# Schema–Independent Scientific Data Cataloging Framework

**Supun Nakandala, Sachith Dhanushka Withana, Dinu Kumarasiri, Hirantha Jayawardena and H.M.N. Dilum Bandara**
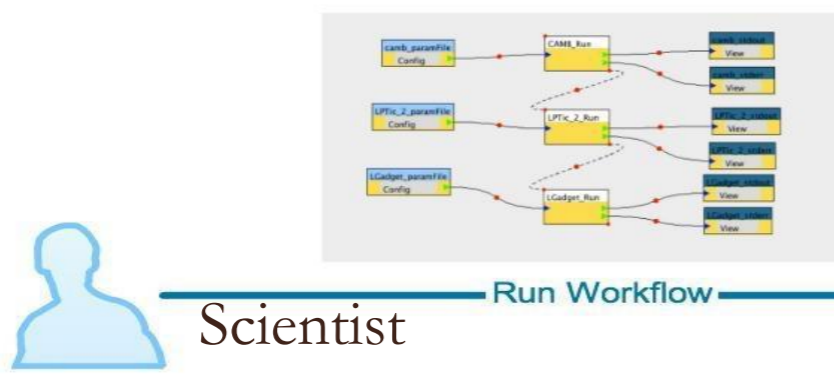
(Department of Computer Science and Engineering, University of Moratuwa, Sri Lanka)

**Srinath Perera**

(Wso2 Inc., Colombo, Sri Lanka)

**Suresh Marru, Sudhakar Pamidighantam**

Indiana University, Bloomington, USA

Scientist

Run Workflow

Cloud

Results + Intermediate Data

Problem

Data Dump

Satellite Data

Instrument Data

Scientific Data
- ✧ Vast Volume
- ✧ Hard to …
  - ❑ Search
  - ❑ Reuse
  - ❑ Share findings

2

# GridChem Usecase

- Gaussian 9 experiments generate vast amount of data in two forms
  - Output file (*.out)
  - Check point file (*.chk)
- Provide efficient searching among these data

