

# PROV-IO<sup>+</sup>: A Provenance Framework for Scientific Data on HPC Systems

Runzhou Han<sup>1</sup>, Mai Zheng<sup>1</sup>, Suren Byna<sup>2</sup>, Houjun Tang<sup>3</sup>, Bin Dong<sup>3</sup>, Dong Dai<sup>4</sup>  
Yong Chen<sup>5</sup>, Dongkyun Kim<sup>5</sup>, Joseph Hassoun<sup>5</sup>, David Thorsley<sup>5</sup>, Matthew  
Wolf<sup>5</sup>

<sup>1</sup>*Iowa State University*   <sup>2</sup>*The Ohio State University*   <sup>3</sup>*Lawrence Berkeley National Laboratory*

<sup>4</sup>*University of North Carolina at Charlotte*   <sup>5</sup>*Samsung*

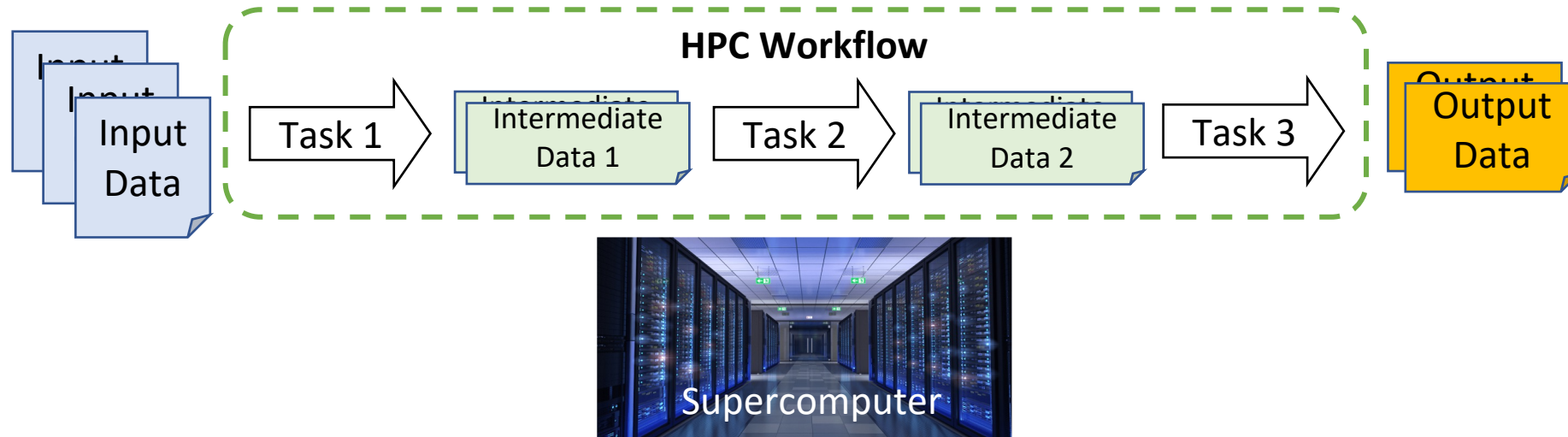


# Motivation

---

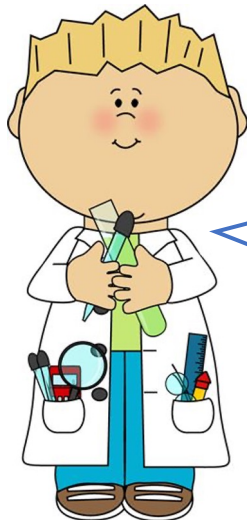
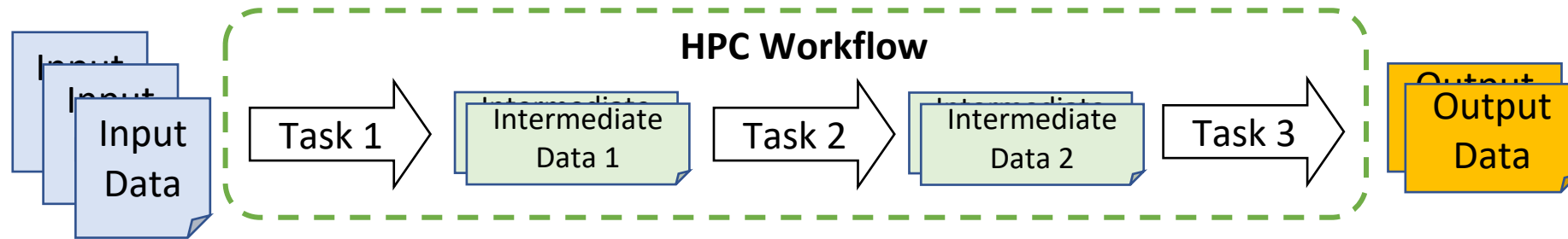
# Scientists Want to Know Their Workflows Better

- Workflows running on HPC systems are complicated



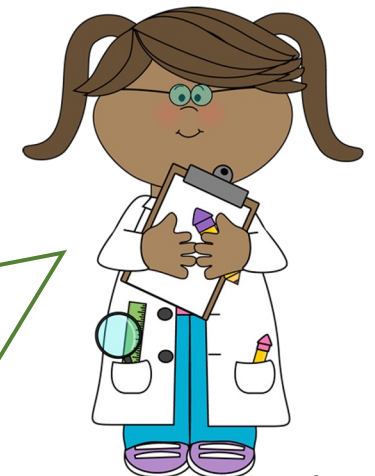
# Scientists Want to Know Their Workflows Better

- Scientists could have a variety of questions about the workflow



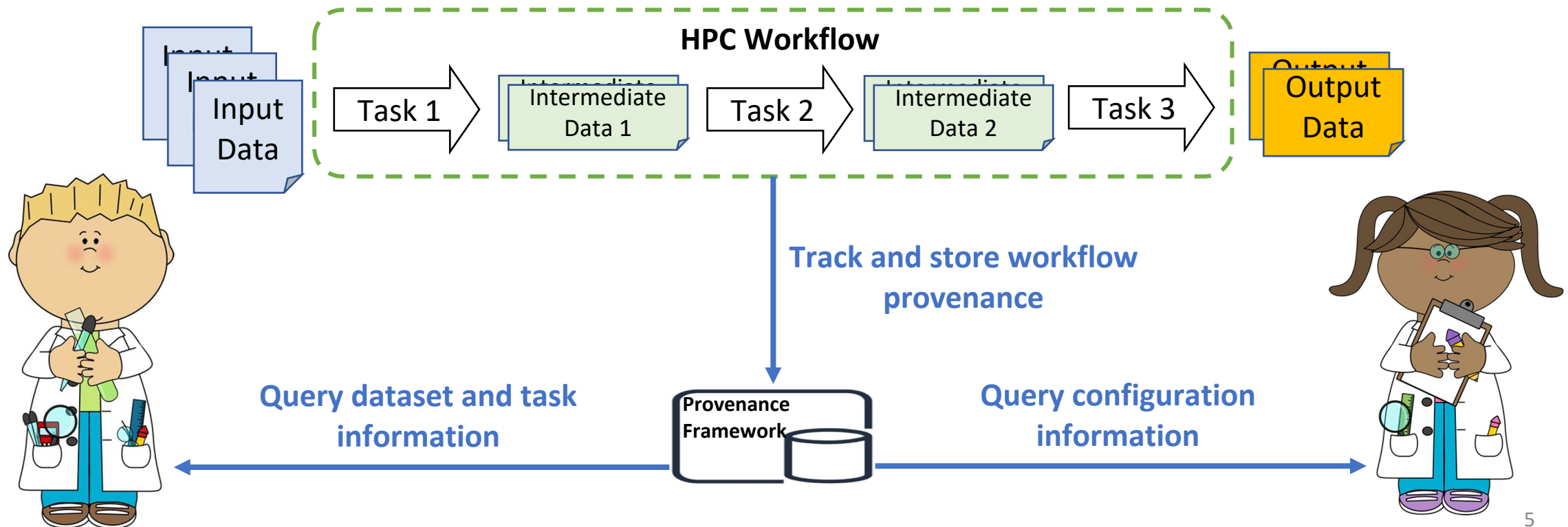
Which dataset slows down the training process?  
Where does the bottleneck take place? ...

Which set of hyperparameters have been used? Which set of data preselection has the best training accuracy? ...



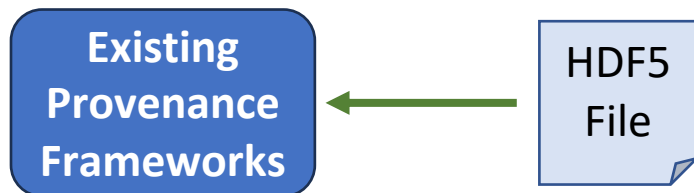
# Provenance Framework Is Designed to Help

- **Provenance Frameworks** are used to collect execution metadata
  - E.g., PASS (ATC'06), PASSv2 (ATC'09), Komadu (JORS'15), ProvLake (eScience'19)



# Limitations of State-of-the-art

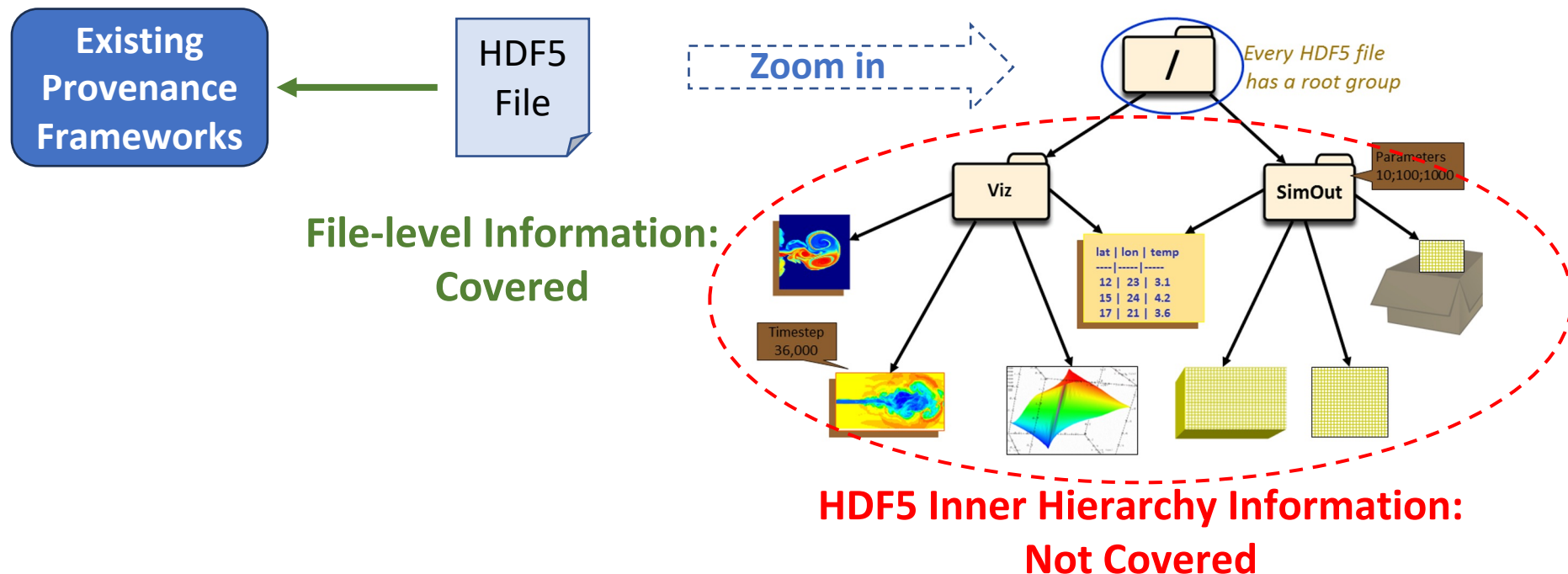
- Limitation 1: Granularity
  - Cannot cover inner hierarchies of scientific data or detailed I/O operations



