# Data Science Infrastructure (DSI) Project

Giovanni Michel
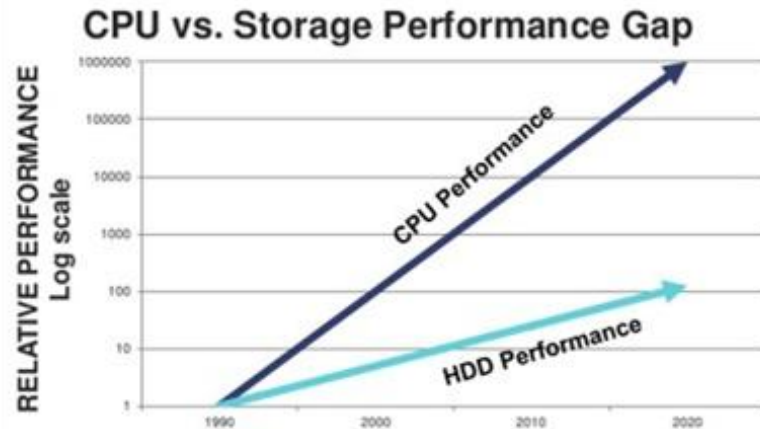
08/04/2022

Mentors                        Jesus Pulido
Terry Turton

Managed by Triad National Security, LLC., for the U.S. Department of Energy's NNSA.

1

# Motivation for DSI

1. Inspired by the time to gather and run simulations.
2. Storing simulation and environment metadata could help to answer LANL mission critical questions.
3. Increasing FLOPS doesn't always help answer scientific questions – only gives you more data that you can't visualize immediately.

## CPU vs. Storage Performance Gap

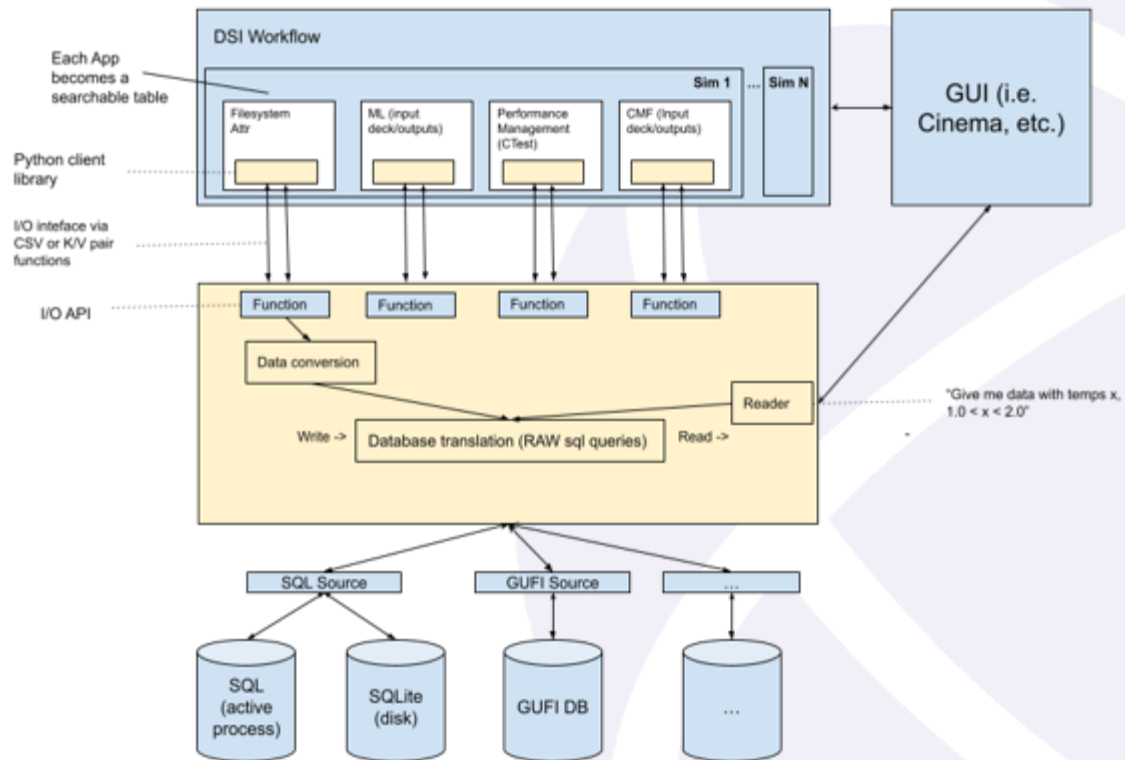# Challenges of Data Storage in a Compute-Centric World

1. POSIX file system permission for access security.
2. Interactivity and Querying is not available, which reduces opportunity for explainability.
3. Permanent data store does not exist. The current mainstream relational databases provide no extensibility to users
4. Lack of requirement & place for storing metadata in current workflows.
5. Experiments could take days, weeks, and months! Extremely computationally expensive to run simulation and reproduce the simulation itself.

