*Exceptional service in the national interest*

Sandia National Laboratories

# EMPRESS—Extensible Metadata PRovider for Extreme-scale Scientific Simulations

**Margaret Lawson**, Jay Lofstead, Scott Levy, Patrick Widener,

Craig Ulmer, Shyamali Mukherjee, Gary Templet, Todd Kordenbrock

SAND2017-12103 C

# Problems Faced

- Simulations with 100s TB per output, run every few minutes
  - Ex. XGC1, Square Kilometer Array Radio Telescope (SKA)

- Storage devices too slow to sift through all output to find "interesting data"

- Scientists have specific data they want to retrieve
  - Ex. "blob" in fusion reactor or a phenomenon in astronomy

# Motivating Question

*How can we facilitate scientific discovery from simulations in the exascale age?*

# EMPRESS' Solution

- Allow users to label data and retrieve data based on labels

- Features:
    - Robust, standard per-process metadata
    - User-created metadata that is fully customizable at runtime
    - Programmatic query API to retrieve data contents based on metadata

# Previous Solutions

- HDF5 and NetCDF – rudimentary attribute capabilities, basic metadata
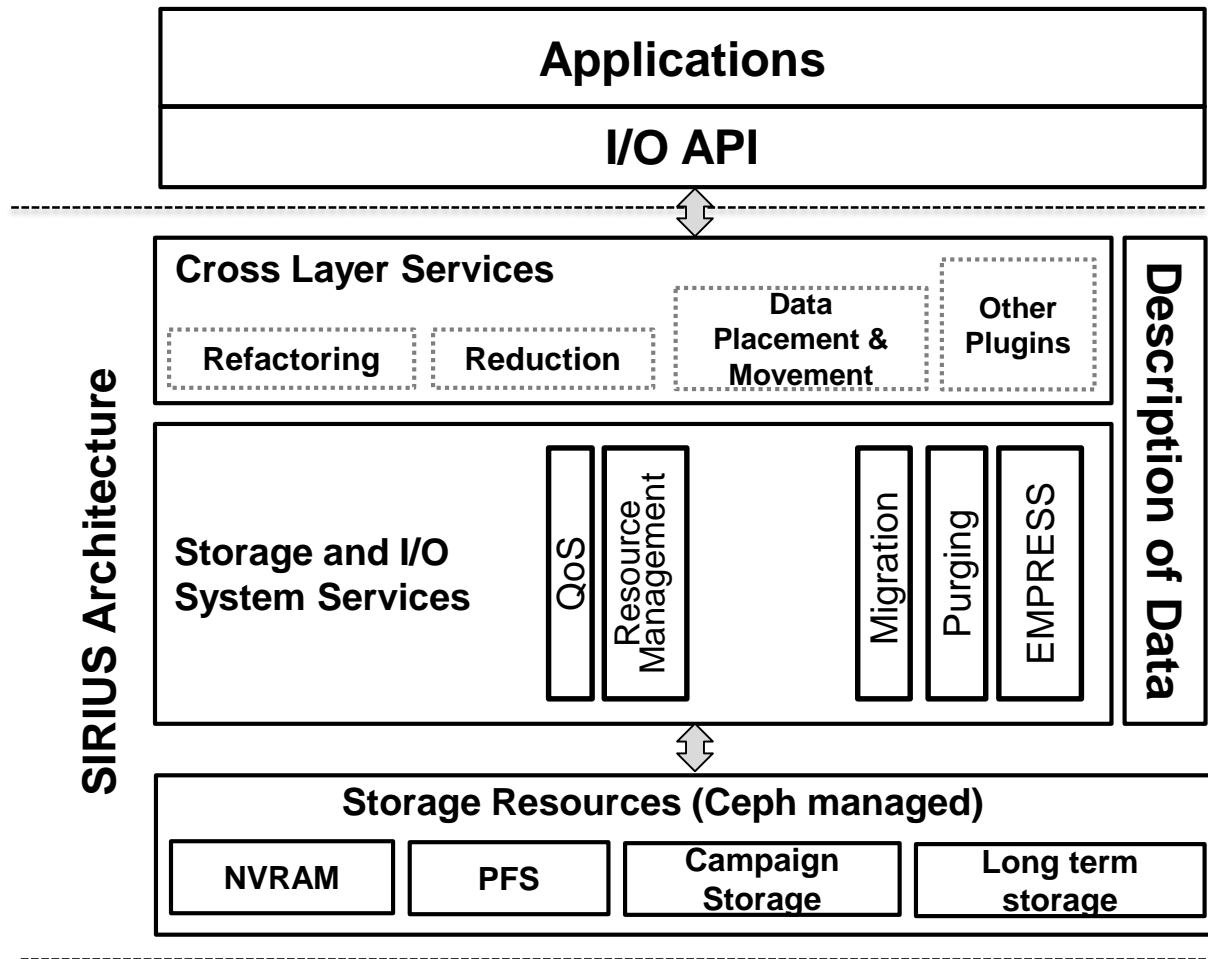
- ADIOS – per-process metadata