**File-level Information:  
Covered**

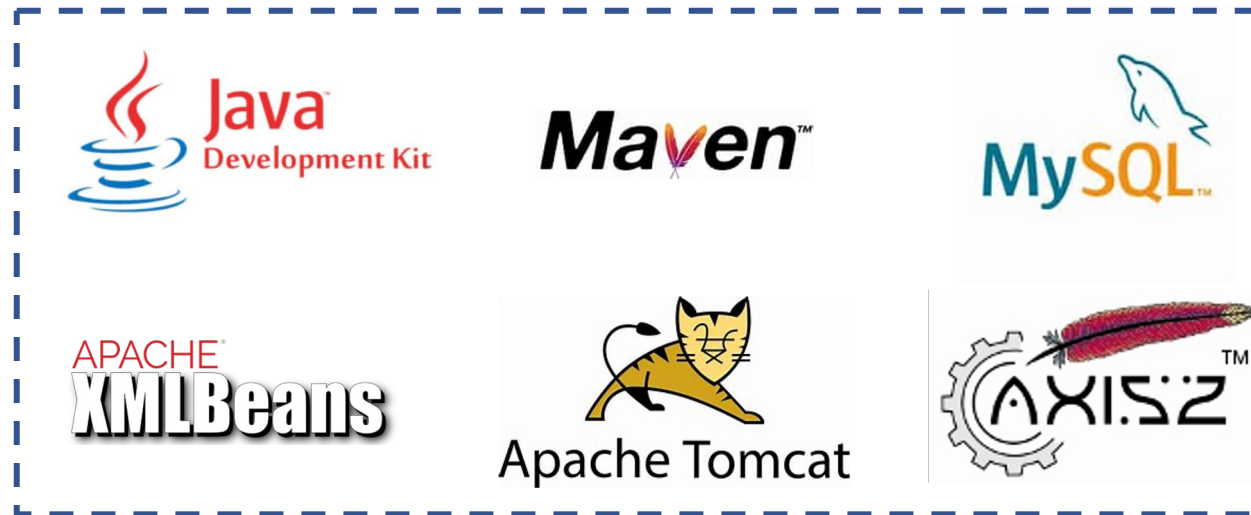
# Limitations of State-of-the-art

- Limitation 1: Granularity
  - Cannot cover inner hierarchies of scientific data or detailed I/O operations



# Limitations of State-of-the-art

- Limitation 2: Compatibility & Portability
  - Heavy dependencies on third party tools which are difficult to port (e.g., “Komadu”)



**Komadu Dependencies**



# Limitations of State-of-the-art

- Limitation 3: Transparency
  - Scientists have to instrument workflows with specific APIs manually (e.g., “ProvLake”)

```
from provlake import ProvLake
from provlake.capture import ProvWorkflow, ProvTask, ProvCycle
```

```
...
prov = ProvLake.get_persister(workflow_name=workflow_name, managed_persistence=False)
workflow = ProvWorkflow(prov).begin()
while current_n > 1:
    with ProvCycle(prov, cycle_name=cycle_name, iteration_id=iteration_id):

        with ProvTask(prov, data_transformation_name=dt_name,
                       parent_cycle_name=cycle_name, parent_cycle_iteration=iteration_id,
                       input_args={"current_n": current_n}) as prov_task:

            result = result * current_n
            current_n = current_n - 1

            prov_task.end({"factorial_result": result})

        iteration_id += 1

print("Finished workflow")
print(result)

workflow.end()
```

Initialize a provenance instance

Insert tracking APIs at multiple hierarchies into the workflow loop

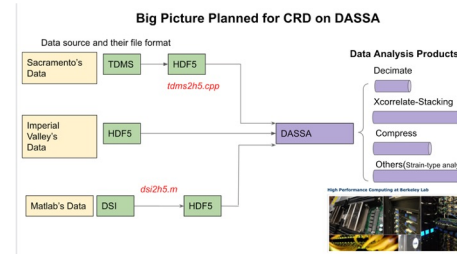
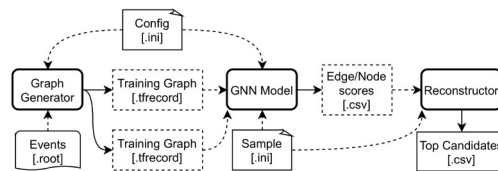
Finalize provenance instance at multiple hierarchies

# Approach

---

# Survey on Practical Needs of Domain Scientists

- Discussed with four research teams from different domains
  - Learn about their workflow & provenance needs



**NVIDIA/Megatron-LM**

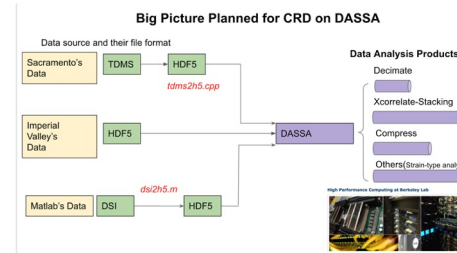
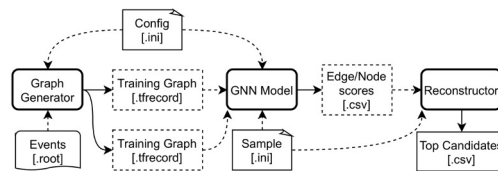


Ongoing research training transformer models at scale

Name	Top Reco	DASSA	H5Bench	Megatron-LM
Domain	GNN for physical emulation	acoustic sensing	synthetic, performance benchmarking	large language model

# Survey on Practical Needs of Domain Scientists

- Discussed with four research teams from different domains
  - Learn about their workflow & provenance needs



NVIDIA/Megatron-LM

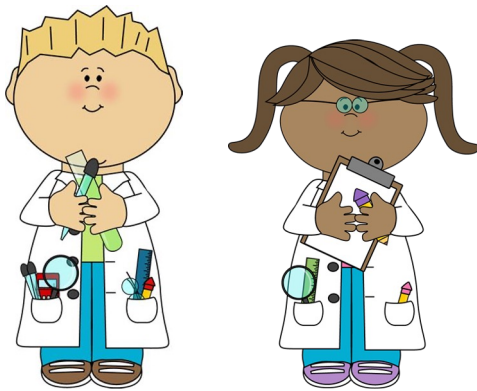


Ongoing research training transformer models at scale

Name	Top Reco	DASSA	H5Bench	Megatron-LM
Domain	GNN for physical emulation	acoustic sensing	synthetic, performance benchmarking	large language model
Provenance need	metadata version control	lineage of data products	I/O stats & performance bottleneck	checkpoint-config. consistency

# Survey on Practical Needs of Domain Scientists

- Summary of provenance needs



We want to know ...

## Data information

End-to-end data information including each intermediate state

## Task information

Information of tasks at multiple granularities, e.g., program, function call

## Configuration information

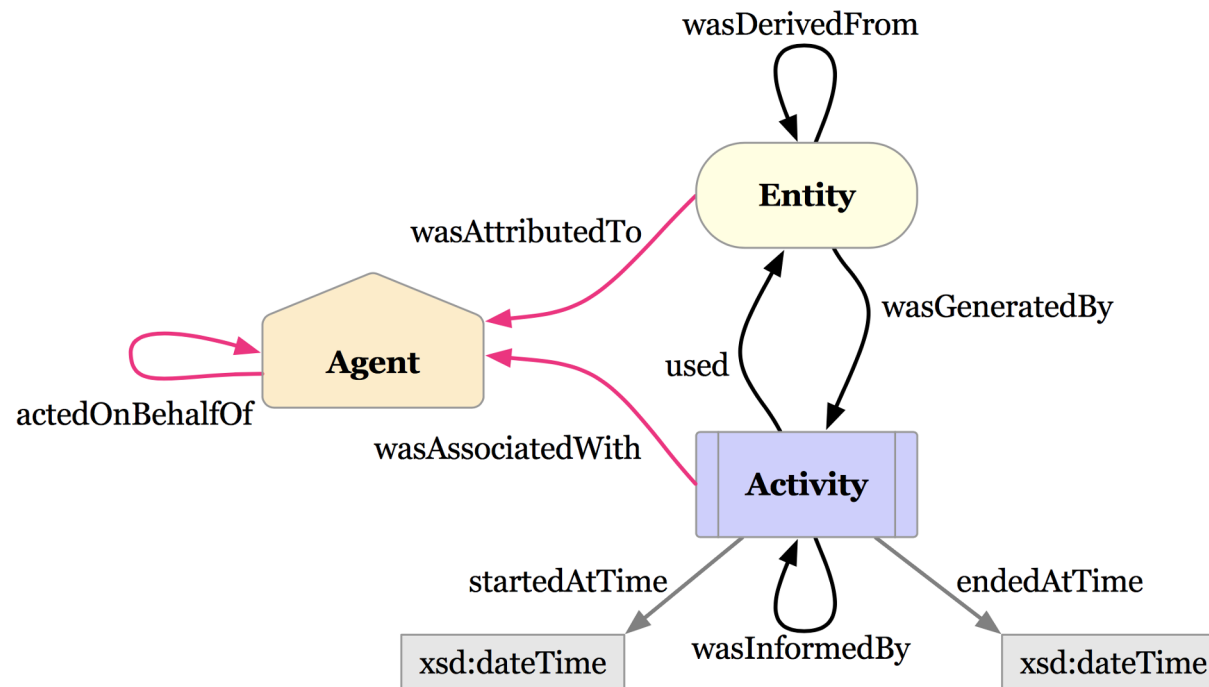
Workflow configurable parameters

## Relation information

Relations between above information

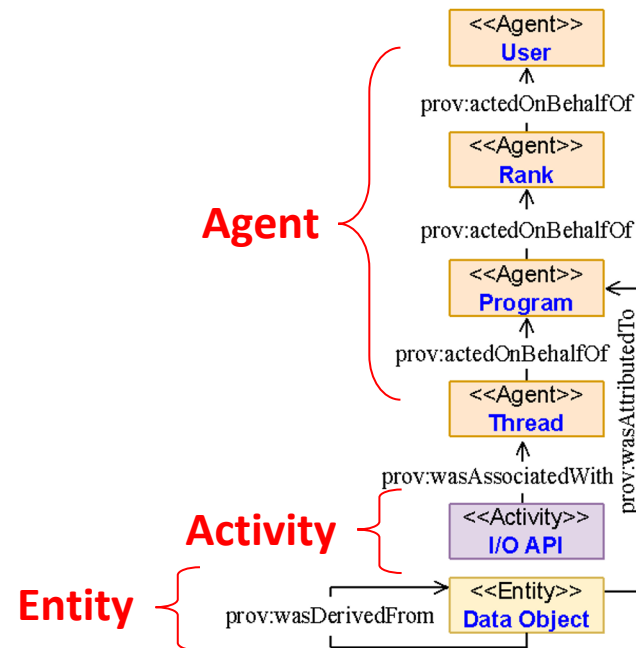
# Design PROV-IO<sup>+</sup> Model Based on Needs

- Derived from W3C provenance data model (PROV-DM)
  - Widely adopted by provenance frameworks (e.g., Komadu, ProvLake)
  - Provides a mapping to Resource Description Framework (RDF) triples



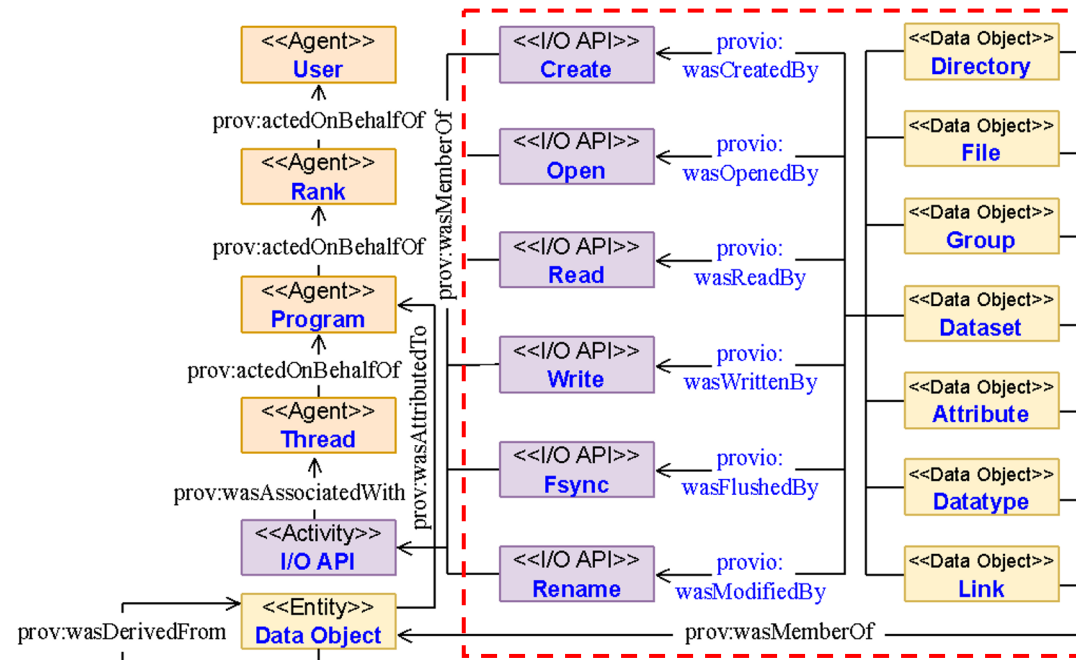
# Design PROV-IO<sup>+</sup> Model Based on Needs

