

Benefits of Time Series Data Tables for HPCInfo

Kenton Romero | kromero2001@tamu.edu | Texas A&M University

Background

What is HPCInfo?

HPCInfo is a user facing web application for LANL HPC users to view Slurm, cluster, and storage data. It uses a self hosted grafana instance with a MariaDB backend. Much of the information shown on HPCInfo is displayed via time windows. These queries have potential to be sped up with an underlying data engine that is more closely aligned, such as a time series database.

Why TimescaleDB?

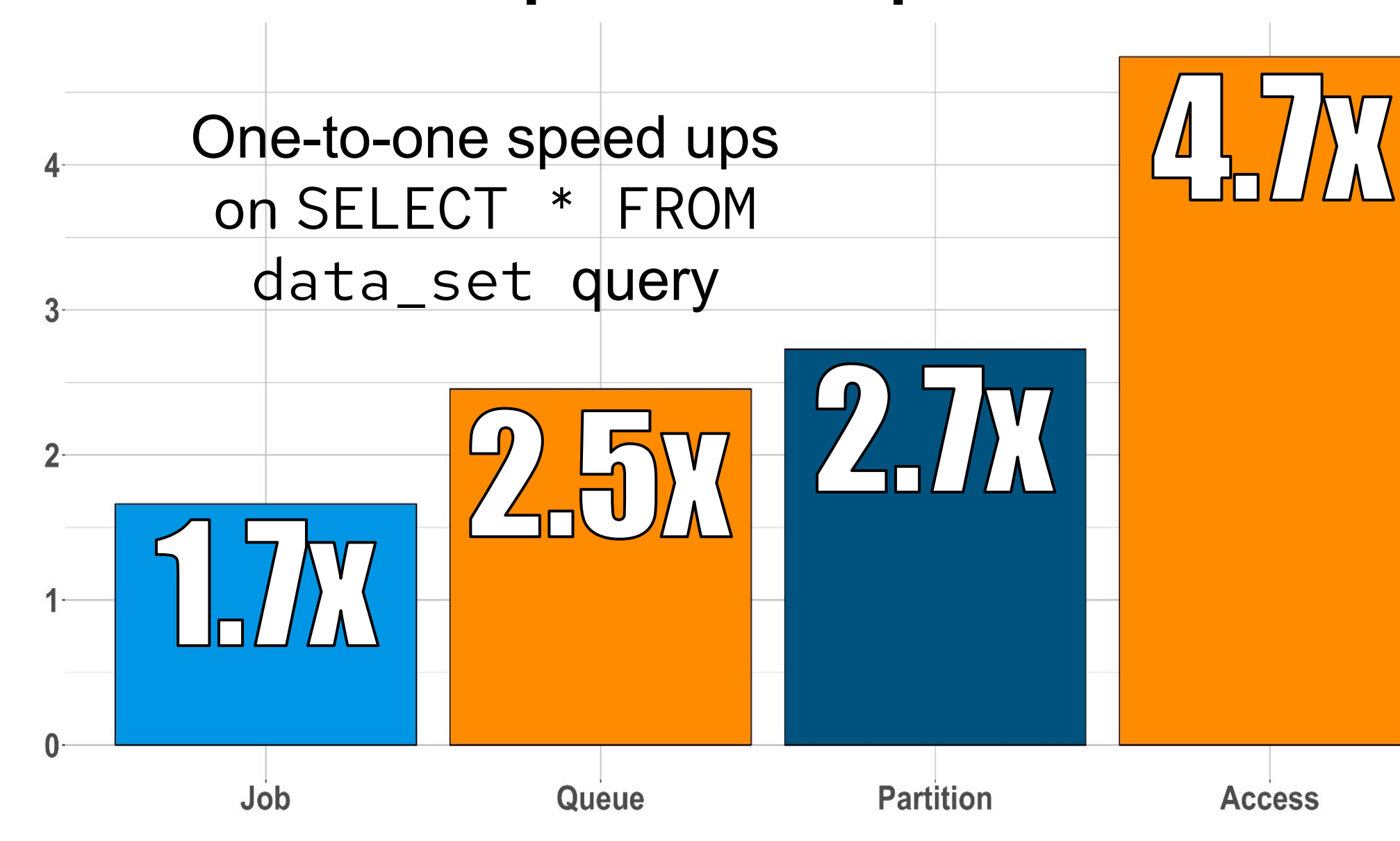
TimescaleDB is a PostgreSQL extension that adds support for time series data tables. TimescaleDB has several advantages compared to other time series databases that made it appealing for this effort.

- SQL syntax reduces changes to HPCInfo queries and code.
- "Hypertables" are enabled per table (not per database)
- Hyperfunctions optimize complex queries on time series data
- Simple migration plan and ease of administration

Data Sets

| Name | frequency | Description |
|-----------|-------------------------------|-------------------------------|
| Job | every time a job is submitted | Grizzly Slurm Job Information |
| Queue | every 5m per cluster | Slurm queue information |
| Partition | every 5m per partition | Partition node status |
| Access | every 5s per login node | User cluster access |

Speed Ups



Experimental Results

Hyperfunctions

TimescaleDB supports hyperfunctions, specialized functions that optimize the analysis of time series data. Average partition usage is calculated with both hyperfunction and SQL syntax and speed ups to the relational SQL query are shown below.