# Why we need a new one ??
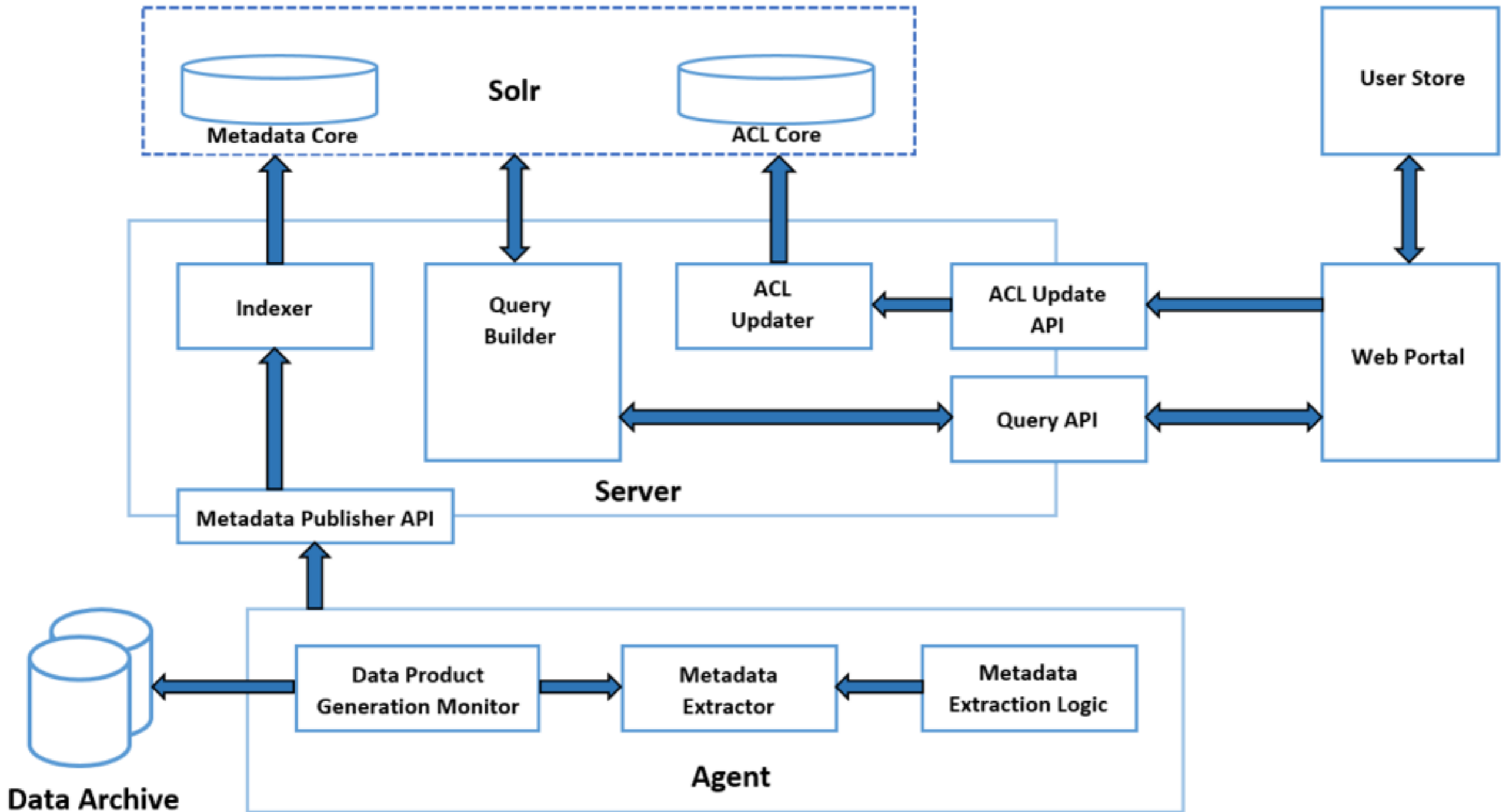
**Existing Solutions**

- Tightly coupled

- Inflexible querying

- Static schemas
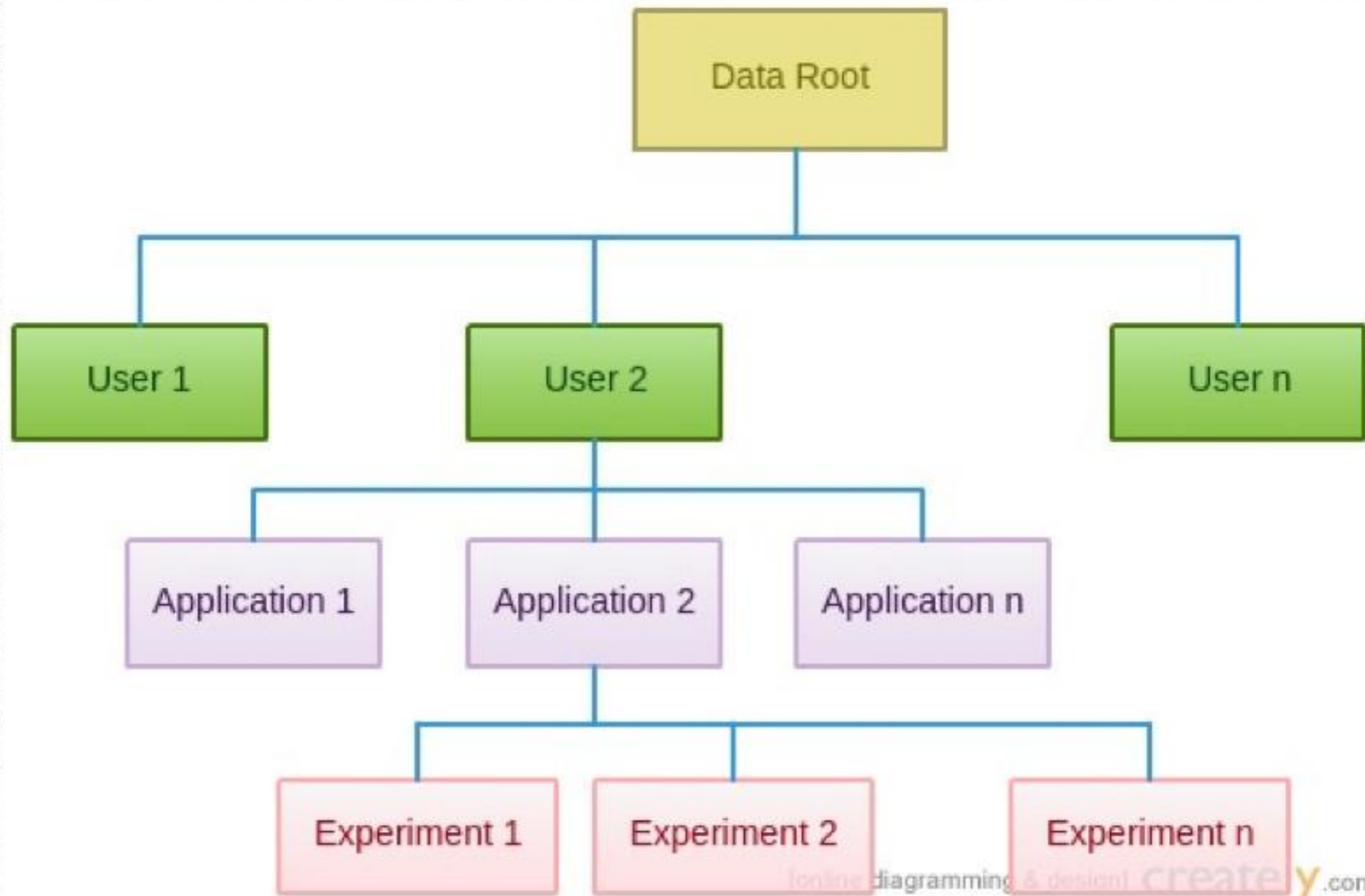
- Eg:-

  - MCS

  - MCAT

  - MyLEAD

**Our Solution**

- Generalizable framework

- Flexible querying

  o Wild card queries

  o Full text queries