# Outline

- The Data Science Infrastructure (DSI) project is focused on providing data management workflows for existing LANL applications such as the Common Model Framework (CMF) via databases (SQL) through defined schema. This approach to data management uses many forms of metadata (simulation, filesystem, performance metrics) to populate a searchable database to enable capabilities for analysis, exploration and visualization.
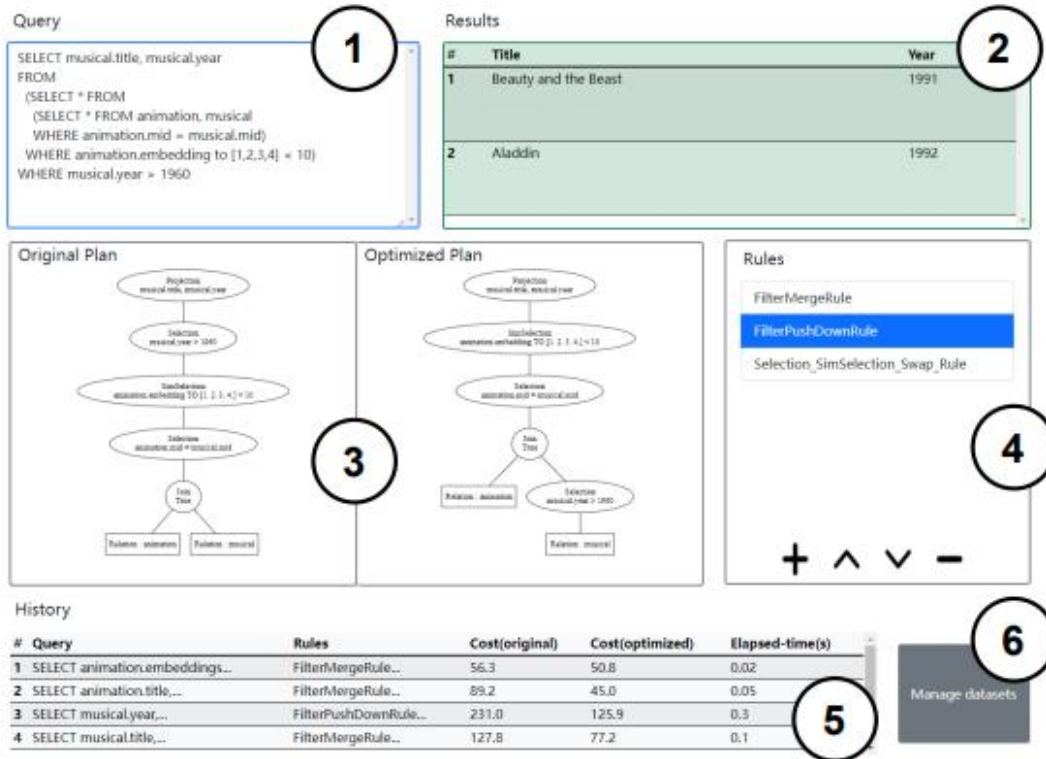
DSI Workflow

Each App becomes a searchable table

Sim 1   Sim N

Filesystem Attr

ML (input deck/outputs)

Performance Management (CTest)

CMF (Input deck/outputs)

Python client library

I/O inteface via CSV or K/V pair functions

I/O API

Function   Function   Function   Function

Data conversion

Reader

Write ->   Database translation (RAW sql queries)   Read ->

GUI (i.e. Cinema, etc.)

"Give me data with temps x, 1.0 < x < 2.0"

SQL Source   GUFI Source   ...

SQL (active process)

SQLite (disk)

GUFI DB

...

