# Study on Handling Dark Data in HPCI Shared Storage System using the WHEEL Workflow Tool

International Data Week 2025
Data and Research & Data Science and Data Analysis
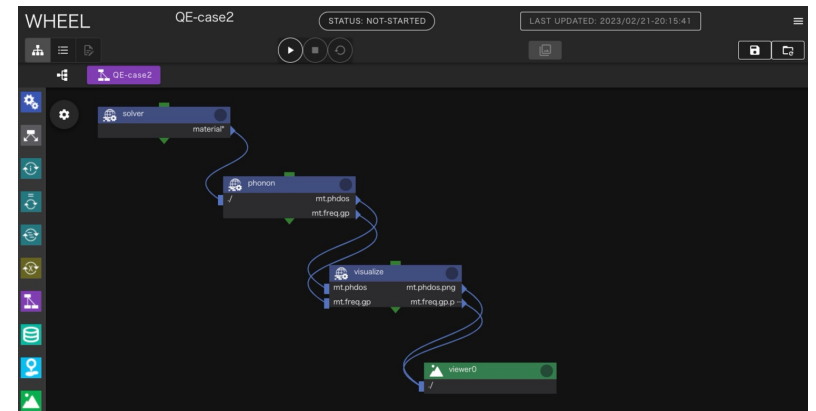2025/10/13

**RIKEN R-CCS**
**- Hidetomo Kaneyama**
**- Tomohiro Kawanabe**
**- Hiroshi Harada**

# Outline



- RIKEN operates the HPCI Shared Storage, providing free storage for HPC/supercomputing research data in Japan.

- Currently storing "35+ PB" and "200M+ research files".

- Many files have not been accessed for over a year and are now considered "cold data.

- Dark data: data unused for analysis/decision-making, with unclear value or risk. And, significant dark data exists in HPCI Shared Storage.

  - https://www.gartner.com/en/information-technology/glossary/dark-data

- Issues:

  - Dark Data blocking "effective use" of storage, causes storage pressure.

- Our approach:

  - Extend WHEEL workflow tools for HPC/supercomputers.

    - Automatically capture information and automatically add extended metadata to HPC datasets

    - Data and extended metadata is stored in HPCI shared storage.

# What is HPCI-SS System?

**HPSS(HPCI Shared Storage) System for Preserving and Sharing Data from Japanese Supercomputers.**

- **Objectives**
  - Data sharing between supercomputers (computational resources)
  - Long-term preservation of research data
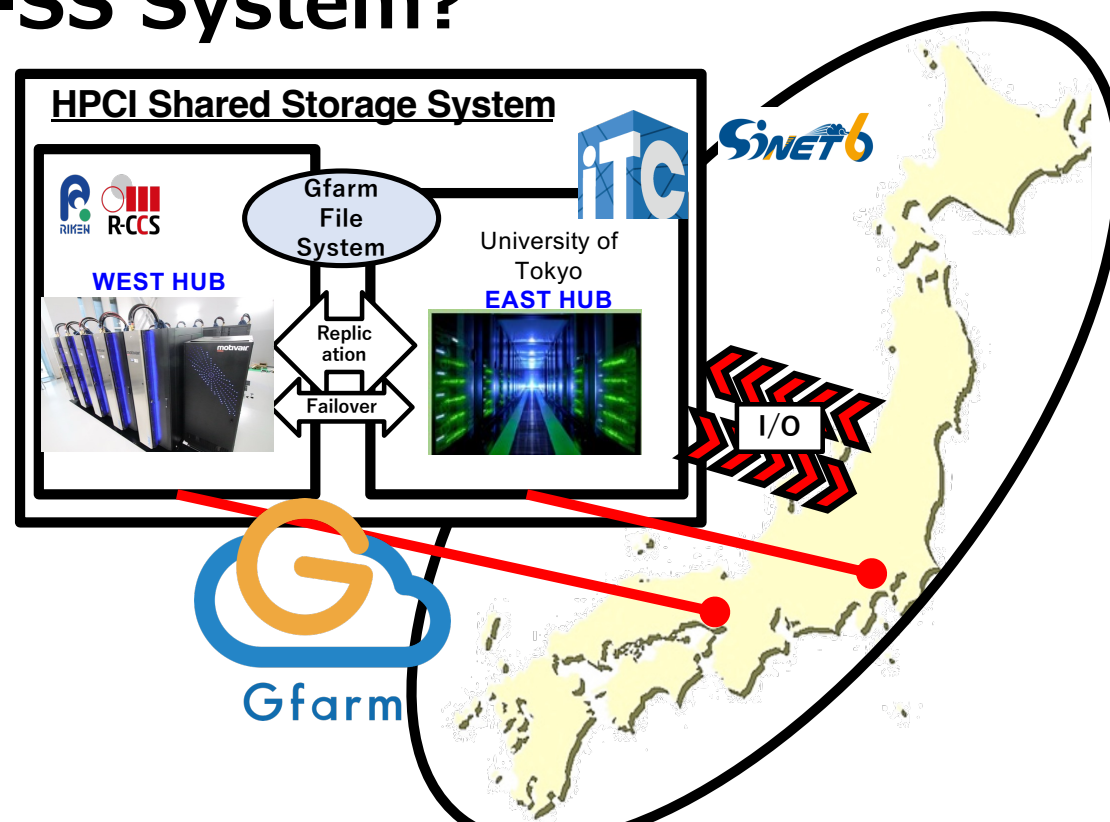  - Public dissemination of research data (utilization of open/public datasets)