  o Substring queries

  o Fielded queries

- Static schema + dynamic fields

# High-level Architecture

# Folder Structure

# What is new in our solution?

- Pluggable metadata extraction logic
- Extensible data product generation monitors
- Use of NoSQL database (Apache Solr)
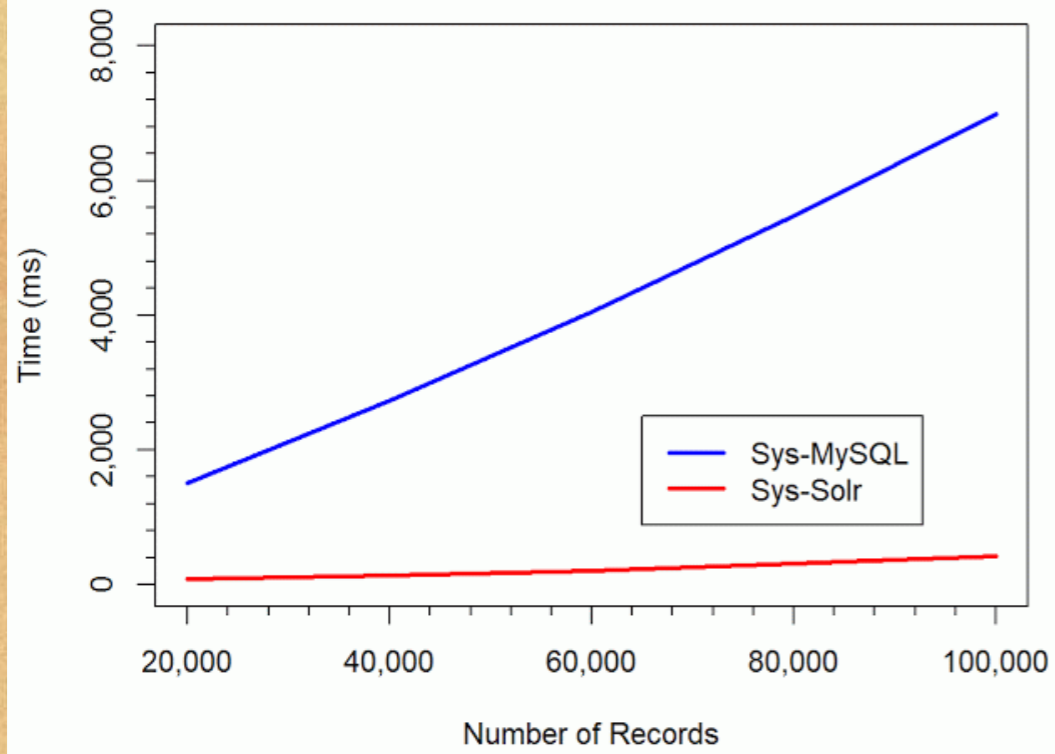- Ability to dynamically add metadata fields