# Related Works

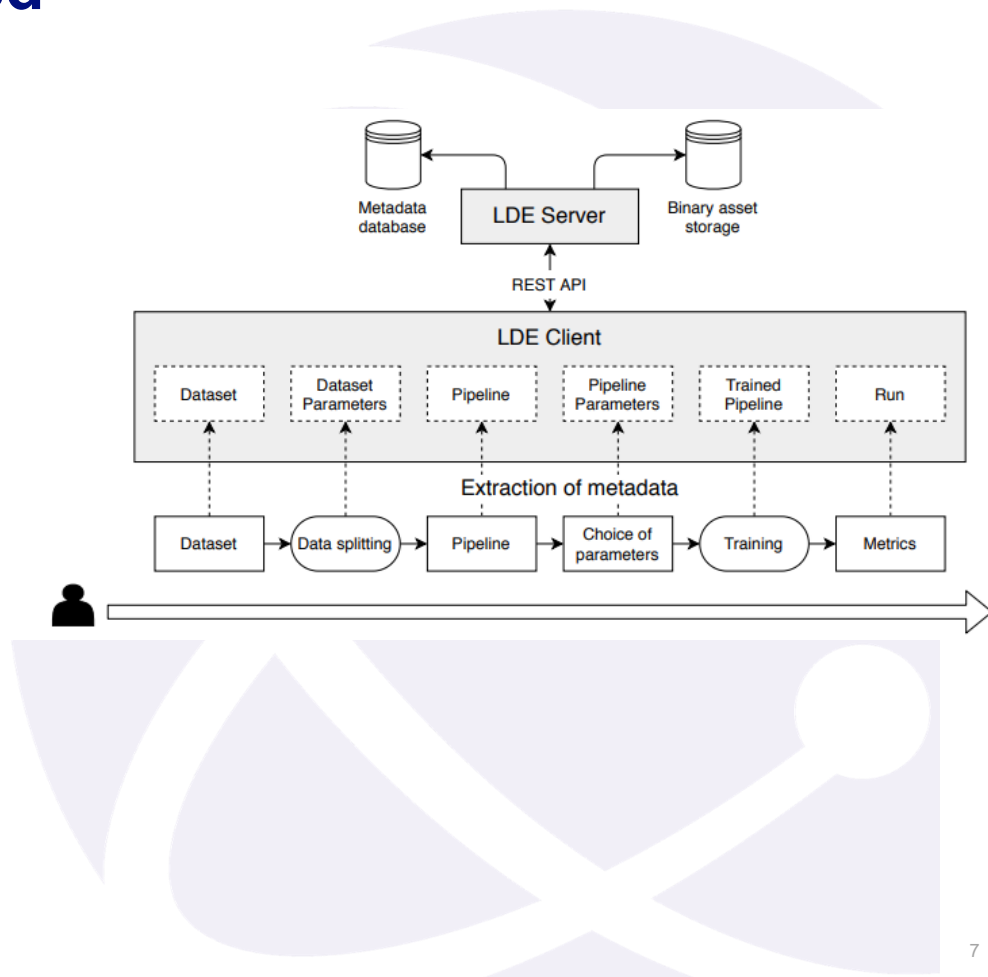- Extensible Database Simulator for Fast Prototyping In-Database DBSIM.
1. Implemented a system for in DB analytics of metadata.
2. GUI for SQL queries in RDBMS.

A. Difficult for LANL scientist unable to use SQL.
B. Unable to visualize analytics of metadata.

# Related Works Continued

- Enabling Reproducibility and Meta-learning Through a Lifelong Database of Experiments.
1. Created a permanent DB that would extract and interact with artifacts from AI experiments from datasets to pipelines.
2. Metadata extracted from experiments contained most of information to reproduce the given experiment.
A. Data captured didn't include environment configuration and modules used.
B. No user permission restrictions.
C. No visualizations plots for statistical analysis of metadata. Unable to perform in database analytics of metadata.
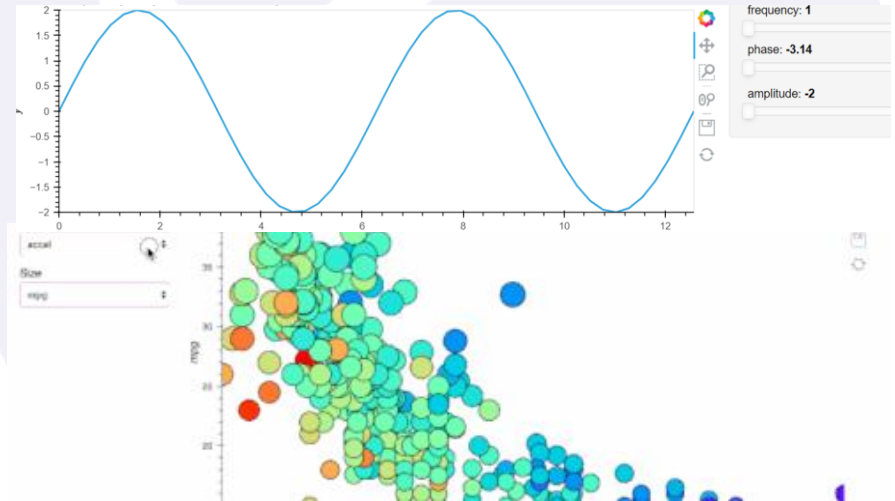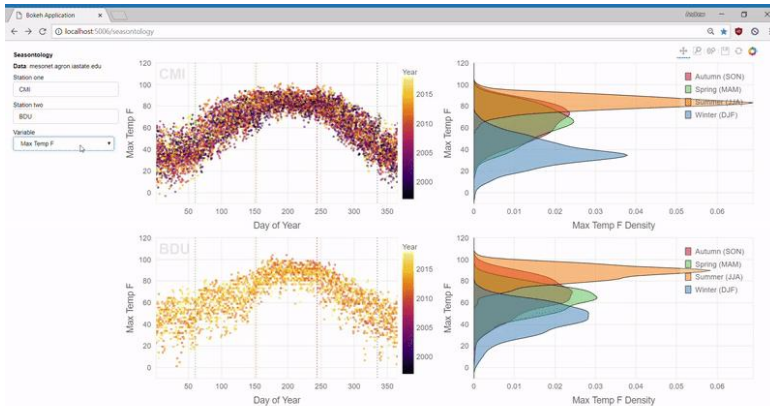
# Expectations for DSI Project

