

# ScODA: SUPERCOMPUTING OPERATIONAL DATA ANALYTICS DATABASE BENCHMARKING

An Emerging Pipeline for Evaluating Distributed Database Performance to Support Operational Data Analytics

Nicholas M. Synovic and George K. Thiruvathukal, Department of Computer Science Loyola University Chicago  
Shilpika, Silvio Rizzi, Doug Waldron, and Michael E. Papka, Leadership Computing Facility, Argonne National Laboratory

## INTRODUCTION

- ScODA (Supercomputing Operational Data Analytics) is an emerging pipeline to evaluate distributed database management solutions (DBMSs) when concurrently reading and writing data from exascale supercomputing systems.
- Evaluating and understanding the performance of distributed DBMSs under different workloads is critical as business intelligence, operations, and high-performance computing (HPC) research teams leverage this data in unique and diverse projects [1, 2].
- The concurrent ingress and egress of Aurora logs consumes ~70% of system resources on currently implemented DBMS systems, and as we look to future exascale systems, the amount of compute resources to store, query, and manage the data is expected to increase.
- ScODA is an emerging pipeline to benchmark relational, document, and time-series DBMSs and lakehouses on operational data analytic (ODA) workflows.
- Projects and workflows that benefit from a performant distributed DBMS include:
  - Predictive and rapid maintenance of supercomputing resources,
  - Business operations and justification for future supercomputing resources,
  - Analysis of hardware failures and use-case patterns, and
  - Environmental System design, engineering, and budgeting