# Performance Test

- MySQL vs Solr
- Data Insert Performance
- Query Performance
  - Exact match queries
  - Range queries
  - Full text queries
  - Prefix match queries
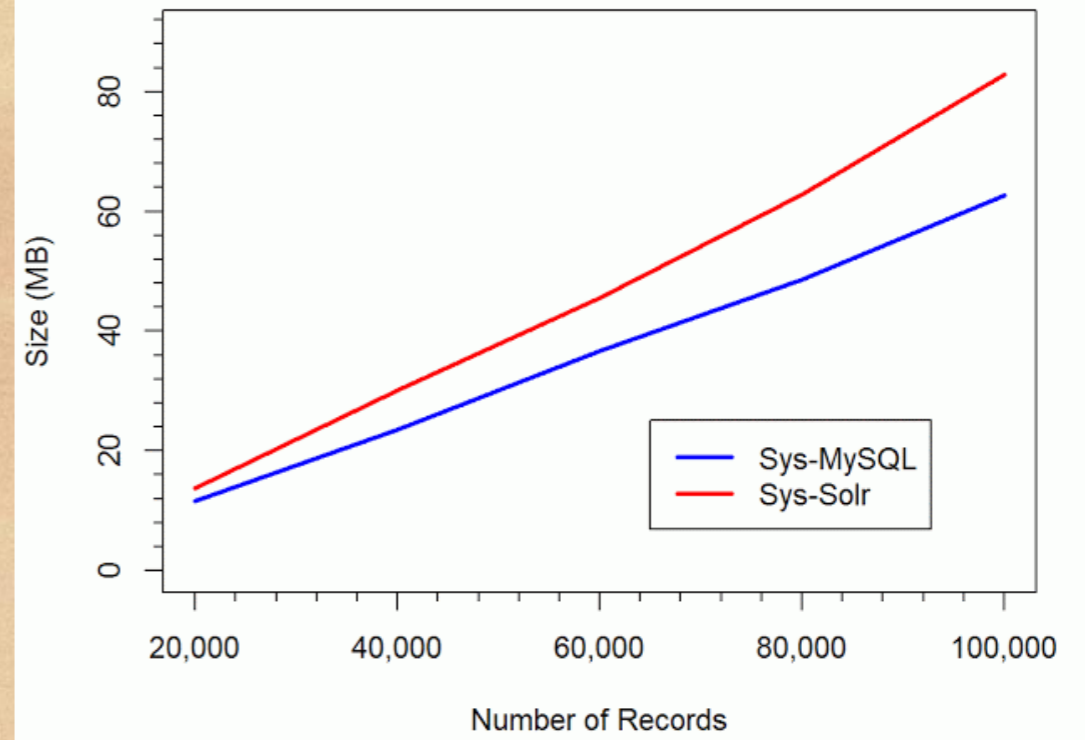  - Suffix match queries
  - Wildcard queries
  - Substring queries