None of these address efficient attribute searching

- FastBit – offers data querying based on values, but very limited support for spatial queries and attributes
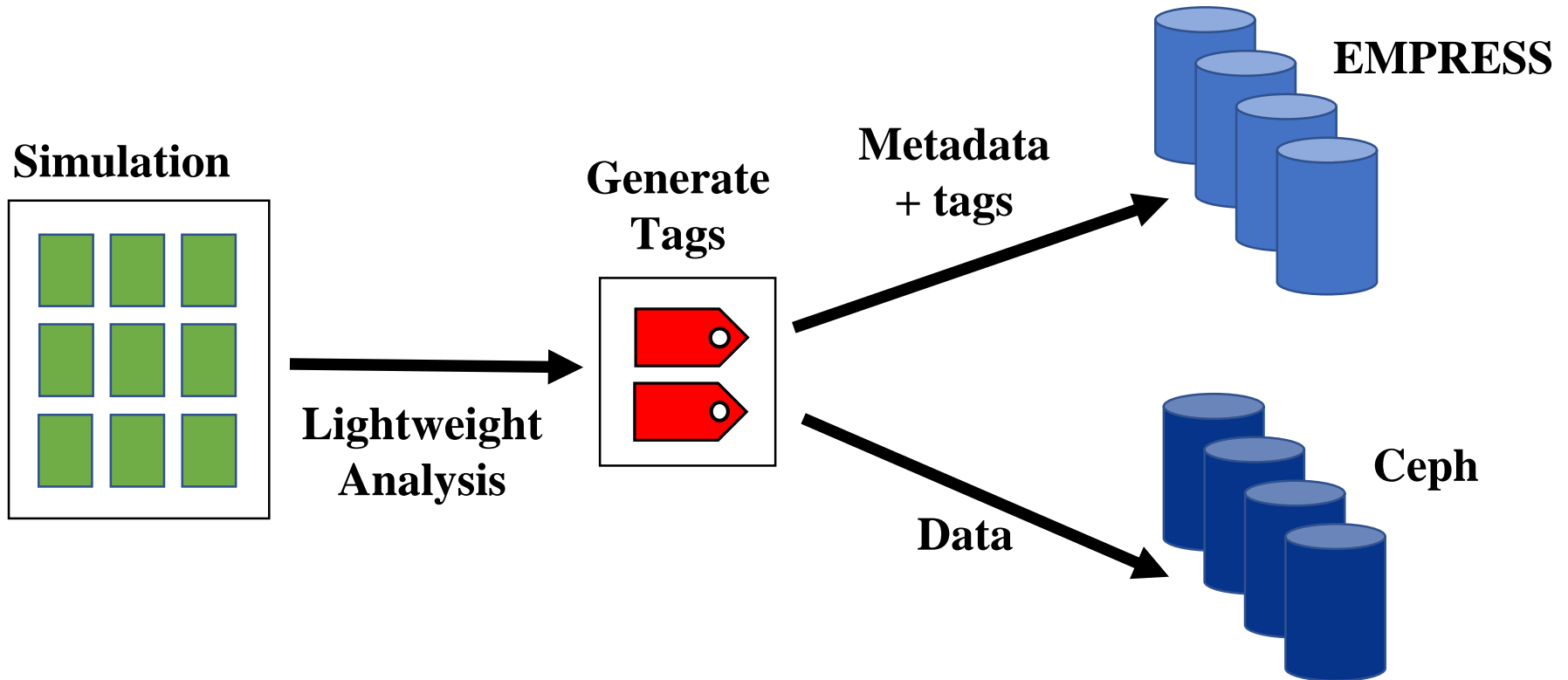
# Why not use a Key-Value Store?