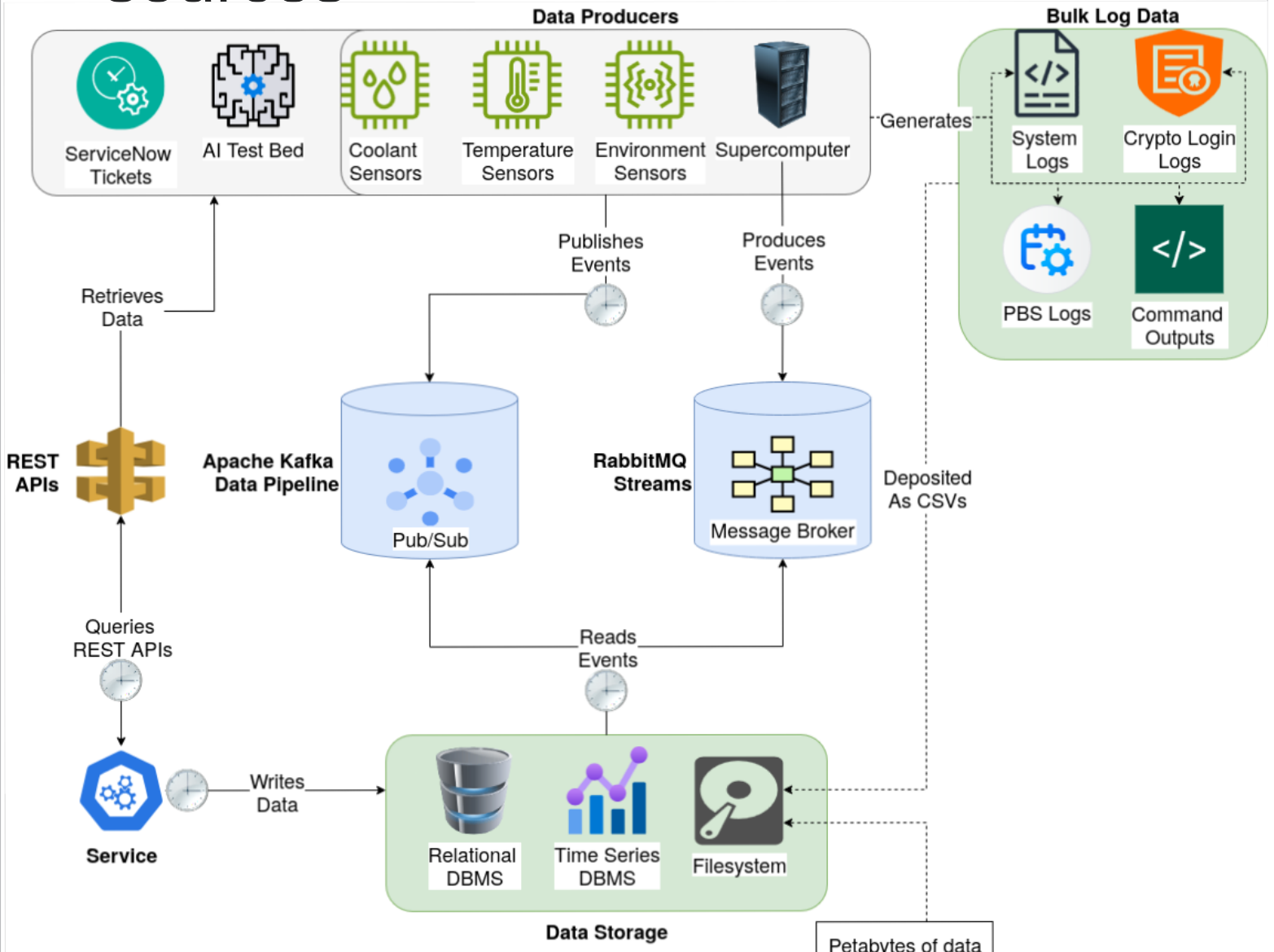
## GOALS

- Benchmark distributed relational, document, time-series, and lakehouse data solutions with environment log data.
- Use existing ODA workflows for benchmarking DBMS including:
  - Independent ingress and egress,
  - Concurrent ingress and egress, and
  - Data aggregation methods

DBMS Name	Data Model	Data Stored
PostgreSQL	Relational	Env. + Sys. Logs
MariaDB	Relational	Env. + Sys. Logs
MySQL	Relational	Env. + Sys. Logs
CouchDB	Document	Env. Logs
MongoDB	Document	Env. Logs
InfluxDB [ANL]	Time series	Env. Logs
VictoriaMetrics [ANL]	Time series	Env. Logs
Apache Iceberg	Lakehouse	Bulk Sys. Logs
Delta Lake [ANL]	Lakehouse	Bulk Sys. Logs