- **Core Software**
  - **Gfarm File System**
    https://github.com/oss-tsukuba/gfarm

- **Key Features**
  - Parallel data I/O transfer from HPCI resources
  - Disaster recovery through inter-site data replication
  - High availability (operational uptime >99%, downtime/unavailability <1%)
  - High-performance network storage capable of transfers exceeding 200 Gbps



| Gen | Periods | Capacity（Logical） |
|---|---|---|
| 1st | FY2012 – FY2018 | 10PB → 15PB |
| 2nd | FY2018 - FY2024 | 45PB → 50PB |
| 3rd(now) | FY2025 - | 95PB+ |

https://www.hpci-office.jp/info/pages/viewpage.action?pageId=111380786

# Cold Data in HPCI-SS



**Number of Files**



**Capacity**

※**Cold data is 'File not Accessed in over ONE YEAR'**

| Total | Capacity | 33PB |
|---|---|---|
| | Number of files | 188 million |
| Cold data ratio (include Dark data) | Capacity | 90.2% |
| | Number of files | 88.8 % |

Cold data grow in HPCI Shared Storage.
- Users face heavy burden in identifying unnecessary files
- Some data has been inherited (e.g., from former members)
- Lack of metadata makes deletion decisions difficult

4

# Discussion and Approach to Dark Data in HPCI-SS

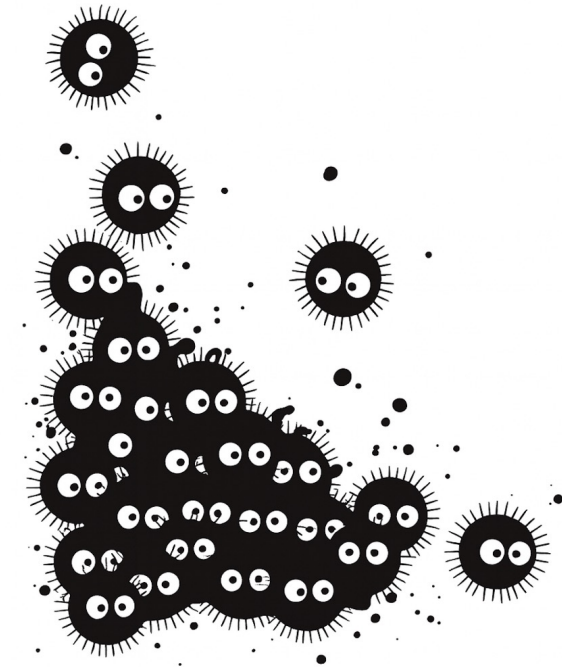**CJK Project(A3)**

**Operational Concerns**
→ **Need to provide environments for extended metadata attachment / methods for automatic enrichment**

**[1] Approach to Cold Data**

- **Ask users to delete very old, unused data**

- **Dark data is harder: unidentified data, high user burden to judge deletion**

- **Actual volume of dark data remains unknown (only discovered when users clean up)**

- **Visualization provides usage/capacity/file counts, but not data content or usefulness**

**[2] Challenges for AI & Open Science**

- **AI/Open Science requires labels and metadata**

- **Current HPC & HPCI storage lacks tagging/labeling mechanisms**

- **Object storage often missing or underutilized**

- **Critical info ("who used what, and how data was obtained") is not preserved**

