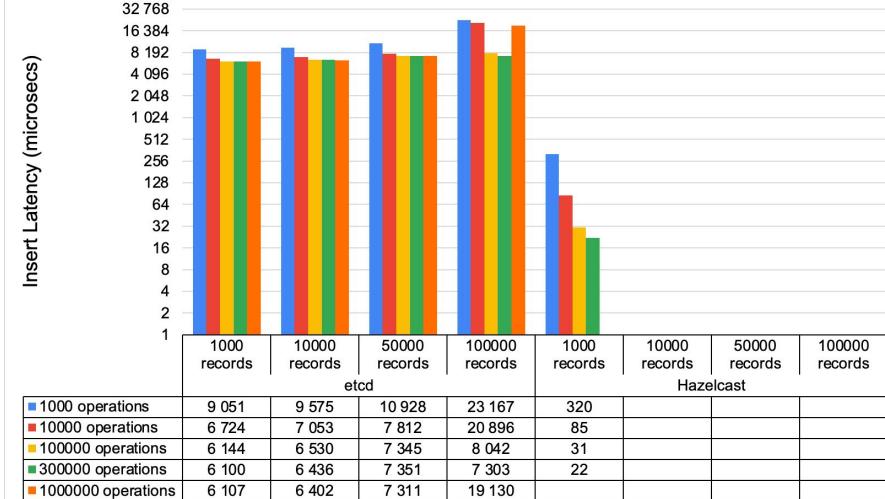
Update latency 3 workloads  

Benchmark Results

Insert Latency distribution of the Custom Workload Run phase for etcd and Hazelcast with different number of operations



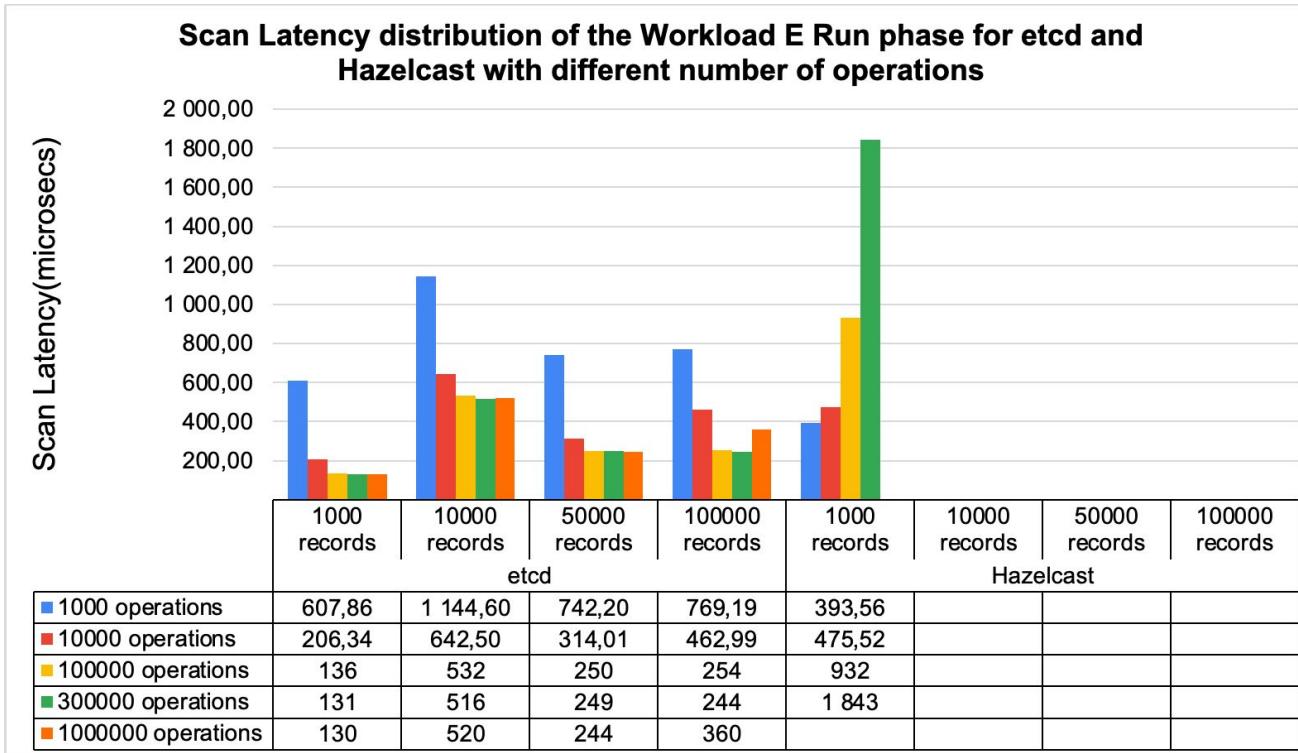
Insert Latency distribution of the Workload E Run phase for etcd and Hazelcast with different number of operations



Insert latency 3 workloads



Benchmark Results



Scan latency for 3 workloads



Conclusion: Comparison with relational model

	Key-Value Stores	Relational Model
Data model	Schema-less	Relational schema
Flexibility	Totally flexible	Not flexible
Applicability	Most use cases, but a unique key identifier is needed for all stored data	Most use cases, but performance varies according to the natural structure of the data
Data retrieval	Faster	Slower
Code	Usually cleaner	Long queries are frequently encountered
Cloud capabilities	Naturally high, CP or AP	The relational model can be extended to cope with distributed operations



Conclusion: Comparison of the two tools

	etcd	Hazelcast IMDG
Data model	Key-Value store	More general object store
Data persistency	Yes	Optional
CAP	CP, and fairly good A	AP, configurable to CP
Physical storage	On disk B+Tree	In-Memory structures: IMap for key-value stores
Index	Lexically ordered index on keys	Implicit hash index on keys
Insertion	Inefficient	Efficient
Read	Efficient	Very efficient
Update	Inefficient	Efficient

Conclusion: Comparison of the two tools

	etcd	Hazelcast IMDG
Delete	Inefficient	Efficient
Scan	Efficient	Impracticable
Targeted applications	Distributed non-changing values: <ul style="list-style-type: none">- Distributed configurations- Discovery services	Distributed highly changing values with high availability: <ul style="list-style-type: none">- Passing messages through the cluster- Storage for session data- Distributed data streaming applications
In our use case	Not a good option: fast changing data	Fairly good option: the IMap provides all necessary features needed for the use case Also, some other structures of Hazelcast could be useful (IList or IQueue)

Conclusion: YCSB

- Good tool for benchmarking cloud applications: particularly KV stores
- Outdated:
 - Many modern tools are not developed
 - Most of the tools are implemented in a very old version (e.g. PostgreSQL v9)
- Delete operation cannot be used for workloads, even though it is useful in some use cases
- The program that runs the workload is developed in Python 2, and does not work with Python 3

Thus, we think that YCSB should be revised and updated, to cope with new trends, technologies and tools!



Thanks for your attention