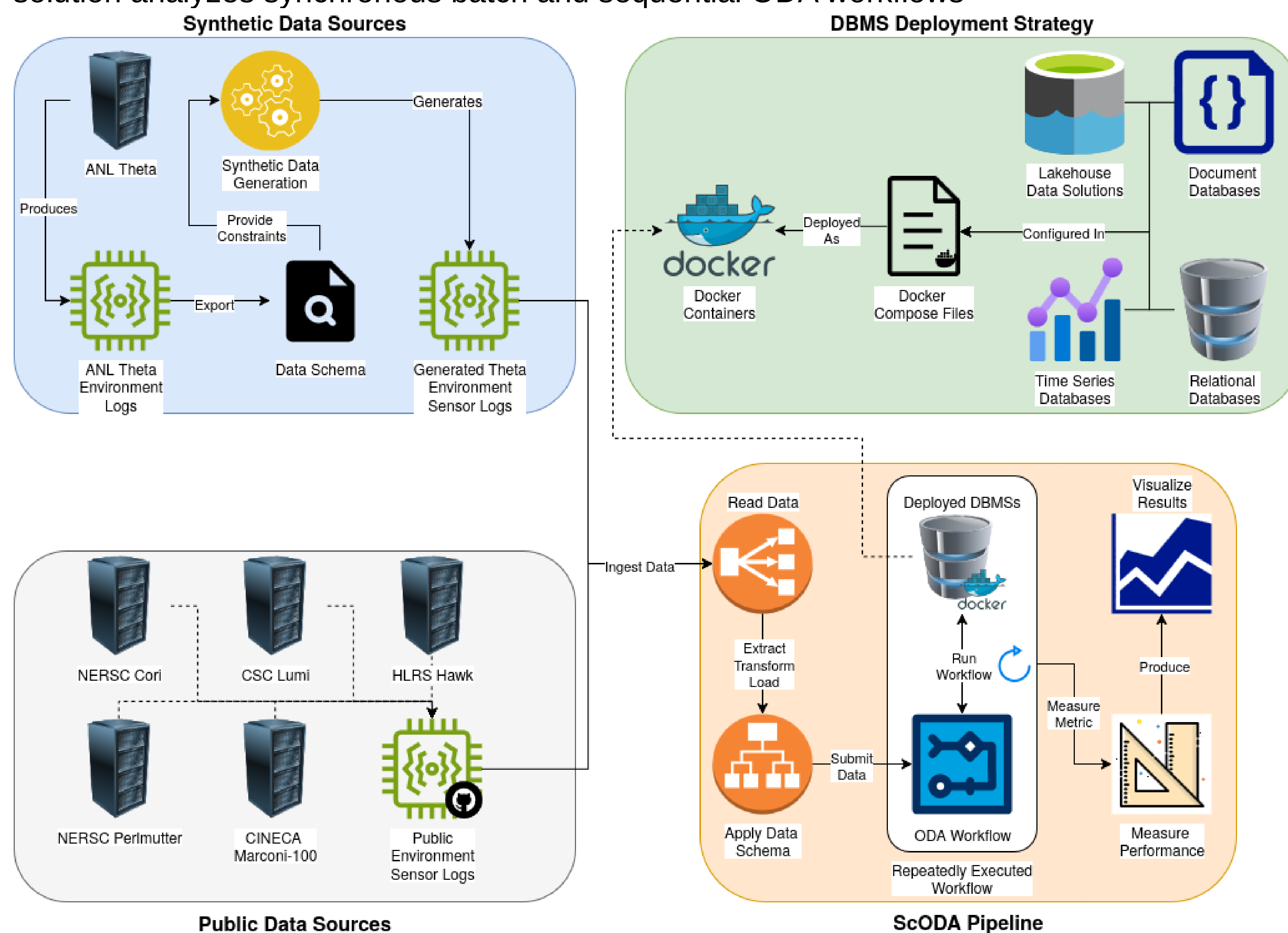
## METHODOLOGY

- Meet with stakeholders and practitioners to understand the current data pipeline
- Model the data pipeline, implement ODA workflows, and evaluate the performance of DBMSs leveraging publicly available and synthetic data sources