- PROV-IO<sup>+</sup> model: a PROV-DM-compliant data model
  - Interoperable with other PROV-DM-compliant data models (e.g., ProvLake/PROV-ML, Komadu model)



# Design PROV-IO<sup>+</sup> Model Based on Needs

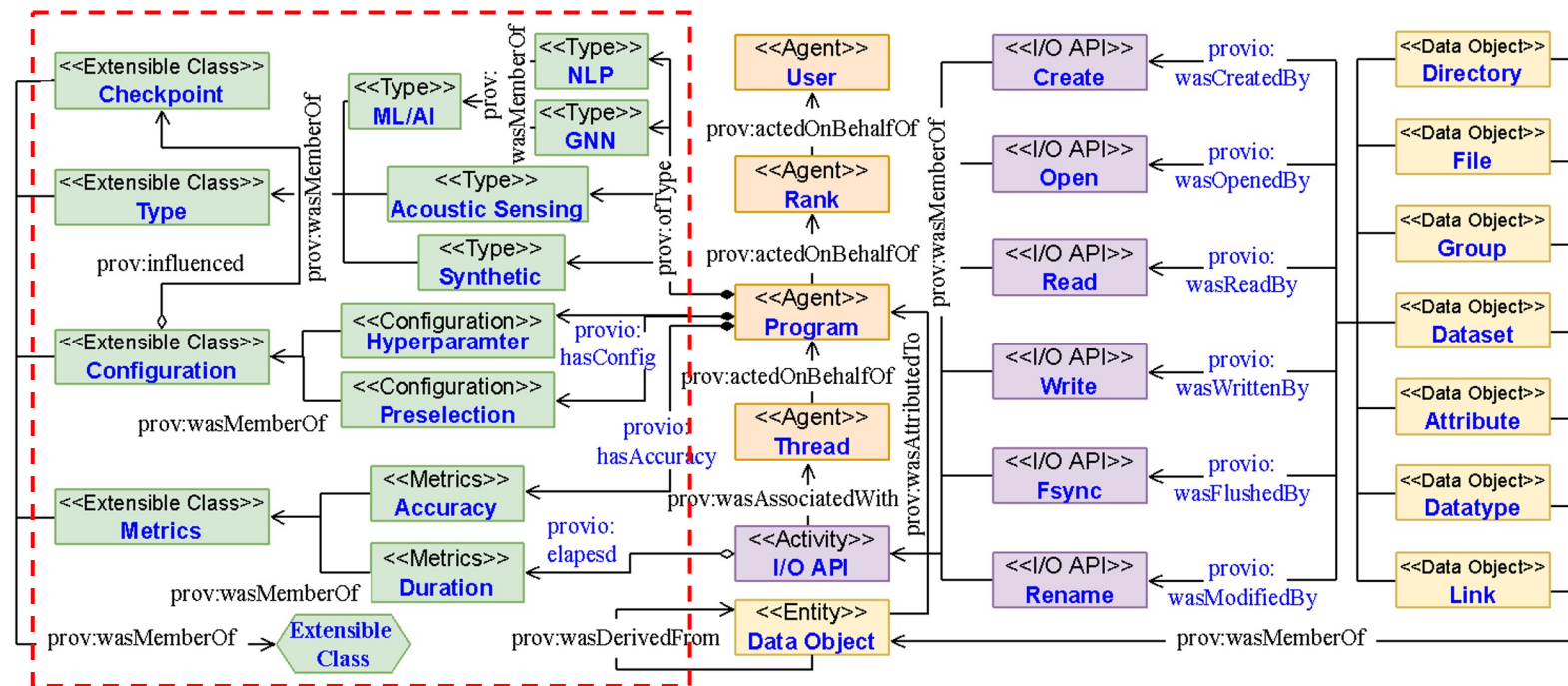
- PROV-IO<sup>+</sup> model: a PROV-DM-compliant data model
  - Covers most of the metadata concepts & I/O operations in popular HPC I/O libraries (e.g., HDF5, POSIX Syscall C library)





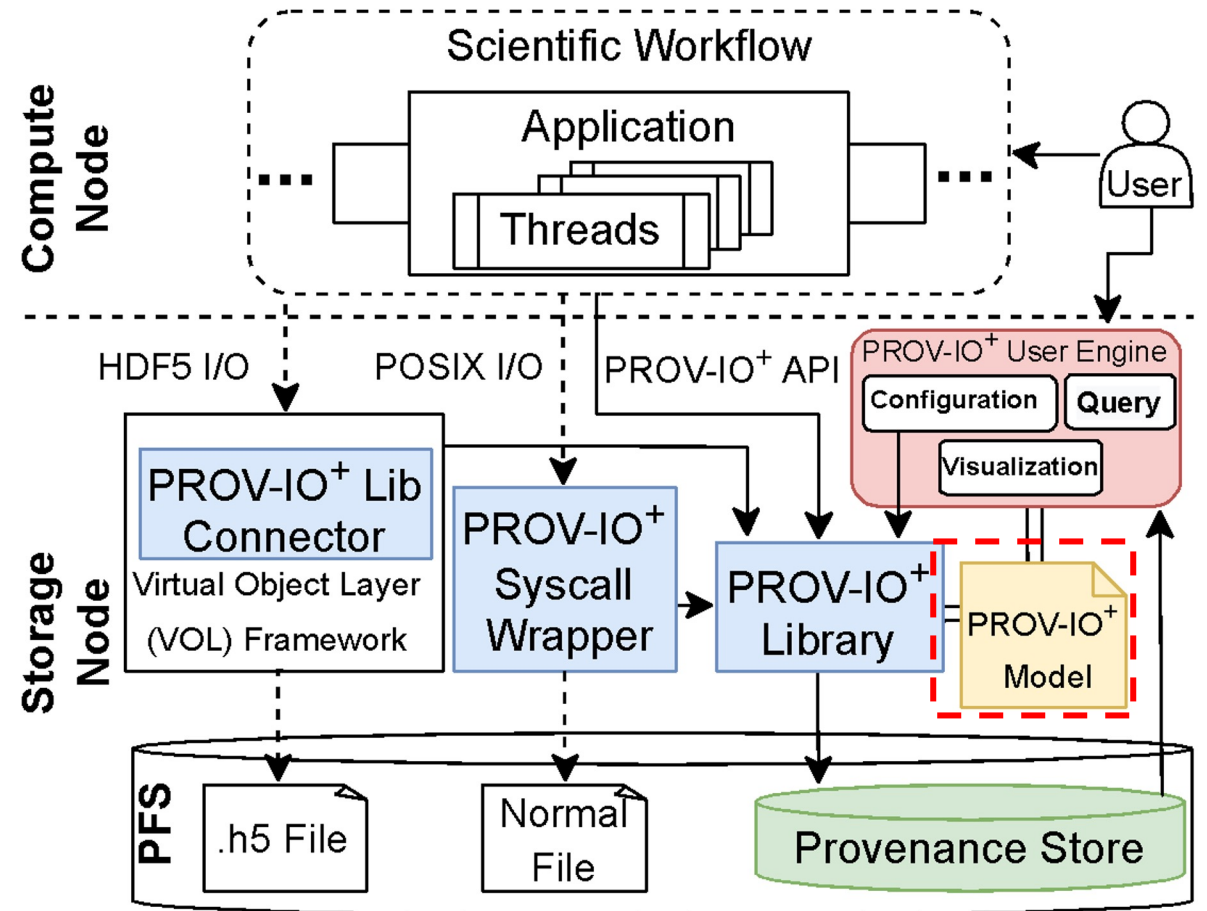
# Design PROV-IO<sup>+</sup> Model Based on Needs

- PROV-IO<sup>+</sup> model: a PROV-DM-compliant data model
  - Provides an interface for users to extend the PROV-IO<sup>+</sup> model with new concepts/relations per their needs



# PROV-IO<sup>+</sup> Framework

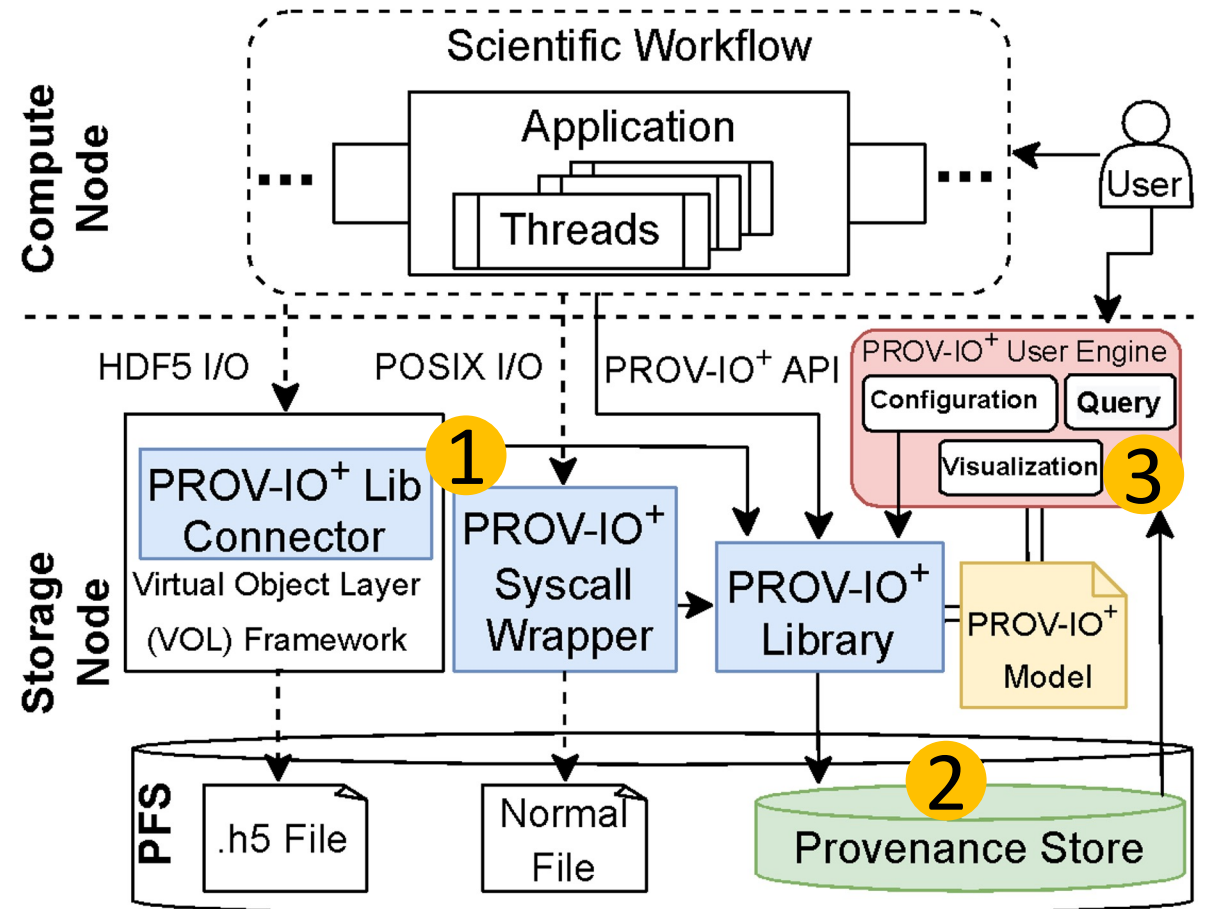
- Overview
  - PROV-IO<sup>+</sup> model