- Custom keys can go a long way, but not far enough
- Two Problems:
    - Inexact matches
    - Custom Metadata

- Relational databases with indices are radically faster at searching like this
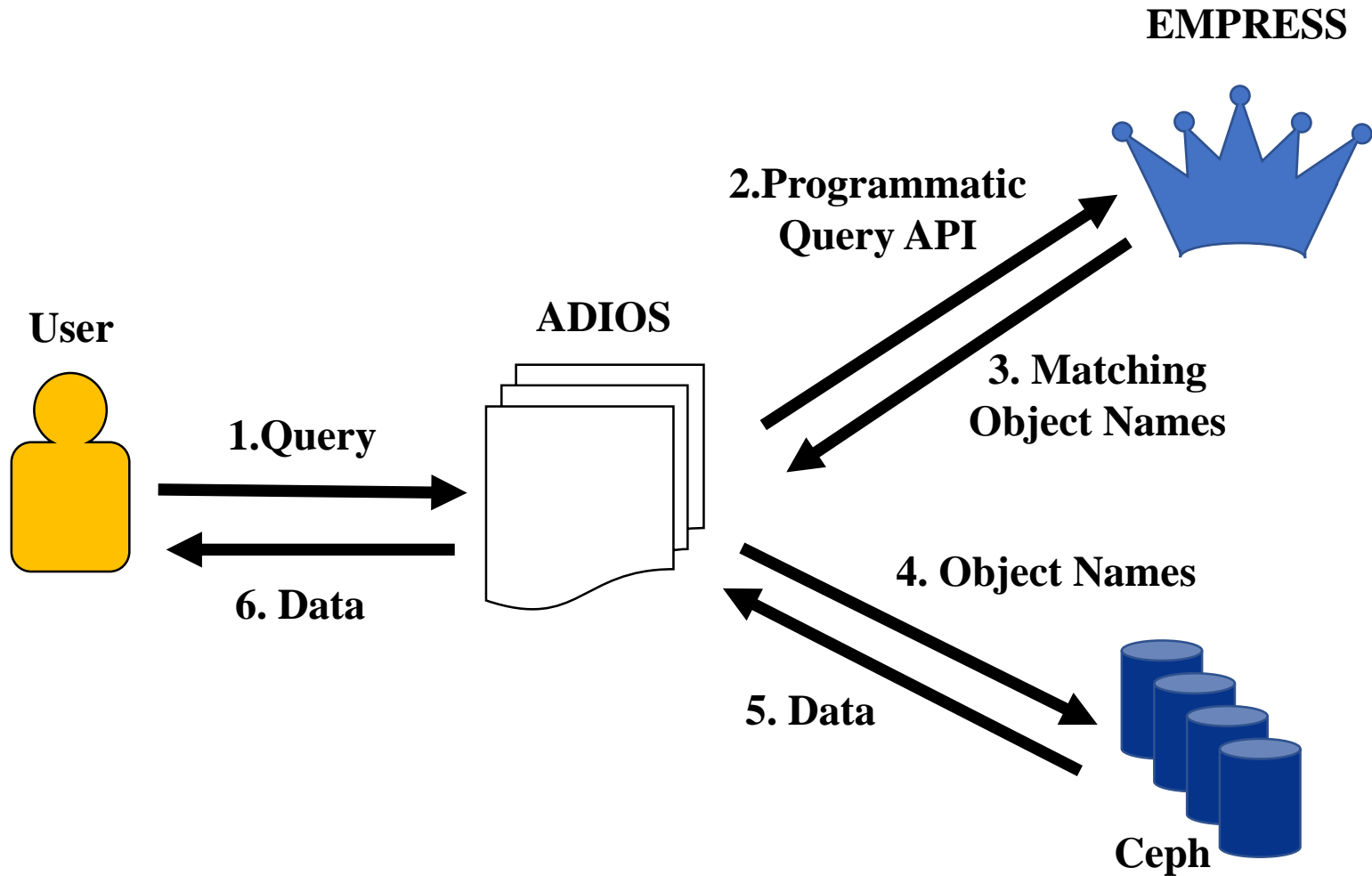