Time Buckets

The time bucket hyperfunction aggregates a number of data rows to represent a bucket of data points as a specific data row. Speed ups are shown below.

In-Place Improvements

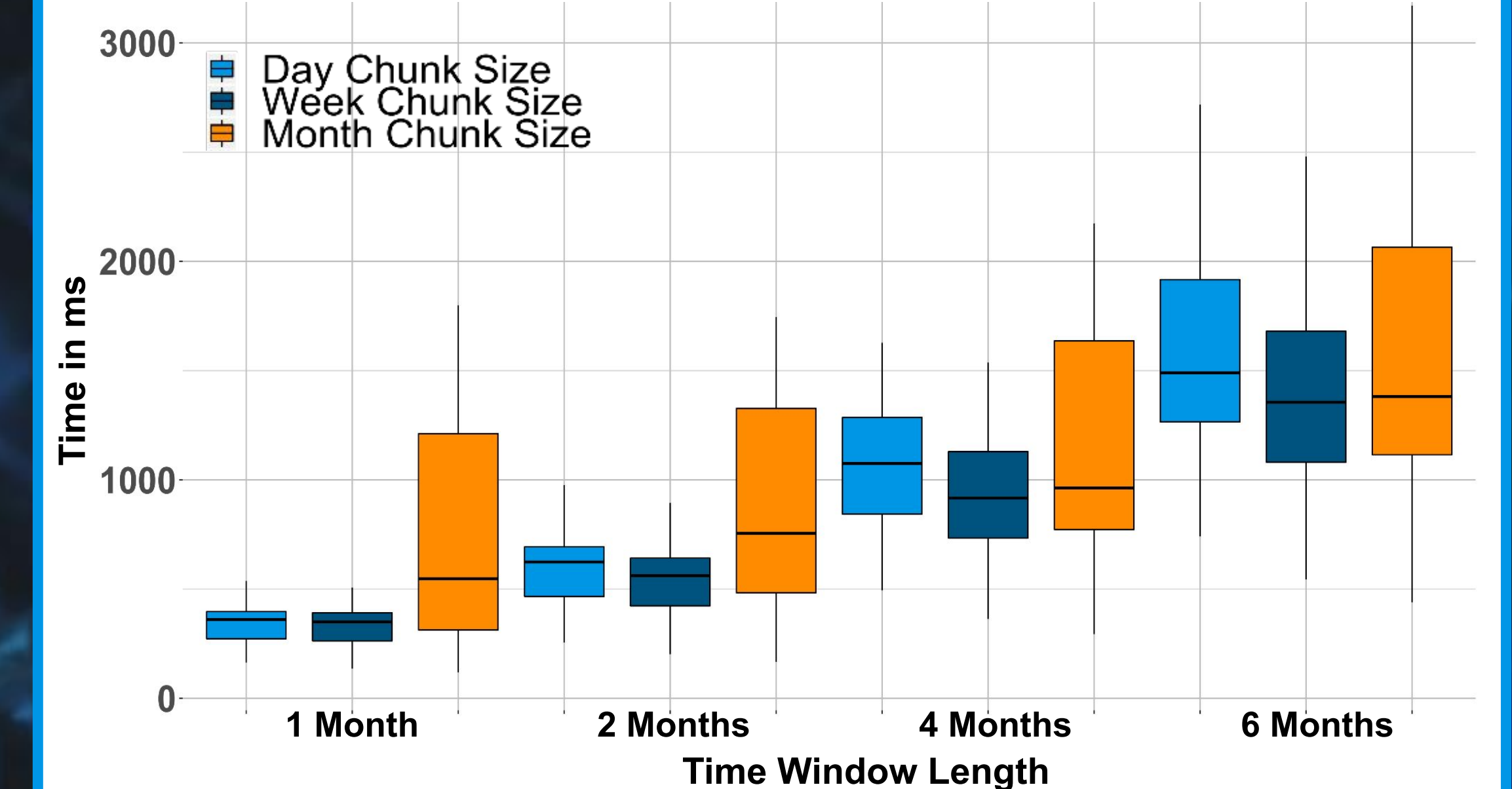
- **10% less** storage with original indexes
- **24% less** storage with reduced indexes
- Uniform speed ups on all data sets
- Increased consistency of query times