# PROV-IO<sup>+</sup> Framework

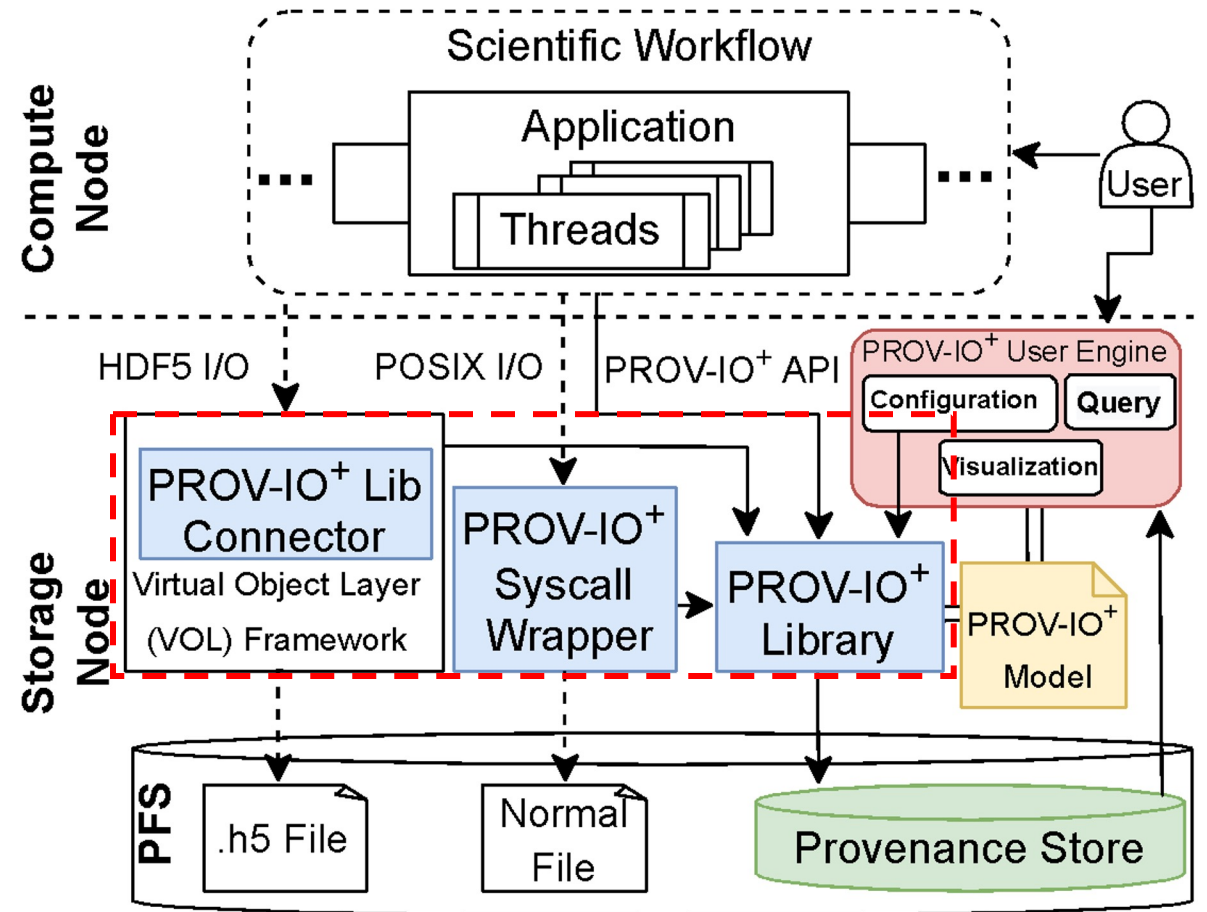
- Overview

- PROV-IO<sup>+</sup> model
- **Three main components based on PROV-IO<sup>+</sup> model**
  1. Tracking (blue)
  2. Store (Green)
  3. User engine (Red)



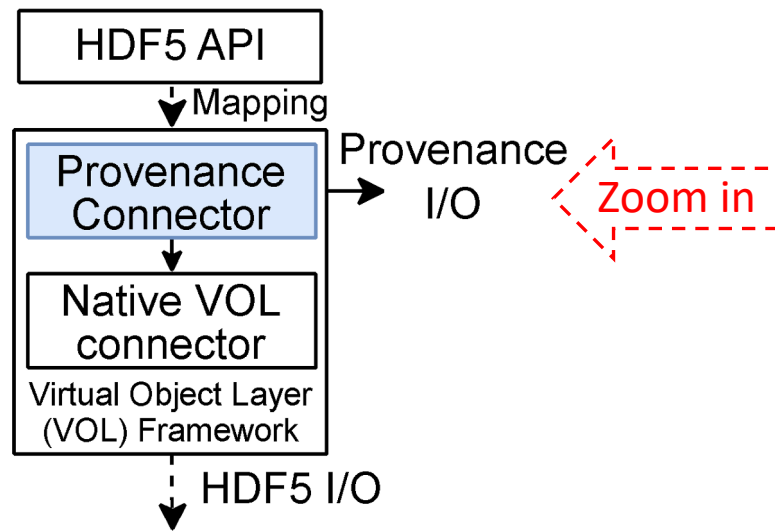
# PROV-IO<sup>+</sup> Framework

- Provenance tracker (blue)
  - Track I/O operations transparently by intercepting library I/O

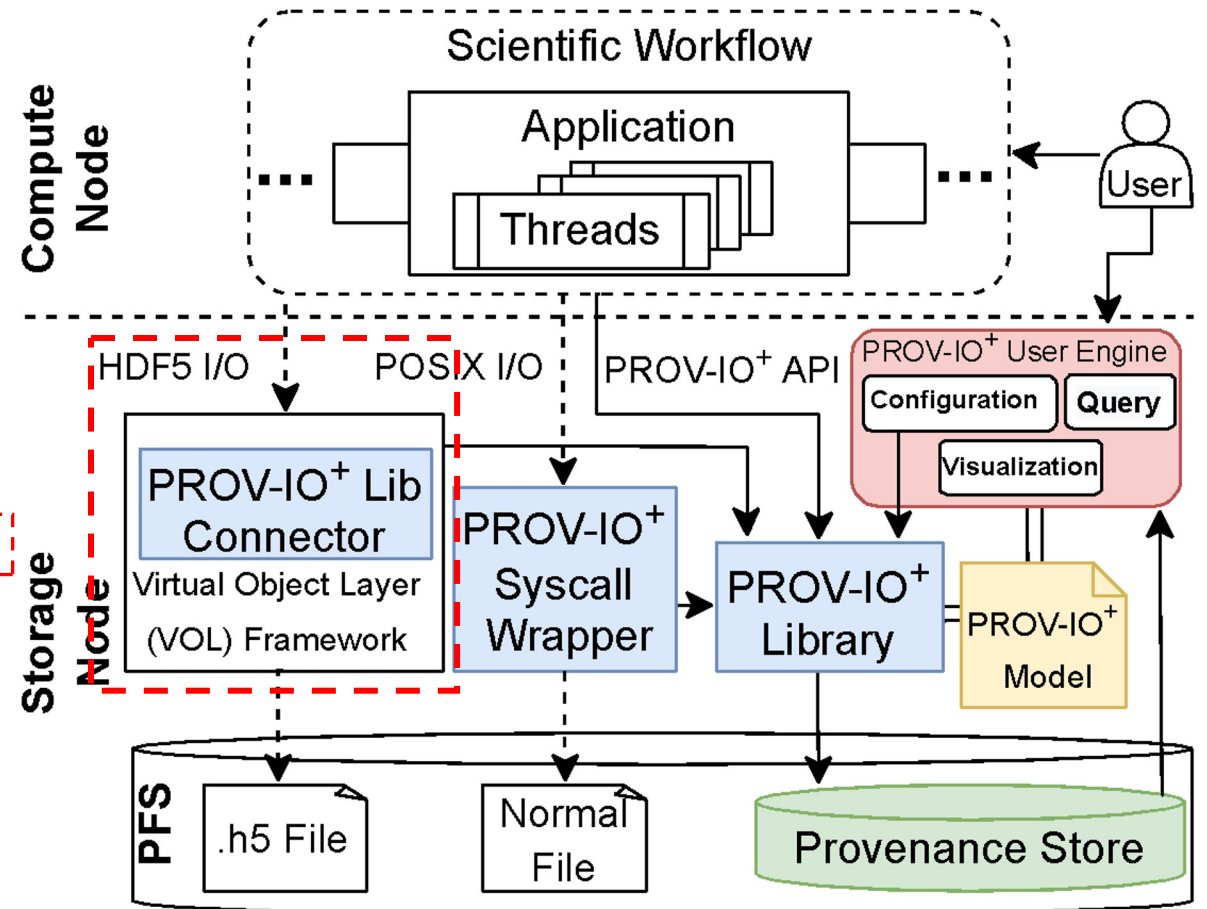


# PROV-IO<sup>+</sup> Framework

- Provenance tracker (blue)
  - Track I/O operations transparently by intercepting library I/O