# SIRIUS Architecture
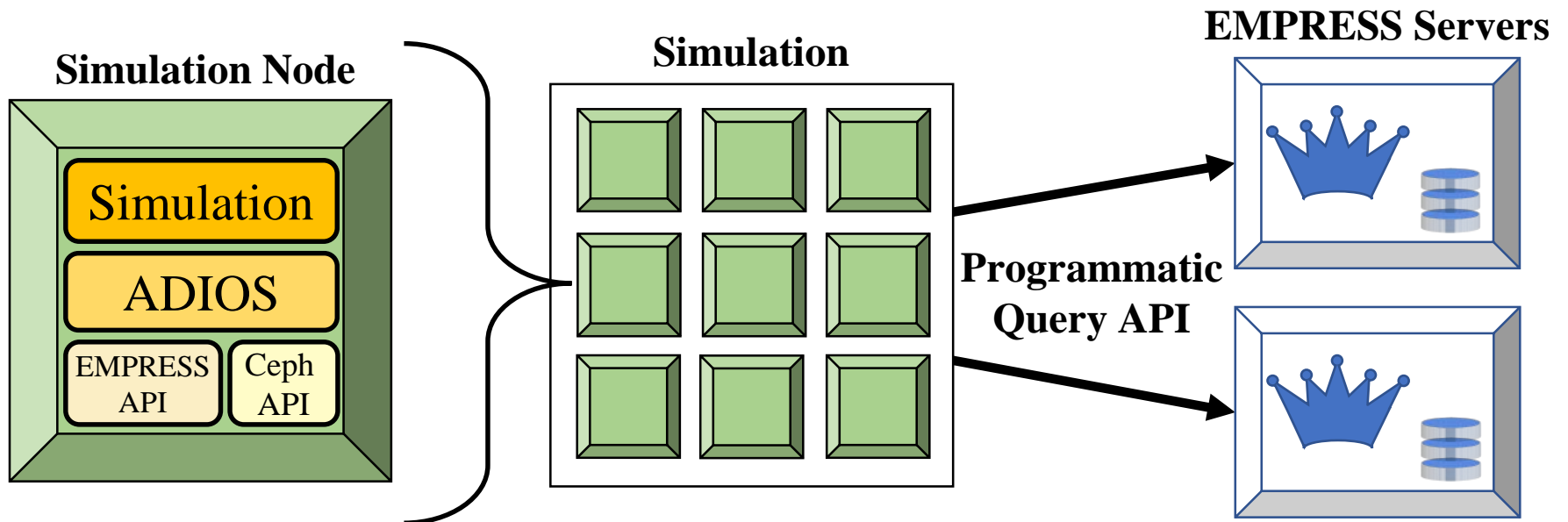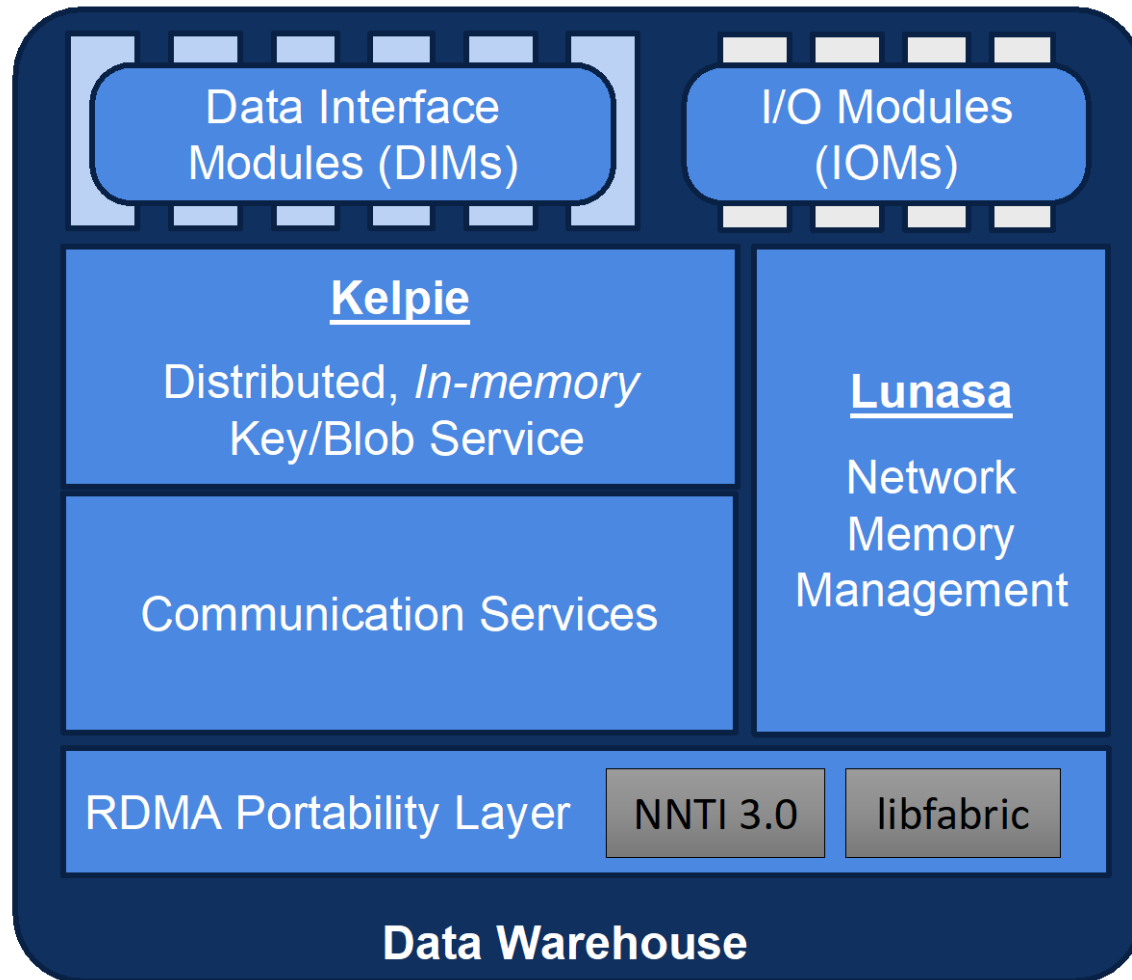
# SIRIUS Workflow – Write Process