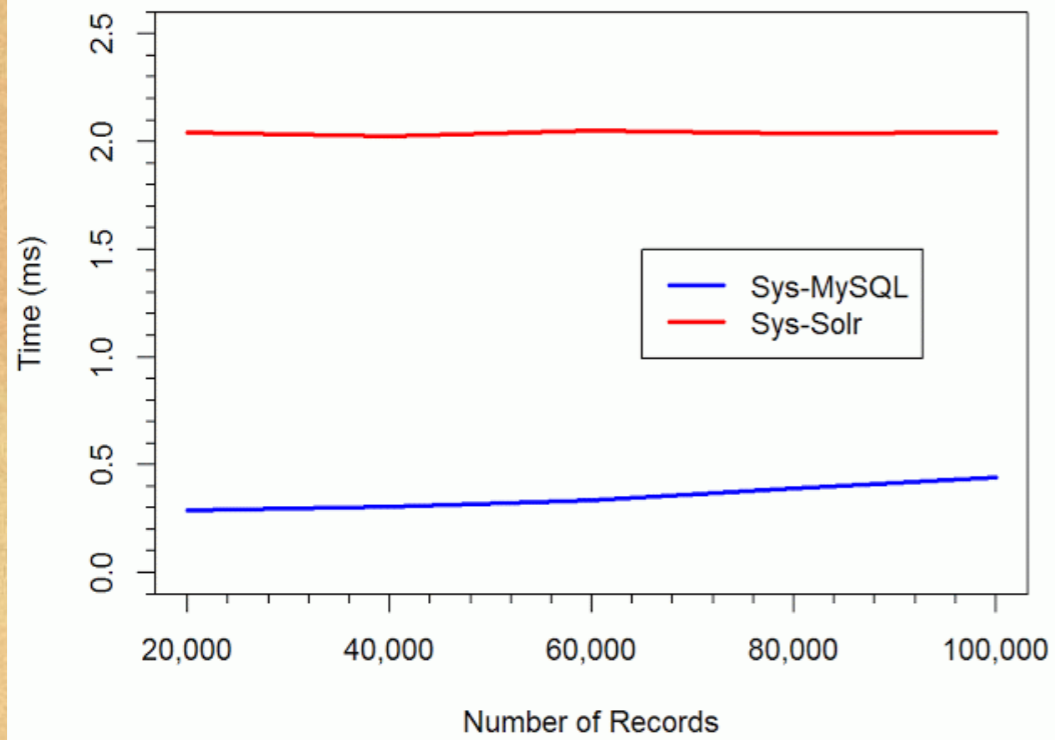Solr resolves more complex queries 91% - 99% faster than a MySQL-based implementation.
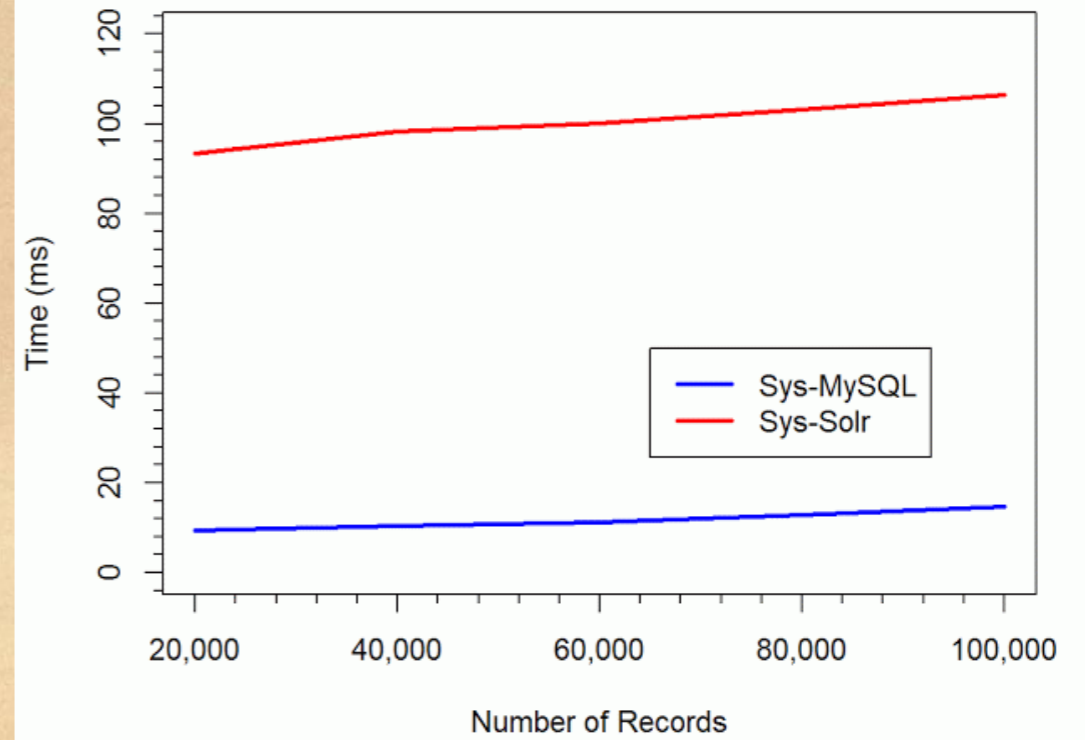
**Data Insertion Time**

Time (ms) vs Number of Records
- Sys-MySQL
- Sys-Solr

**Storage Space Utilization**

Size (MB) vs Number of Records
- Sys-MySQL
- Sys-Solr

**Exact Match Query Execution**

Time (ms) vs Number of Records
- Sys-MySQL
- Sys-Solr

**Range Query Execution**

Time (ms) vs Number of Records
- Sys-MySQL
- Sys-Solr

# Summary

- What we did: A schema-independent scientific data catalog with pluggable parser logic and Solr backend

- Future work: Airavata integration and provenance aware execution

# Thank You …