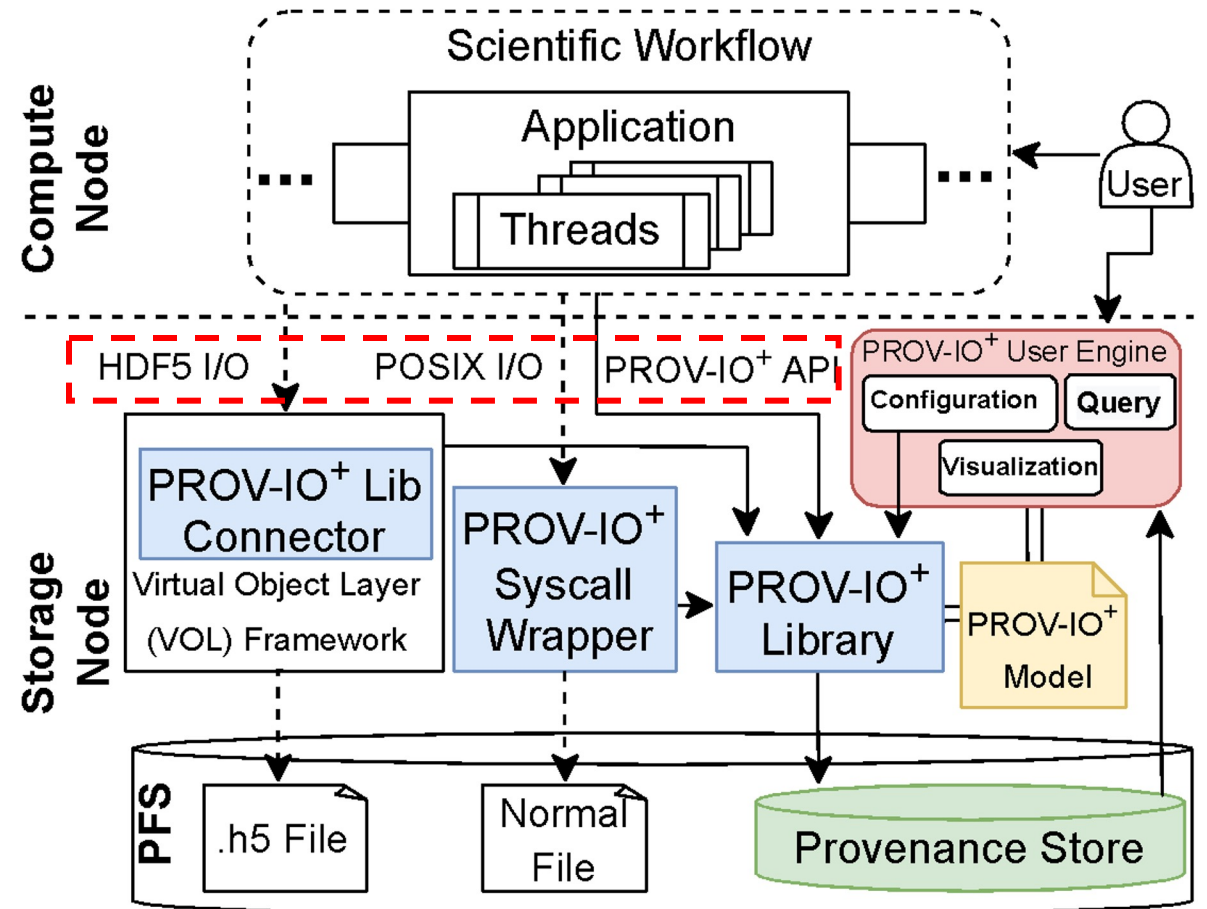


HDF5 Provenance VOL connector



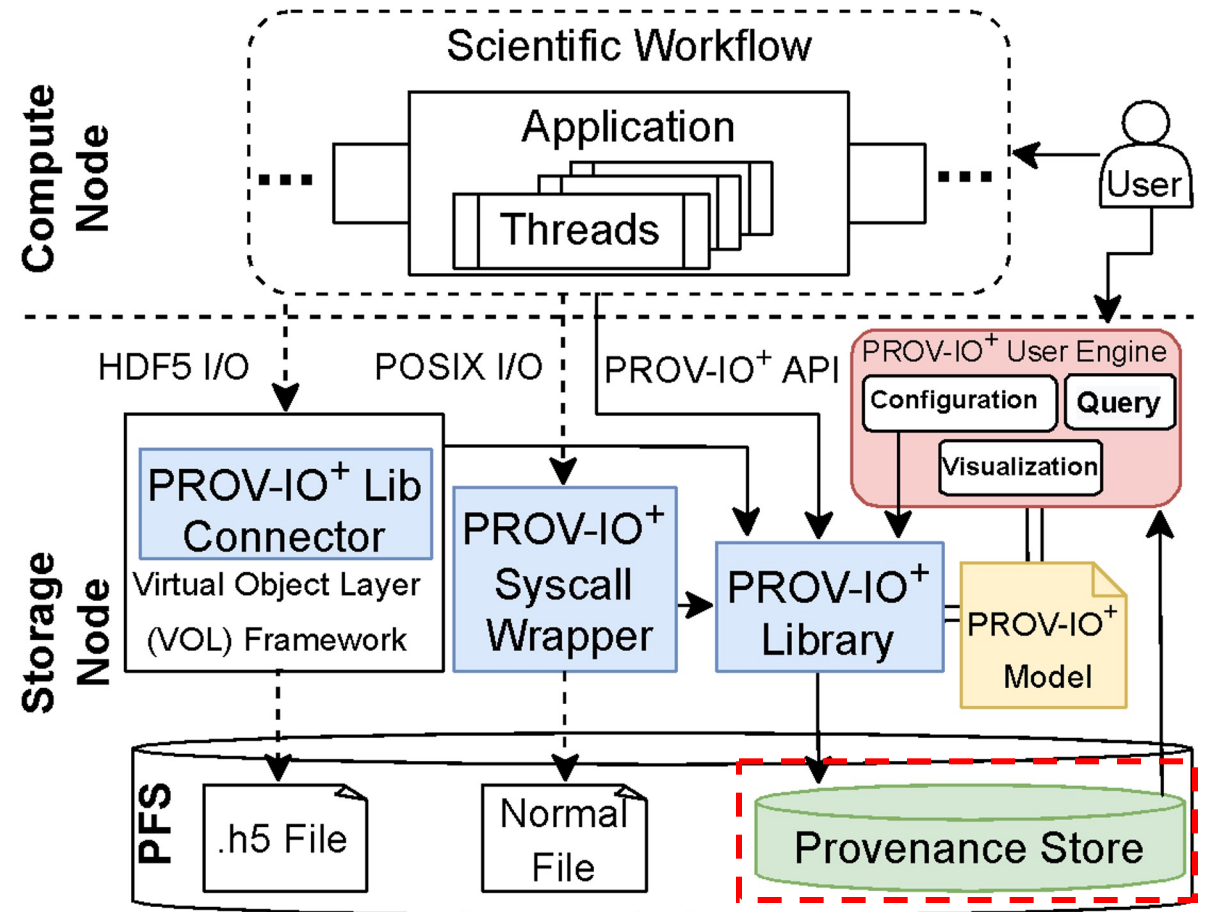
# PROV-IO<sup>+</sup> Framework

- Provenance tracker (blue)
  - Track I/O operations transparently by intercepting library I/O
  - Support for popular I/O libraries use by HPC workflows
  - Provide a Python interface for manually instrumentation



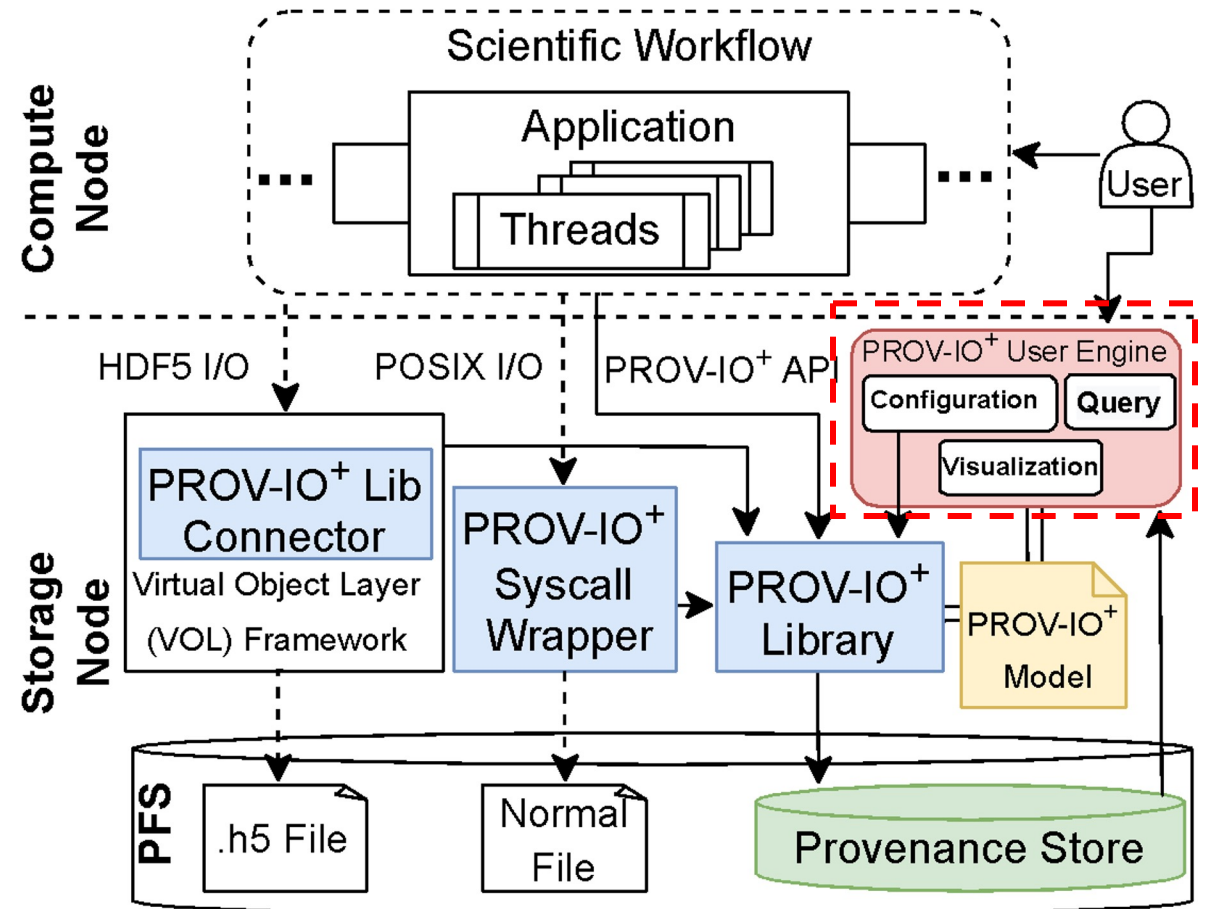
# PROV-IO<sup>+</sup> Framework

- Provenance store (green)
  - Serialize provenance as RDF
  - Avoid concurrent provenance serialization by having each thread write to its own file
  - Consolidate provenance files offline



# PROV-IO<sup>+</sup> Framework

- User engine (red)
  - A configuration interface for user
  - Provenance query with SPARQL
  - Provenance visualization

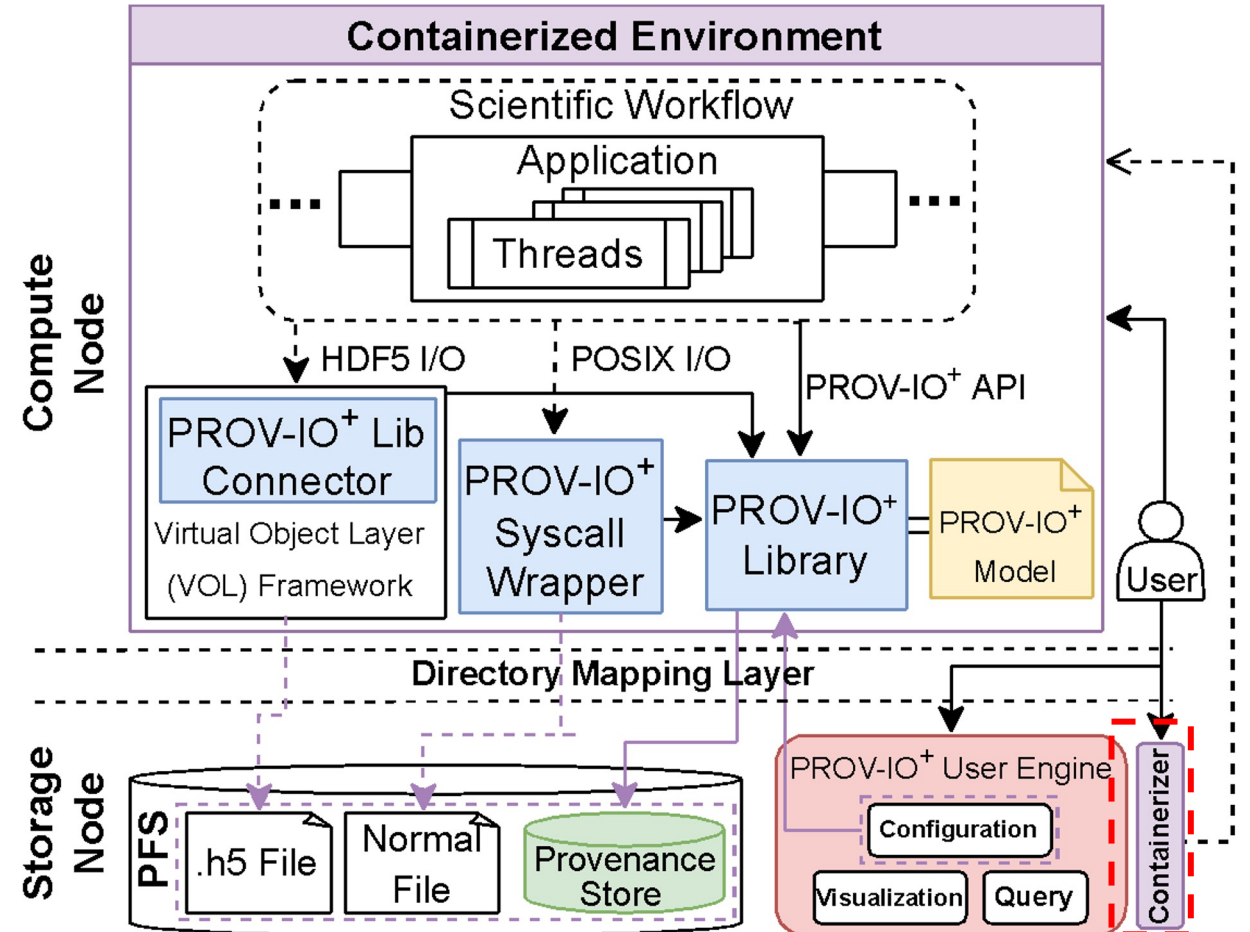




# PROV-IO<sup>+</sup> Framework

- Overview

- PROV-IO<sup>+</sup> model
- Three main components based on PROV-IO<sup>+</sup> model
- **Support for containerized environment**
  - Newer HPC systems may have containerized job management (e.g., Singularity)

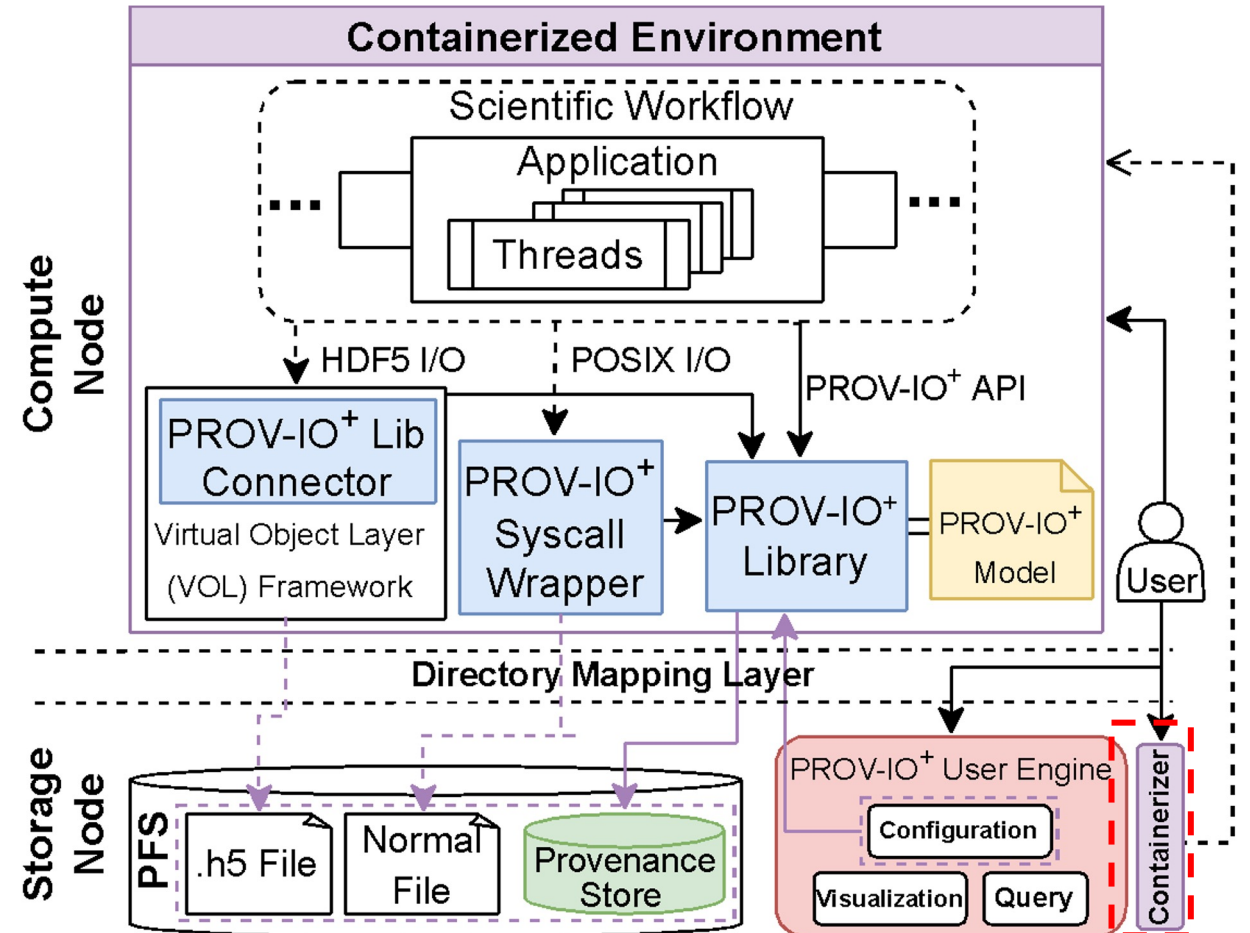


# PROV-IO<sup>+</sup> Framework

- Overview

- PROV-IO<sup>+</sup> model
- Three main components based on PROV-IO<sup>+</sup> model
- **Support for containerized environment**
  - Newer HPC systems may have containerized job management (e.g., Singularity)