**Simulation**

**Generate Tags**

**Metadata + tags**

**EMPRESS**

**Lightweight Analysis**

**Data**

**Ceph**

# SIRIUS Workflow – Read Process

**EMPRESS**

**User**

**ADIOS**

**2.Programmatic Query API**

**3. Matching Object Names**

**1.Query**

**6. Data**

**4. Object Names**

**5. Data**

**Ceph**

# High Level Design

# Faodail

# Storage - Tracked Metadata

- Dataset information
  - Application, run, and timestep information
- Variable information
  - Catalogs types of data stored for an output operation
- Variable chunk information
  - Subdivision of simulation space associated with a particular variable
- Custom metadata class
  - Metadata category the user adds for a particular dataset
  - Ex. Max
- Custom metadata instance
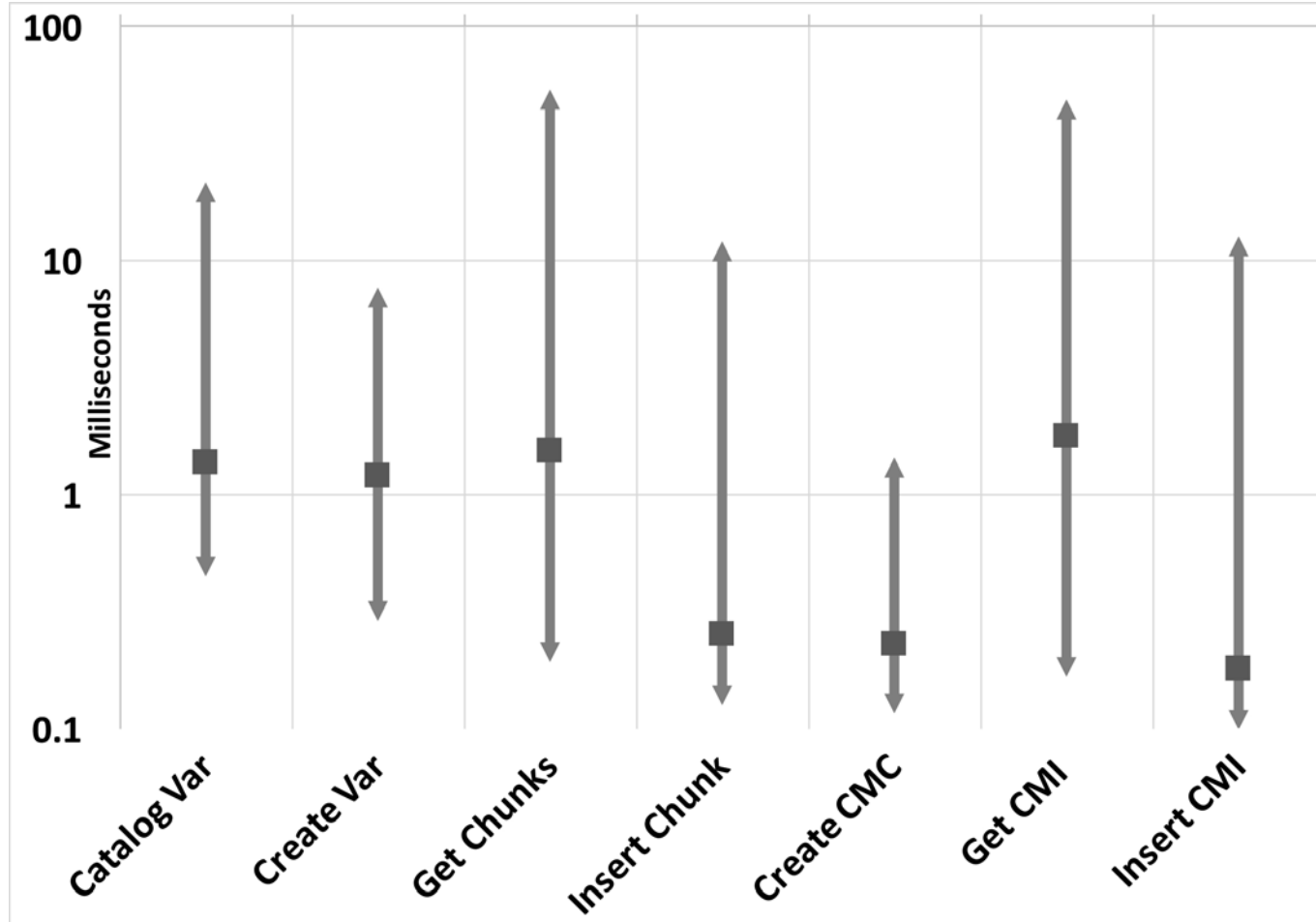  - Ex. Flag for chunk or a bounding box spanning chunks

