**Evidence Based Big Data Benchmarking to Improve Business Performance**

# DataBench Toolbox Demo
BDV Meet-Up
27 June 2019, Riga

Tomas Pariente Lobo, Iván Martínez ATOS

**Toobox Goals & Objectives**

### Holistic benchmarking approach for big data

- The DataBench Toolbox will be a **component-based system** of both **vertical** (holistic/business/data type driven) **and horizontal** (technical area based) **big data benchmarks**. **following** the layered architecture provide by **the BDVA reference model**.

### Not reinventing the wheel, but use wheels to build a new car

- It should be able to **work** or, if possible, integrate **with existing benchmarking initiatives** and resources where possible.

### Homogenising metrics

- The Toolbox will implement ways, to emerge **Big Data benchmarking technical metrics and business insights**

### Web user interface

- It will include a web-based visualization layer to **assist to the final users to specify their benchmarking requirements** to help them to search, select, deploy, run and getting benchmarks technical results and business insights.

# Toolbox usage: General Overview

**DataBench**

## Toolbox for Benchmark providers

Big Data Benchmark
Registration/update

Benchmark
registration process,
metadata and filters

Integrating
Big Data Benchmark

Benchmark registration
of deployment &
execution process

Business Benchmark
Samples Registration

Registration of
business benchmark
and examples

## Toolbox for end users

### Toolbox for developers

| Deployment | Execution |
| Selection | Getting results |
| Recommendation | Displaying results |
| Displaying Tech. Metrics | Displaying comparatives |

Search

### Toolbox for business users

| Recommendation | Best practices |
| Business Insights | Displaying comparatives |

# Alpha version of the Toolbox already available for Alpha-testers

**DataBench**

| Searchable |
| --- |
| • HiBench |
| • SparkBench |
| • YCSB |
| • TPCx-IoT |
| • Yahoo Streaming Benchmark |
| • BigBench V2 |
| • TPC-H |
| • TPC-DS |
| • Hadoop Workload Examples |
| • PigMix |
| • Social Network Benchmark |
| • WatDiv |
| • Sanzu |
| • BigDataBench |
| • CLASS Benchmark |

| Runnable |
| --- |
| • HiBench |
| • YCSB |
| • Yahoo! streaming |
| • CLASS (in progress) |

# CherryData Use Case

- The company needs to benchmark the following low latency databases:
  - **Arango**
  - **Orient**
  - **Couchbase**
  - **Redis**

## Yahoo! Cloud Serving Benchmark (YCSB)

### Description

The YCSB framework is designed to evaluate the performance of different "key-value" and "cloud" serving systems, which do not support the ACID properties. The benchmark is open source and available on GitHub. The YCSB++ , an extension of the YCSB framework, includes many additions such as multi-tester coordination for increased load and eventual consistency measurement, multi-phase workloads to quantify the consequences of work deferment and the benefits of anticipatory configuration optimization such as B-tree pre-splitting or bulk loading, and abstract APIs for explicit incorporation of advanced features in benchmark tests.

### Reference:

https://github.com/brianfrankcooper/YCSB

### Benchmark characteristics

Micro-Benchmark  Inhouse/On-Premise  Cloud  Gigabytes  Terabytes  Petabytes  Exabytes  Fault tolerance  Execution time  Throughput  Synthetic data  Tables, files or structured data processing(OLTP)  Databases/RDBMS  NoSQL  NewSQL/In-Memory  Distributed  Interactive/near/Real-time  Volume  Execution performance  Fixed-sized records  Timeseries report

# Benchmark Provider

# Registering a benchmark in the Toolbox

# Adding configuration for benchmark deployment and run



## Preparing an Ansible Playbook

Steps:

1) Ansible template to be filled by benchmark providers
2) Upload the playbook to Toolbox Git
3) Create a job template in Ansible AWX for that playbook
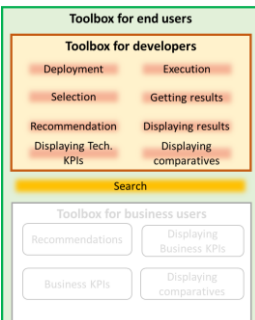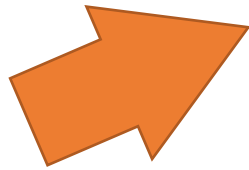4) Link the benchmark with the template so it can be run from the platform

# Technical End User

I am Paolo Ravanelli, CTO of **Cherrydata**. My company needs to benchmark the following low latency databases:
- **Arango**
- **Orient**
- **Couchbase**
- **Redis**

The company is a strong believer in benchmarking and your suggestion to use **YCSB** has already been very useful.

# Selecting and executing benchmarks: YCSB

- For benchmarks ready to run:

1. Search and choose the benchmark you want to run from the list
2. Fill in the variables with the data of the target system (i.e. host IPs)
3. Provide credentials to log into the target system (public key)
4. Let the system run the playbook (deployment and running)

# Sharing results after executing YCSB

# Visualizing Benchmark Results

- Ongoing work (for the Beta version): Investigating visual paradigms to homogenize and show the results of the a given run, comparison with other runs or with other benchmarks...

# Visualizing YCSB execution results

# Summary

- Next Toolbox releases
  - Beta Toolbox by December 2019
  - Final release by June 2020

- Generation of a Benchmarking Knowledge Graph supporting technical and business aspects

- Find relations between technical metrics and business insights based in use cases

# More info

- Check our website: https://www.databench.eu/

# Contacts

✉ info@databench.eu

🐦 @DataBench_eu

f DataBench

in DataBench Project

DataBench

▶ DataBench Project

**DataBench**

Evidence Based Big Data Benchmarking to Improve Business Performance

tomas.parientelobo@atos.net

ivan.martinez@atos.net