For more design/implementation details, please refer to PROV-IO+ paper at: <https://arxiv.org/abs/2308.00891>



# Evaluation

---

# Experimental Methodology

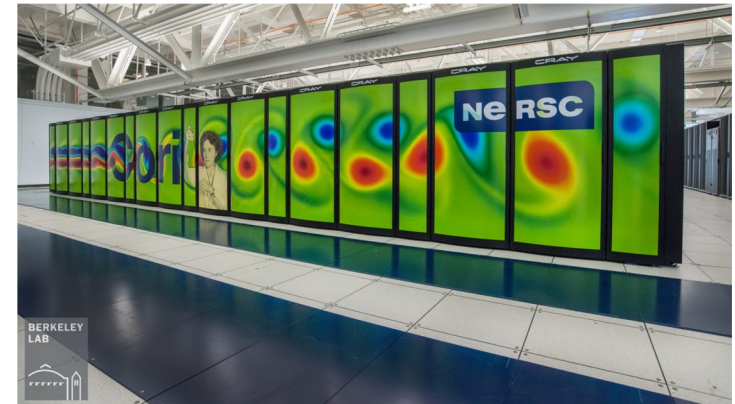
- Platforms

- Cori @LBNL (traditional workflows)

- Up to 64 nodes (4096 MPI ranks)
    - Measured workflows
      - Top Reco (GNN for physical emulation)
      - DASSA (acoustic sensing)
      - H5Bench (synthetic)

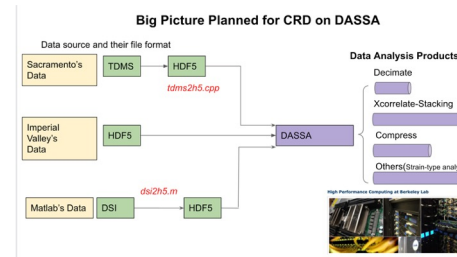
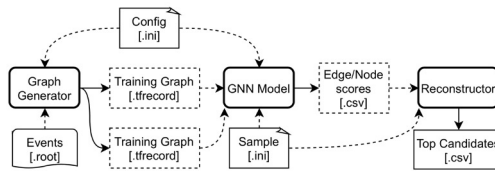
- Samsung supercomputer (containerized workflows)

- 8 A100 GPUs (due to strict quota)
    - Measured workflow
      - Megatron-LM (large language model)



# Experimental Methodology

- Information tracked for each workflow



NVIDIA/Megatron-LM

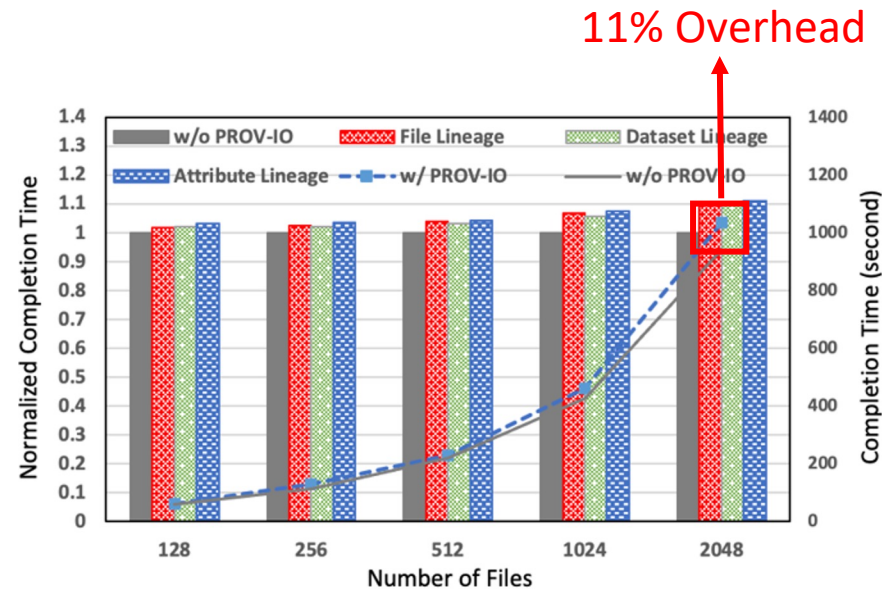


Ongoing research training transformer models at scale

Name	Top Reco	DASSA	H5Bench	Megatron-LM
Domain	GNN for physical emulation	acoustic sensing	synthetic, performance benchmarking	large language model
Provenance need	metadata version control	lineage of data products	I/O stats & performance bottleneck	checkpoint-config. consistency
Information tracked	<ol style="list-style-type: none"> <li>hyperparameter</li> <li>data preselection</li> <li>training accuracy</li> </ol>	<ol style="list-style-type: none"> <li>program name</li> <li>I/O API (HDF5)</li> <li>file/dataset/attr</li> </ol>	<ol style="list-style-type: none"> <li>I/O API (HDF5)</li> <li>I/O API duration</li> <li>user/rank/program/file</li> </ol>	<ol style="list-style-type: none"> <li>checkpoint info</li> <li>training loss</li> <li>model configuration</li> </ol>

# Tracking Overhead

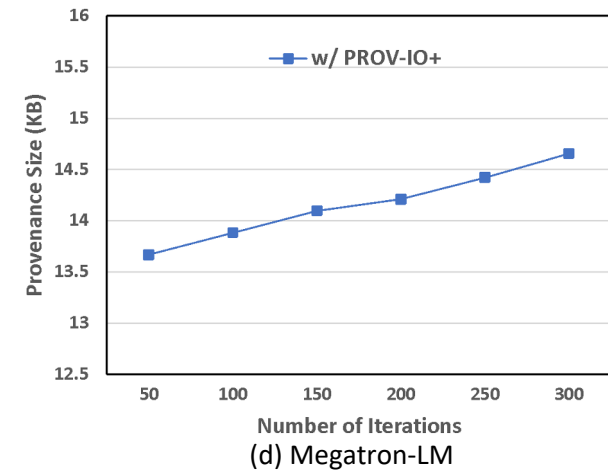
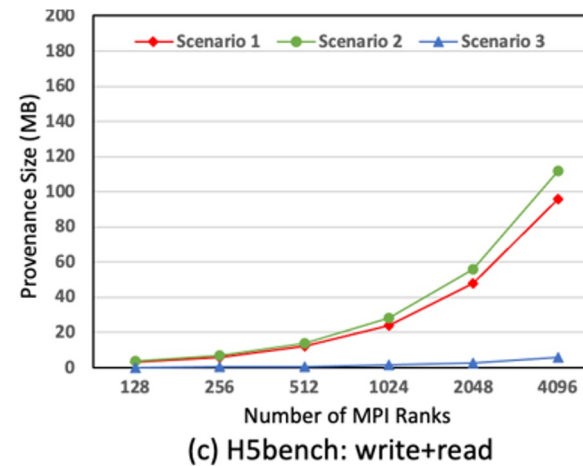
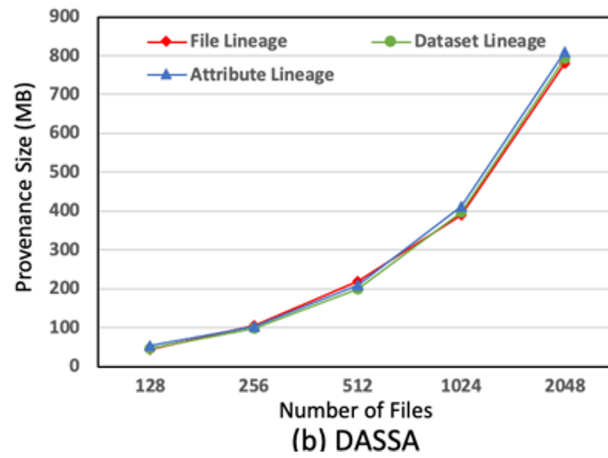
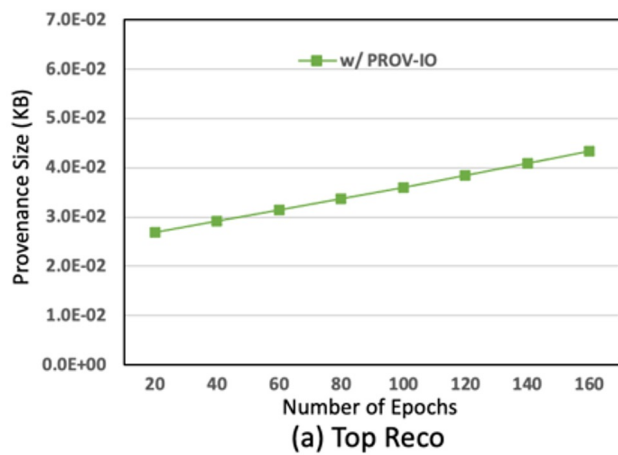
- Over all experiments, tracking overhead is **11% at maximum**
- More than **97%** of the experiments has overhead **less than 3%**



Tracking overhead on DASSA workflow

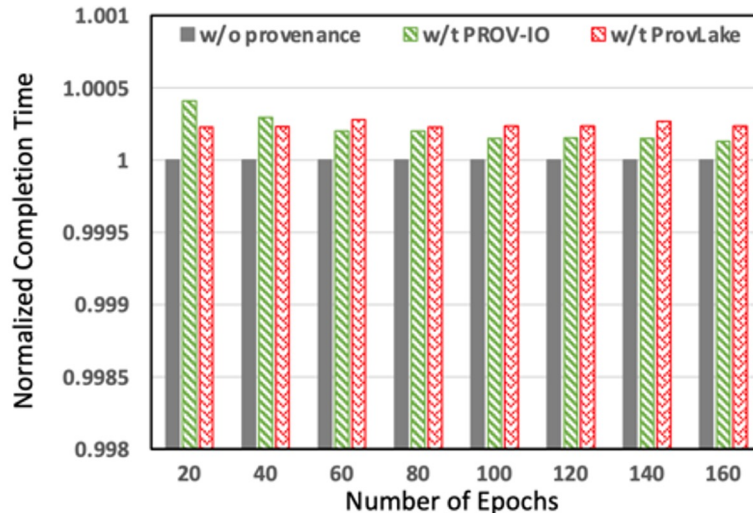
# Storage Overhead

- Provenance size *increases linearly* with experimental scale

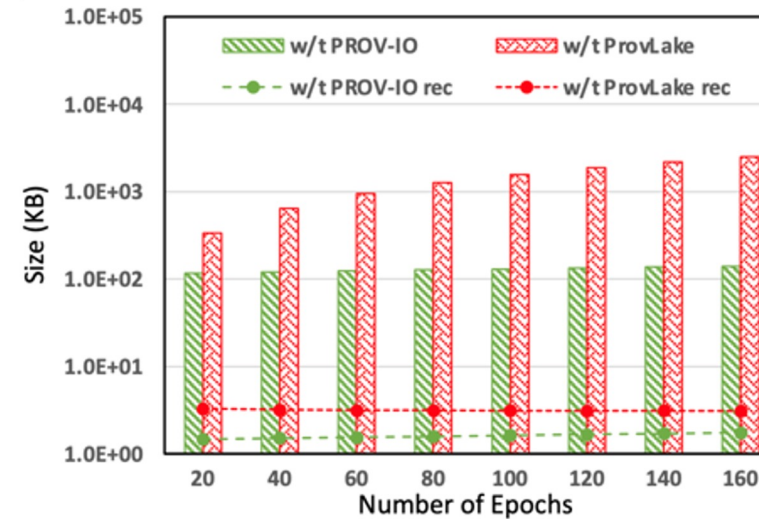


# Comparison with IBM ProvLake

- PROV-IO<sup>+</sup> has *lower tracking overhead* in experiments *with more training epochs*
- PROV-IO<sup>+</sup> always has *less storage overhead*