# Testing Goals

- Scalable?
  - Number of client processes: 1024-2048
- Effect of client to server ratio
  - Ratios tested: 32:1 – 128:1
- Overhead of including a large number of custom metadata items
  - Number of custom metadata classes: 0 or 10
  - On average 2.641 custom metadata instances per chunk
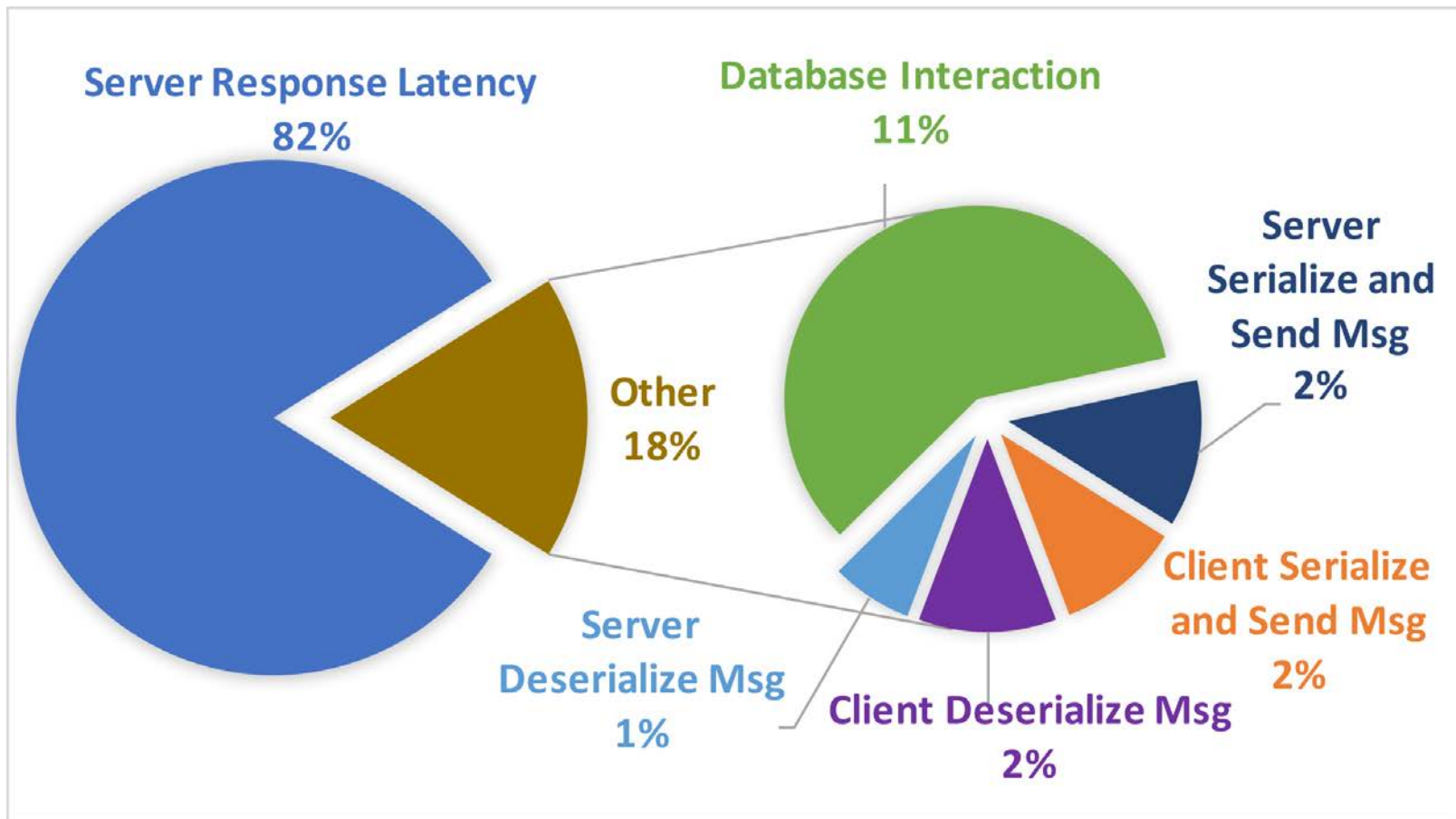
# Testing Goals (Continued)

- Proof of concept, can EMPRESS efficiently support:
  - Common writing operations
    - 2 datasets written, each with 10 globally distributed 3-D arrays
  - Common reading operations
    - 6 different read patterns that scientists frequently use (Lofstead, et al. "Six Degrees of Scientific Data")
  - A broad range of custom metadata
    - 10 custom metadata classes including max, flag, bounding box (two 3-D points)
- Scientific validity
  - A minimum of 5 runs per configuration on 3 computing clusters:
    - Serrano (total nodes: 1122)
    - Skybridge (total nodes: 1848)
    - Chama (total nodes: 1232)

