# Metadata Management to Aid Data Discovery

Jay Lofstead

NSDF All Hands

12 April 2023

Exceptional

service

in the

national

interest

Sandia National Laboratories

# Problem Space

- Long term data archiving
  - Which data set(s) contain what I want to study?

- Short/medium term data set identification
  - What does each set contain and which ones to save?

- Provenance largely focuses on environment
  - Critical, but insufficient/inefficient

*Three current generations of tools with a fourth being developed*

# First Generation Tools

- File-level metadata tagging

- POSIX extended attributes

- HPSS

- Starfish (https://starfishstorage.com/)

- JAMO – Joint Genome Institute Archive and Metadata Organizer

# Second Generation Tools

- Raw data indexing
  - Exact value or binning

- FastBit
  - Does a value in the bin range exist in this file/dataset?

- SciDB
  - Multi-dimensional array data model with rich query facilities

- IO libraries
  - HDF5, ADIOS, NetCDF attribute capabilities

# Third Generation Tools

- Feature tagging
    - Region/var/run
    - Bounding box, simple tag

- Key-value based
    - SoMeta - encode tag information into key and use value for data location (focused on object stores)
    - TagIt – distributed, shared nothing storage integrated for faster data searching

# Third Generation Tools

- RDMBS-based
  - EMPRESS – independent database(s) of tags related to run, timestep, or var; flexible query interface into data directly via logical locations
  - BIMM – image database with tags

# Fourth Generation Tool

- Derived quantity information-based tagging
  - For a climate model, where is the pressure gradient greater than a particular value

- Problems!
  - Derived quantities take time to compute
  - Derived quantities can be as large as the original data

# Thank you

- [gflofst@sandia.gov](mailto:gflofst@sandia.gov)