(a) Tracking overhead



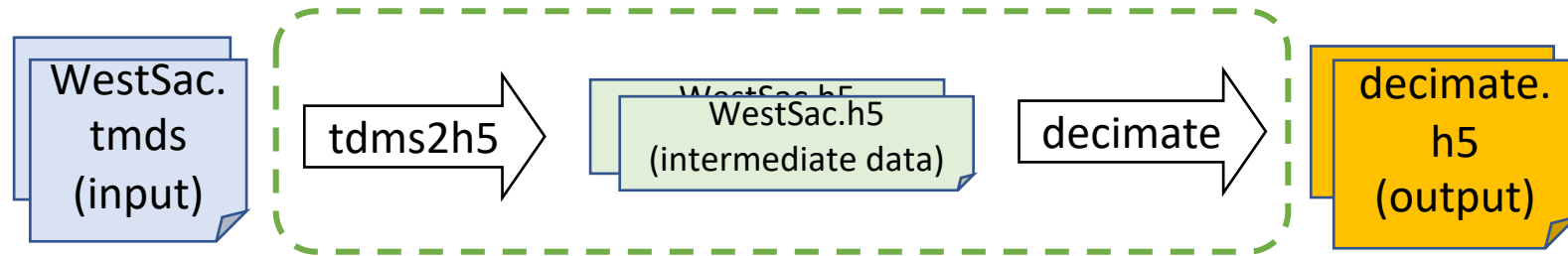
(b) Storage overhead

Comparison on Top Reco Workflow



# Provenance Query & Visualization

- Data lineage backward tracing example with DASSA

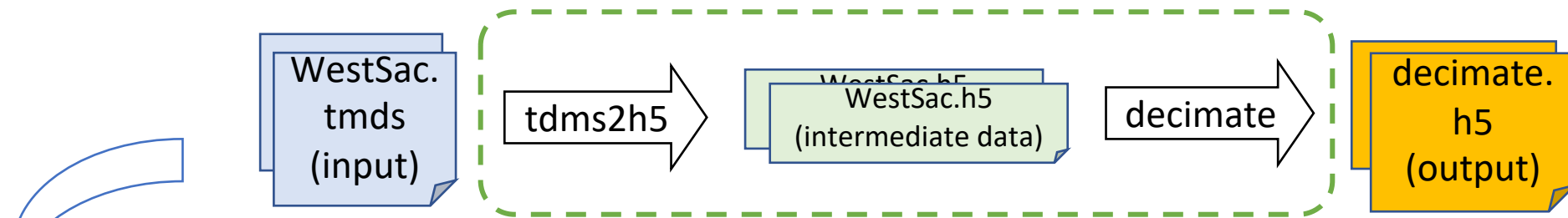


What's the origin  
of the output data  
**decimate.h5**?

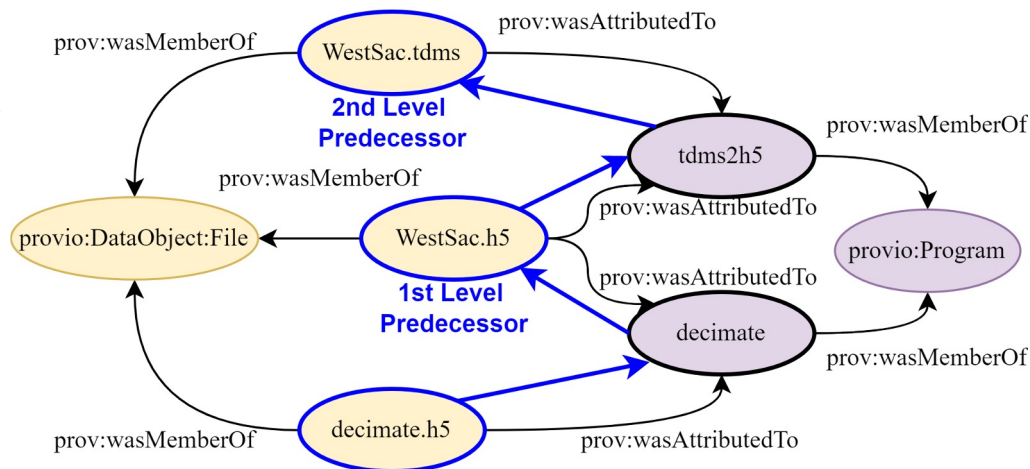


# Provenance Query & Visualization

- Data lineage backward tracing example with DASSA



**Visualized Provenance**



Example SPARQL query to locate 1<sup>st</sup> level predecessor:

**Q1: Decimate.h5 prov:wasAttributedTo ?**  
**Q2: ? prov:wasAttributedTo decimate**  
**provio:wasReadBy H5Fopen**

2<sup>nd</sup> level predecessor can be obtained with similar query.

# Conclusion & Future Work

- Conclusion
  - Identified domain scientists' real provenance needs
  - Built PROV-IO<sup>+</sup> framework under the guidance of PROV-IO<sup>+</sup> model
  - Evaluated PROV-IO<sup>+</sup> framework on two HPC systems
    - PROV-IO<sup>+</sup> can address domain scientists' concerns effectively & efficiently
- Future work
  - More efficient provenance post processing
  - Advanced query API to help users analyze workflows more efficiently

# Conclusion & Future Work

- Conclusion
  - Identified domain scientists' real provenance needs
  - Built PROV-IO<sup>+</sup> framework under the guidance of PROV-IO<sup>+</sup> model
  - Evaluated PROV-IO<sup>+</sup> framework on two HPC systems
    - PROV-IO<sup>+</sup> can address domain scientists' concerns effectively & efficiently
- Future work
  - More efficient provenance post processing
  - Advanced query API to help users analyze workflows more efficiently

Source code is available at:

<https://github.com/data-storage-lab/prov-io>

Docker image is available at:

<https://hub.docker.com/repository/docker/rzhan/prov-io/general>

**Thank You  
&  
Questions?**

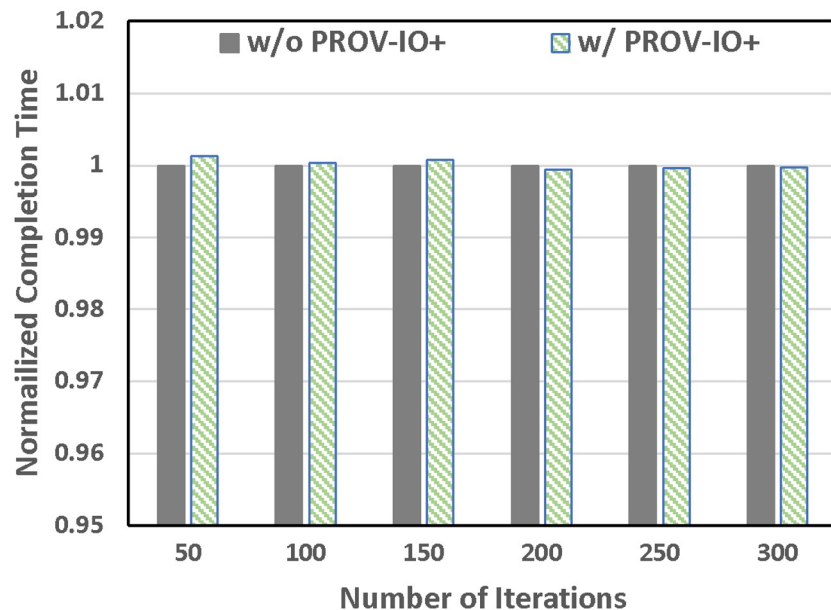


# Backup Slides

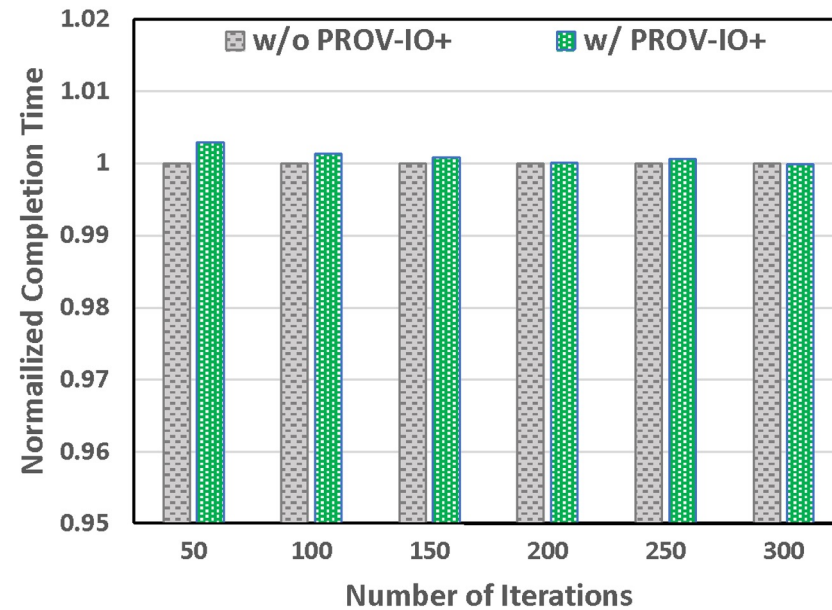
---

# Containerization Overhead

- *Negligible containerization overhead* observed in Megatron-LM use case on Samsung Supercomputer