# Testing – Query Times



- EMPRESS efficiently supports a wide variety of operations including custom metadata operations
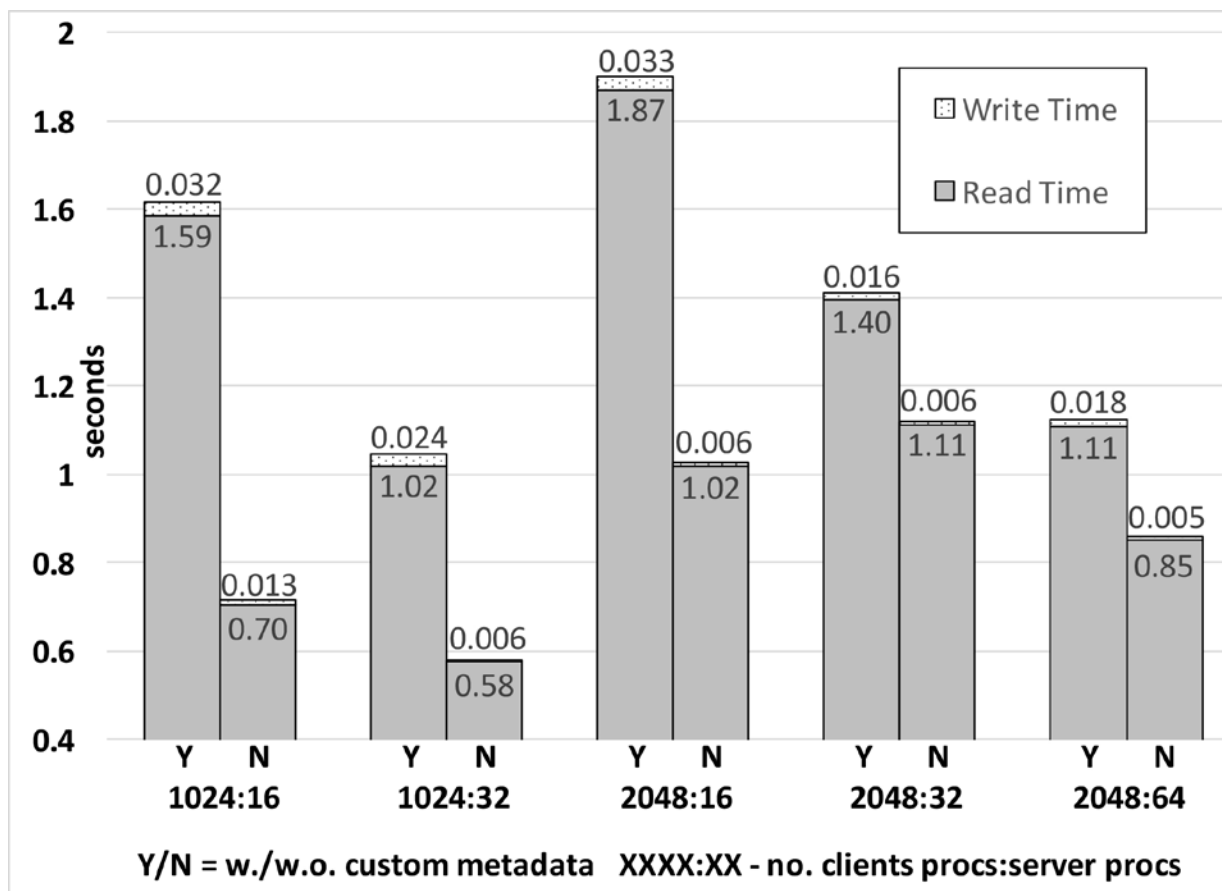
# Testing – Chunk Retrieval Time



- Most time is spent waiting for the server to respond
    - Room for improvement in the Faodail infrastructure

# Testing – Writing and Reading Time



Y/N = w./w.o. custom metadata    XXXX:XX - no. clients procs:server procs

- Good scalability for fixed client-server ratio
- No significant overhead for adding custom metadata
- Client-server ratio greatly affects performance

# Future Work

- Increasing EMPRESS' flexibility, efficiency, and scalability
  - Support more queries
  - Different metadata distribution?

# Acknowledgements

**Algorithm 1** Writing algorithm

1: **procedure** WRITETIMESTEP ▷ Each process does this
2:     **for all** variables assigned **do**
3:         md_create_var (...)
4:     **end for** ▷ Write portion of all vars
5:     **for all** custom metadata classes assigned **do**
6:         md_create_type (...)
7:     **end for** ▷ Write portion of all custom md types
8:     **for all** variables **do**
9:         md_insert_chunk (...) ▷ Add a var chunk; get the ID
10:         **for all** custom metadata desired **do**
11:             md_insert_attribute (...) ▷ Add custom md instance
12:         **end for**
13:     **end for**
14: **end procedure**

**Algorithm 2** Reading algorithm

1: **procedure** READDATA                                   ▷ Each Process Does this
2:     md_catalog_vars (...)        ▷ Get list of vars from any server
3:     **for all** metadata servers needed **do**
4:         md_get_chunk(...)   ▷ get all chunks in area of interest
5:         **for all** chunks returned **do**
6:             md_get_attribute (...) ▷ get the custom md instances
7:         **end for**
8:     **end for**
9: **end procedure**