Chunk Interval

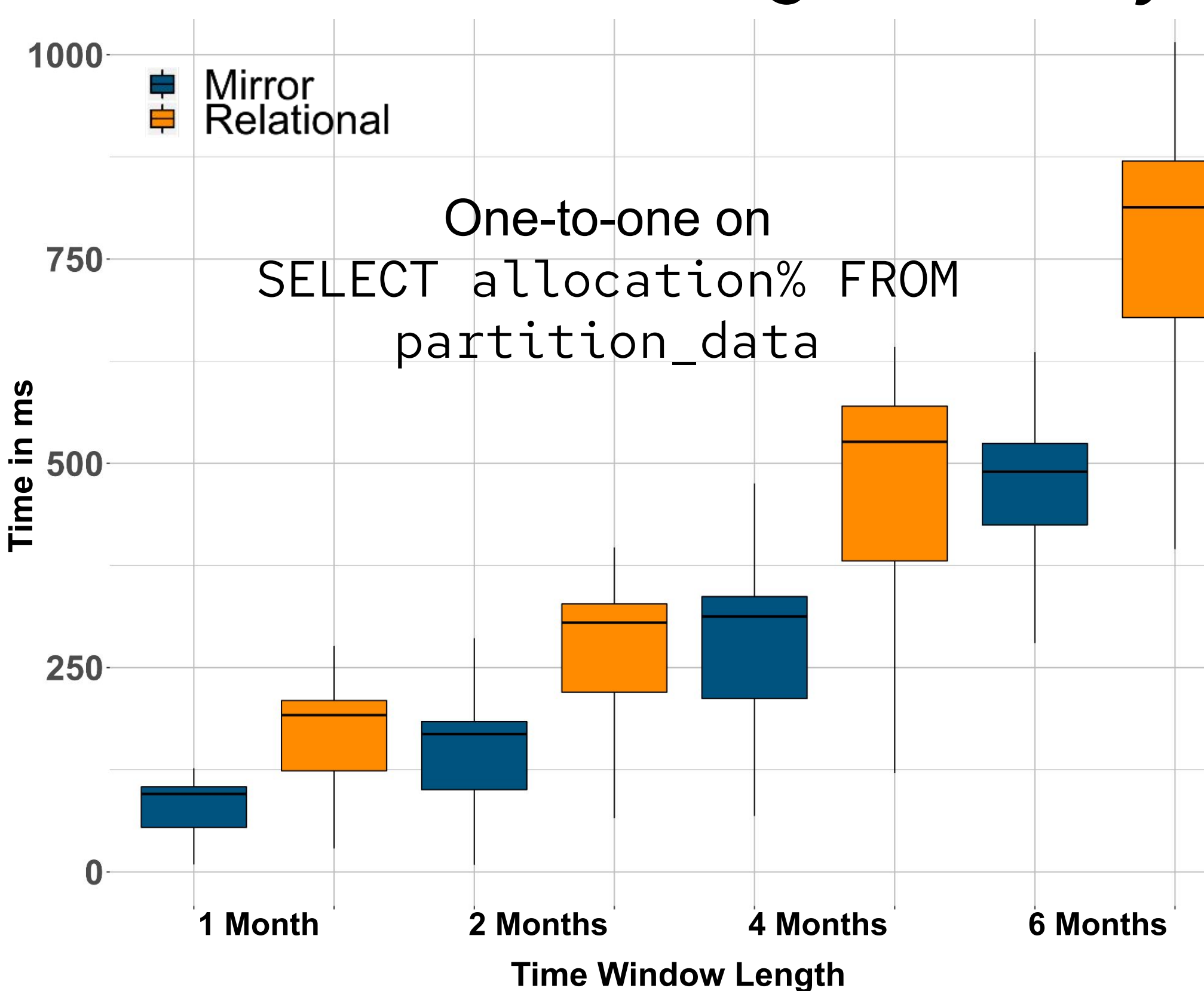
Chunk interval determines the smallest subset of data that a TimescaleDB query will iterate over. Testing shows performance of multiple chunk intervals with SELECT * FROM access_data queries over multiple lengths of time windows.

- Adjusted per table
- Default is 7 day chunk sizes

Chunk Size Performance



Partition Usage Query



Methodology

| Software | Version |
|--------------------------|---------|
| TimescaleDB | 2.11.0 |
| TimescaleDB Toolkit | 1.16.0 |
| PostgreSQL | 14.8 |
| Red Hat Enterprise Linux | 8.8 |

Database Transfer & Benchmark Workflow

Scripts have been developed to smooth out the transition from MariaDB to PostgreSQL. Tables of interest are imported into "Hypertables" with identical indexes to HPCInfo. Benchmarking clears cache between every query to avoid caching of hypertable chunks. All queries are over a random time window of fixed length.

Conclusions

Overall, these results show that TimescaleDB is well suited for our applications with HPCInfo. Here are the following outcomes for HPCInfo from this research:

- Developed tools will ease the transition of MariaDB to PostgreSQL
- Demonstrated performance of SQL queries in TimescaleDB, allowing for minimal code and query changes to HPCInfo
- Reduced potential storage footprint
- Determined the optimal chunk size of 7 days for HPCInfo Queries
- Found hyperfunctions are unnecessary for simple queries used by HPCInfo

Special Thanks to my mentor, Hunter Easterday



References

Photo by [Pietro Jeng](https://unsplash.com/de/@pietrozj?utm_source=unsplash&utm_medium=referral&utm_content=creditCopyText) on [Unsplash](https://unsplash.com/photos/n6B49lTx7NM?utm_source=unsplash&utm_medium=referral&utm_content=creditCopyText)

[Cheddar icons](https://www.flaticon.com/free-icons/cheddar "cheddar icons") created by berkahicon - Flaticon