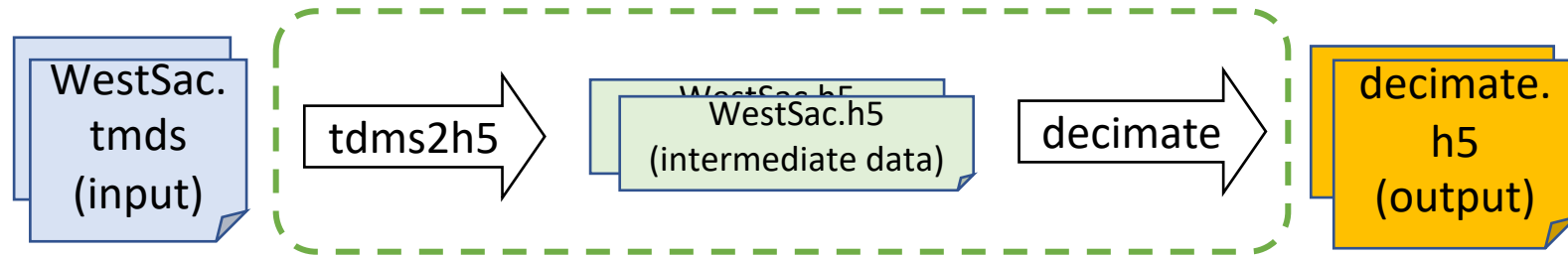
(a) Non-containerized workflow



(b) Containerized workflow

# How to Query Provenance

- Data lineage backward tracing example with DASSA

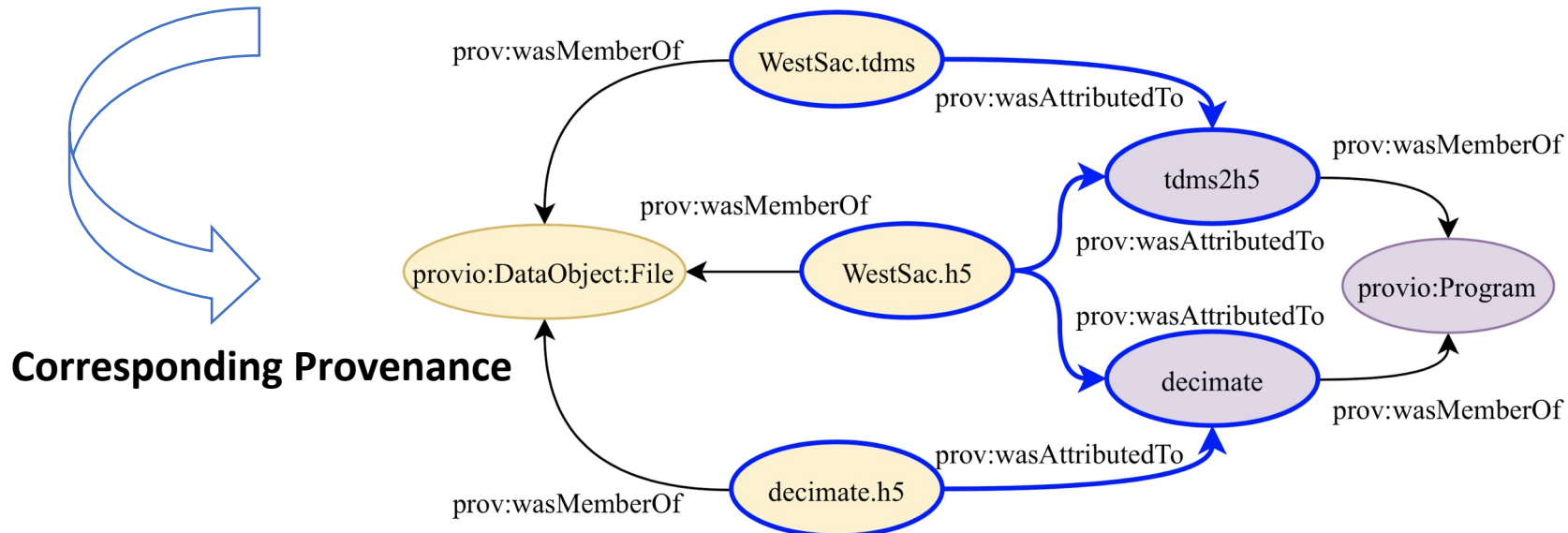
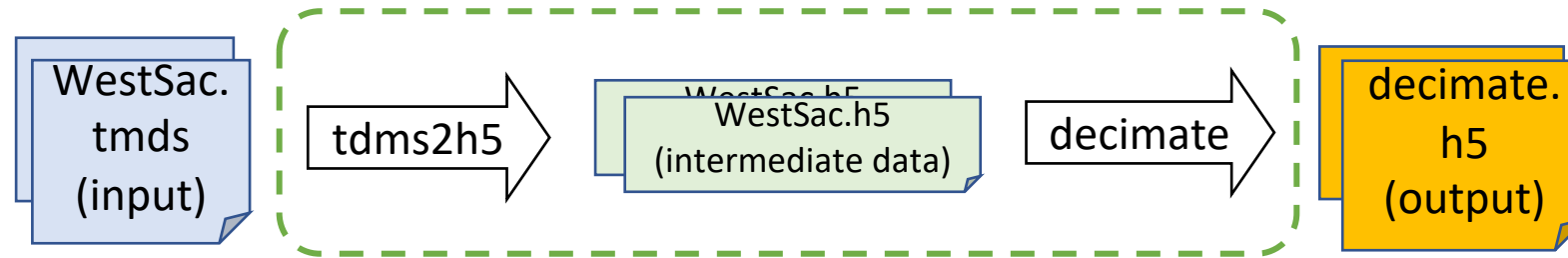


What's the origin  
of the output data  
**decimate.h5**?



# How to Query Provenance

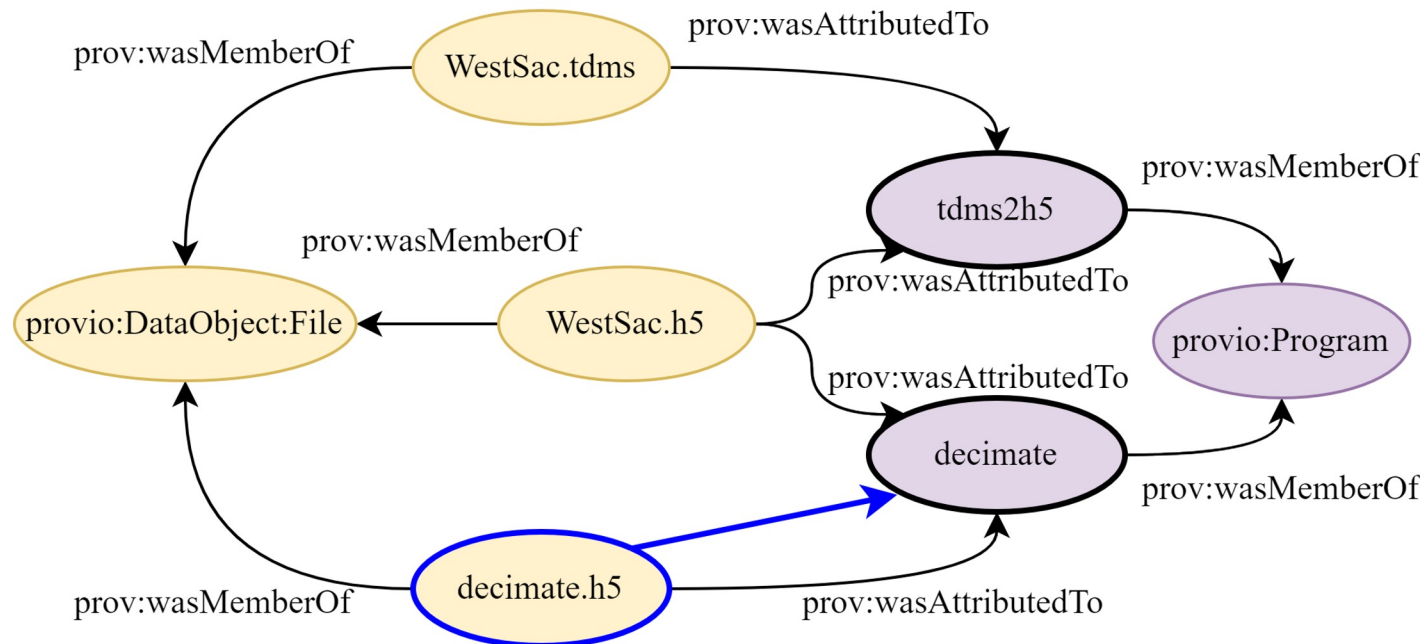
- Data lineage backward tracing example with DASSA





# How to Query Provenance

- Data lineage backward tracing example with DASSA



**Step 1.** Search the program whose output is decimate.h5

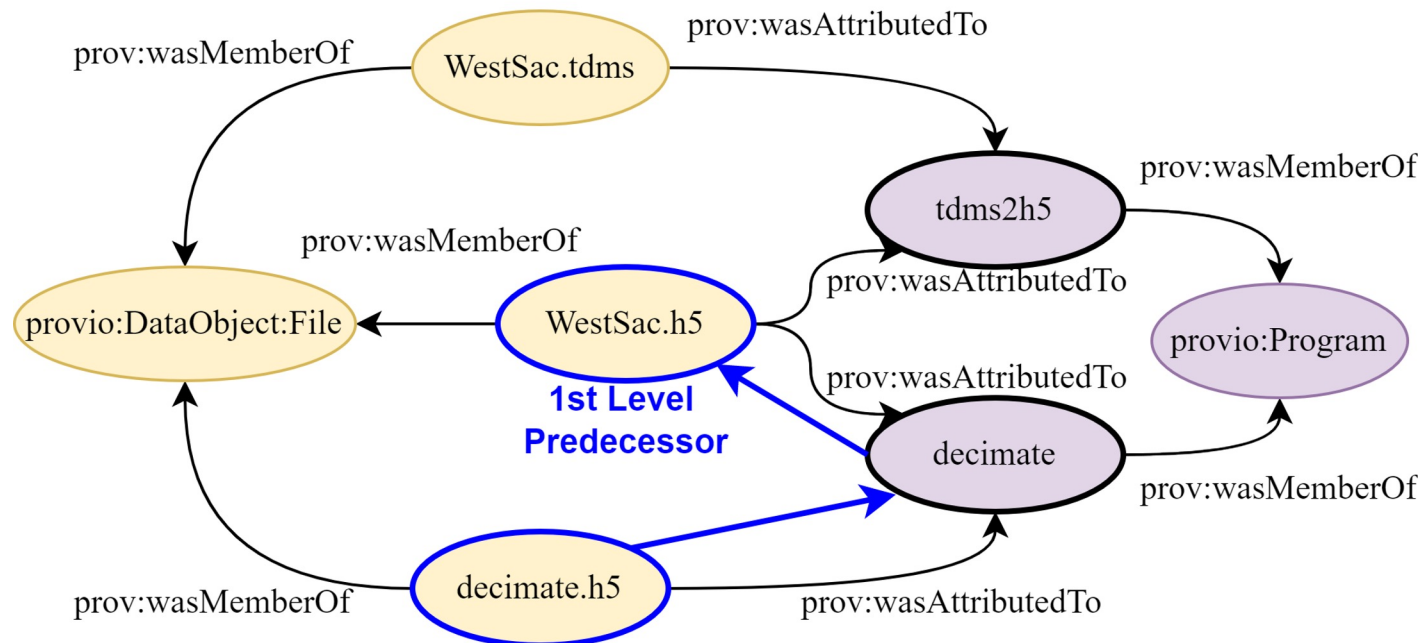
Example SPARQL query:

**Decimate.h5 prov:wasAttributedTo**

?

# How to Query Provenance

- Data lineage backward tracing example with DASSA



Step 2. Search the input data of decimate

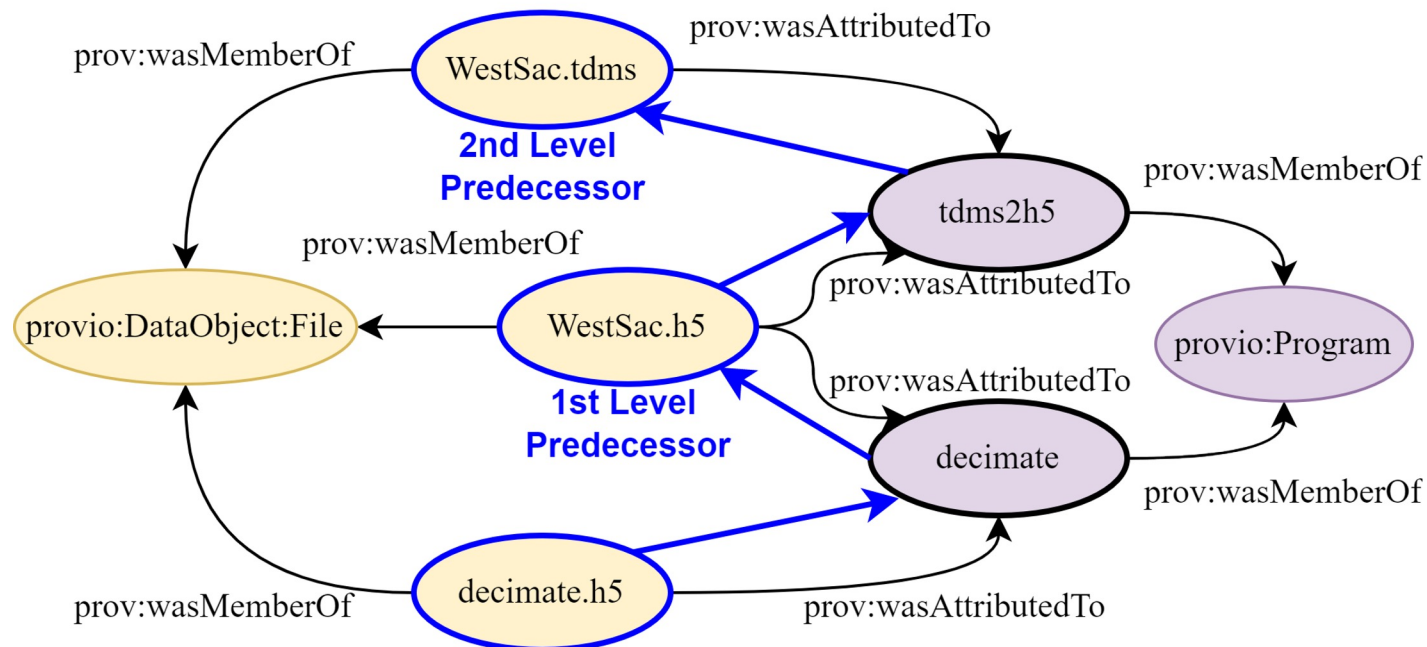
Example SPARQL query:

**? prov:wasAttributedTo  
decimate**

prov:wasReadBy **HEFopen**

# How to Query Provenance

- Data lineage backward tracing example with DASSA



**Step 3&4.** Repeat Step 1&2. Look for the program which created WestSac.h5 and then search the input of that program

LoC of query =  $N * 3$   
( $N$  is level of predecessor data object)