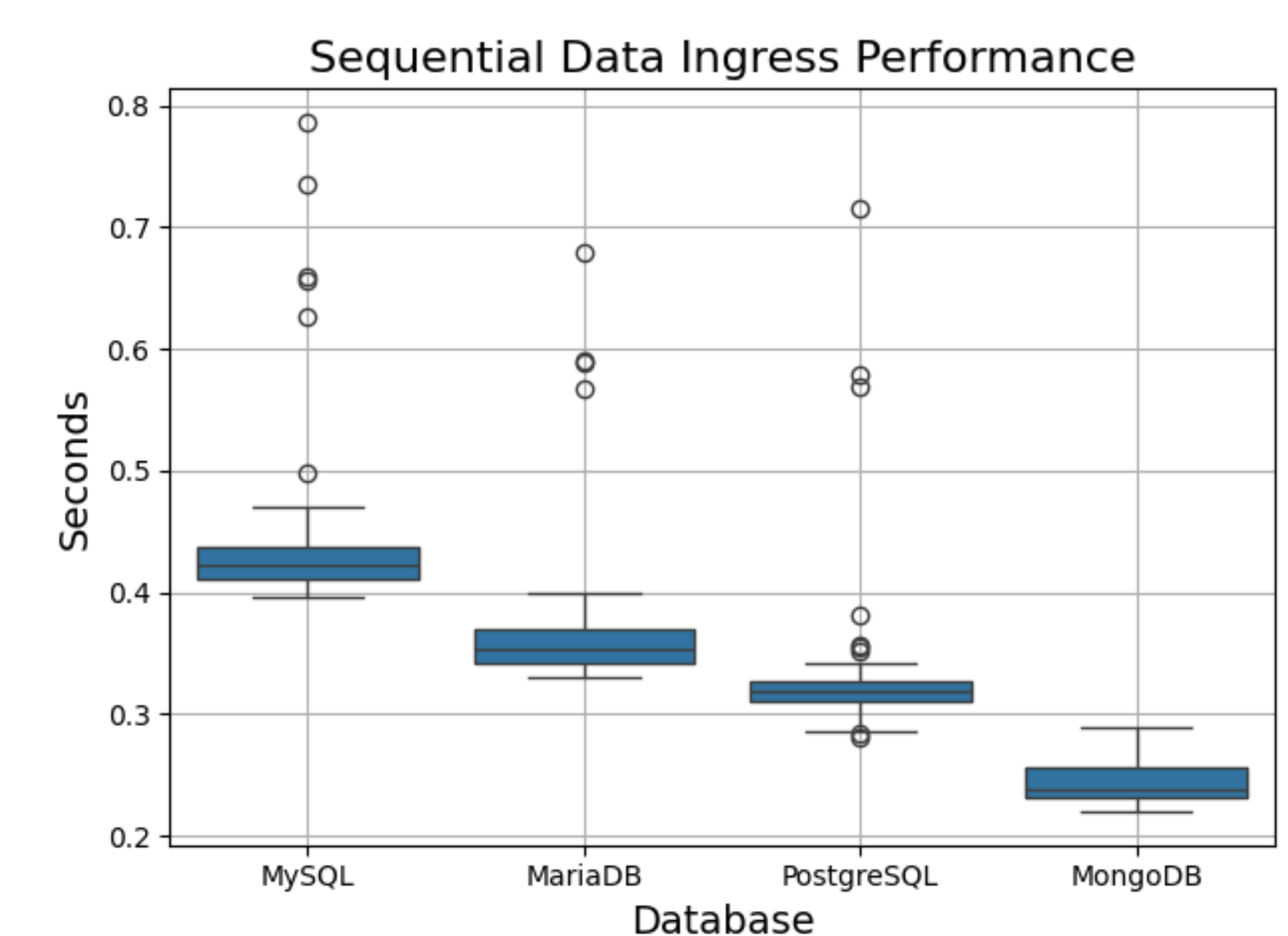
## PIPELINE SOLUTION

- We leverage publicly available environment log data [3] and generate synthetic data from log schemas
- DBMSs are deployed as Docker containers orchestrated with Docker Compose
- Our solution analyzes synchronous batch and sequential ODA workflows



## PRELIMINARY RESULTS

- DBMSs were deployed as Docker containers.
- We benchmark the DBMSs regarding data ingress of publicly available and synthetic supercomputer environment logs.
- We report the top four (4) performing DBMSs regarding their performance below on a subset of the publicly available log data.



## CONCLUSIONS

- We have benchmarked three (3) relational, two (2) document, and two (2) time-series DBMSs, and two (2) lakehouse data solutions with respect to their performance on identified ODA workflows
- On data ingress, document databases outperform relational databases on both sequential and batch workflows
- We have released all of our source code and methodology on GitHub available at this QR code:



## NEXT STEPS

- Add support for proprietary DBMSs (e.g., IBM DB2)
- Support ODA workflows benchmarks including:
  - In-flight data compression and decompression,
  - Machine and deep learning, and
  - Real-time data streaming and visualization
- Enable a hyperparameter-style search over DBMS configuration spaces to identify optimal settings

## REFERENCES

- B. Lenard, *et. al.*, "An Approach for Efficient Processing of Machine Operational Data," in Database and Expert Systems Applications.
- E. Pershey *et. al.*, "A Big Data Approach for Efficient Processing of Machine Operational Data," in Proceedings of the 37th International Conference on Scalable Scientific Data Management.
- T. Patki, *et. al.*, "A Global Perspective on Supercomputer Power Provisioning: Case Studies from United States and Europe," in Proceedings of the 37th International Conference on Supercomputing