- Automated Workflows and Improved metadata management for LANL mission relevant science workflows.
- Data and Experiment Reproducibility by in-database analytics and learning.
- Create a persistent database for sharing and collaboration.
- Storage and logging of parameters, code versions, metrics, environment configurations with specific modules, and output files.

# Potential Front-Ends for DSI

- R&D on different GUI options for DSI such as ModelDB, Apache Superset, Trame, and ML-Flow.
- Data and Experiment Reproducibility by in-database analytics and learning.
- Create a persistent database for sharing and collaboration.

# Features Implemented

- In Database analytics for comparing datasets and interactively analyze dataset statistics.
- Visualization plots for comparing metadata and sharing results.
- SQL queries implemented in backend to interact with GUI for Querying and data access API.
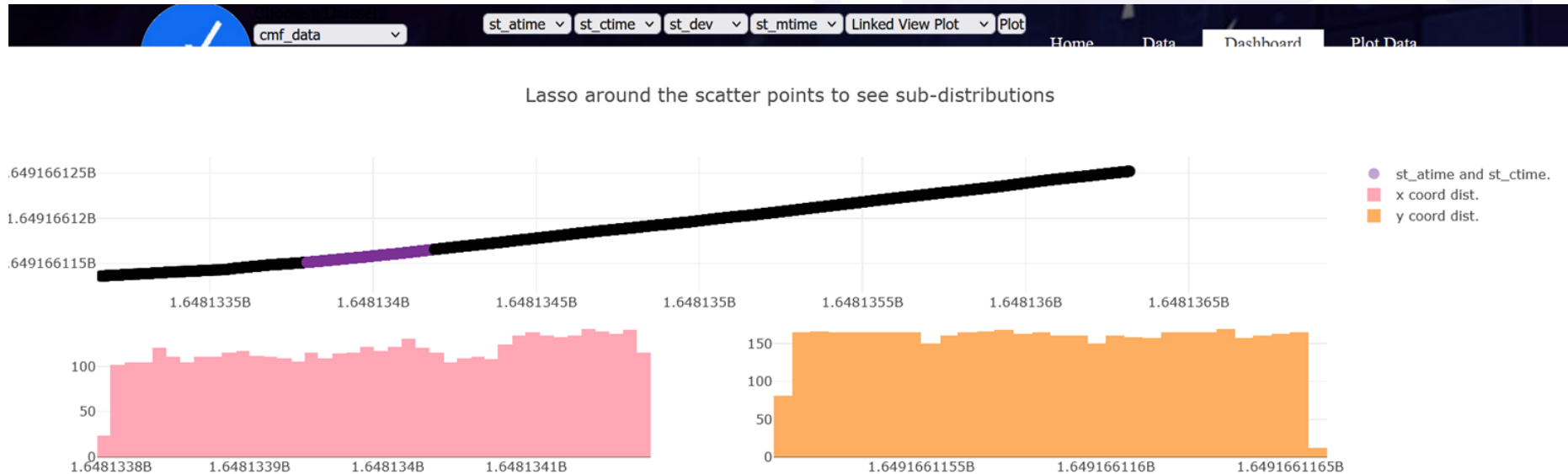- Access to metadata, data via Backend APIs.

Choose a Dataset: cmf_data | st_atime | MAX | Compute

*1648136323.44526*

# More Features Implemented

- Linked view of charts with point selection using the 'lasso' tool. Users can select points of interest on the scatter plot for an in-depth analysis and analysis of CMF data.

# And More Features Implemented…

- Parallel Coordinates Plot Viewer. Selected entries from the scatter plot are highlighted in a parallel coordinates plot viewer.