- **Some users keep records manually (e.g., Excel) → high management burden**

https://sca25.sc-asia.org/wp-content/uploads/2025/04/pos121s1.pdf

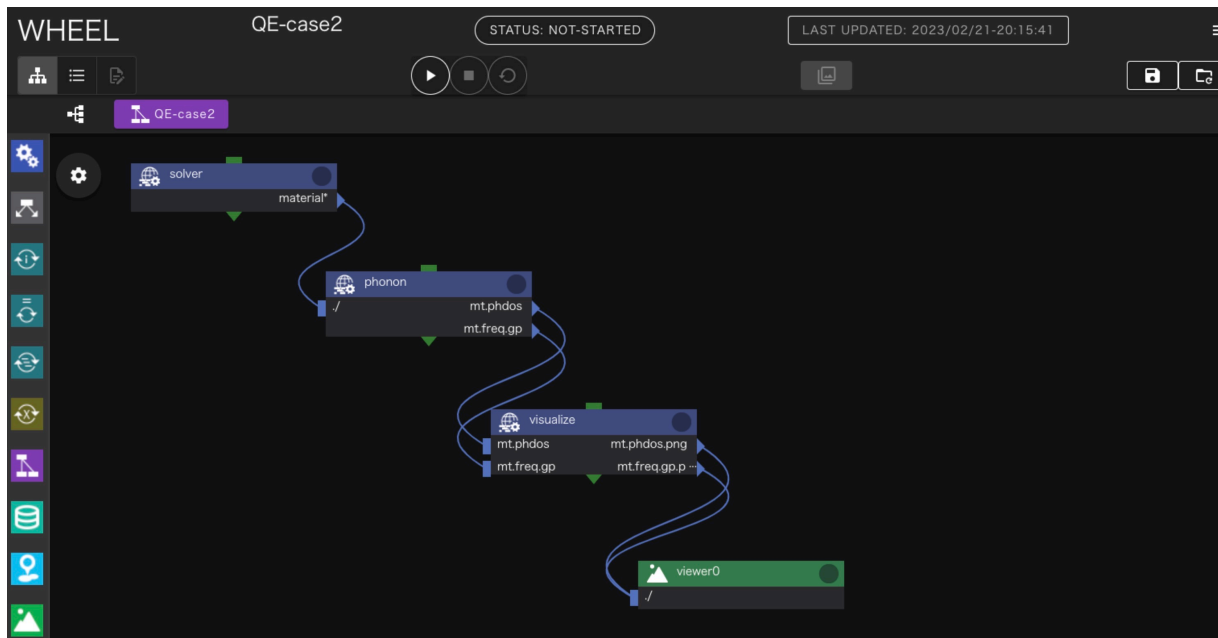# Discussion and Approach to Dark Data in HPCI-SS

- **Why Workflow Tools?**
- Data in HPCI-SS is mainly generated on external HPC/supercomputers
- Data management starts **at the point of generation**
- HPC storage (e.g., Lustre, BeeGFS) is fast but lacks metadata functions
- Replacing entire storage systems is **too costly**
- Commercial metadata products (e.g., Starfish, IBM AFM) require:
  - Management authority at each site
  - High license and operational costs
- → Implement metadata management **on the workflow side**
- **What does this solve?**
- Automatically collect information during computation, pre-/post-processing
- Associate **extended metadata** with datasets
- Compress and register datasets in HPCI-SS as manageable units
- Enable metadata-based **search and management** in HPCI-SS (via Gfarm)

# Overview of WHEEL

- **A web-based GUI workflow build and execution tool**

- **JavaScript application using Node.js**

- **Open-source software distributed under the BSD-2 license**

  - **https://github.com/RIKEN-RCCS/OPEN-WHEEL**

# Method

[A] Automatic extended metadata in HPC environments

- Adding functions to job schedulers or HPC storage (e.g., Lustre) is outside our authority

- Commercial tools (e.g., Starfish) involve high licensing/operational costs → Requires discussion at HPCI-wide level
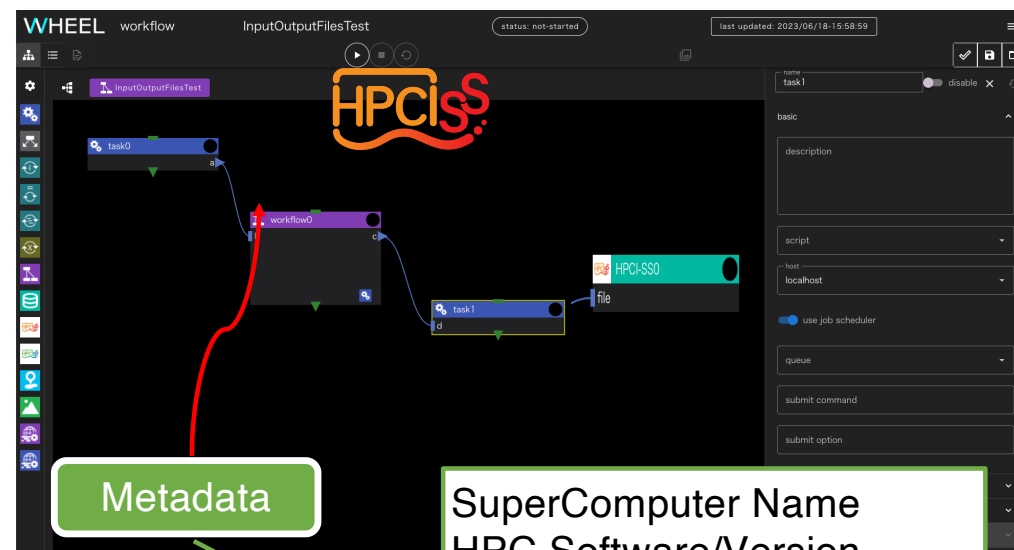
[B] Management software on HPCI-SS side

- 3rd-generation system already designed (capacity-first policy)

- High cost, difficult to shift to products like VAST Catalog / IBM AFM

Our Approach

- Seek improvements without direct system modifications

- Collaborate with WHEEL workflow tool (R-CCS, Senior Engineer Kawanabe)

- Automatically collect extended metadata during workflow execution

- Store enriched metadata + datasets into HPCI-SS (Gfarm)
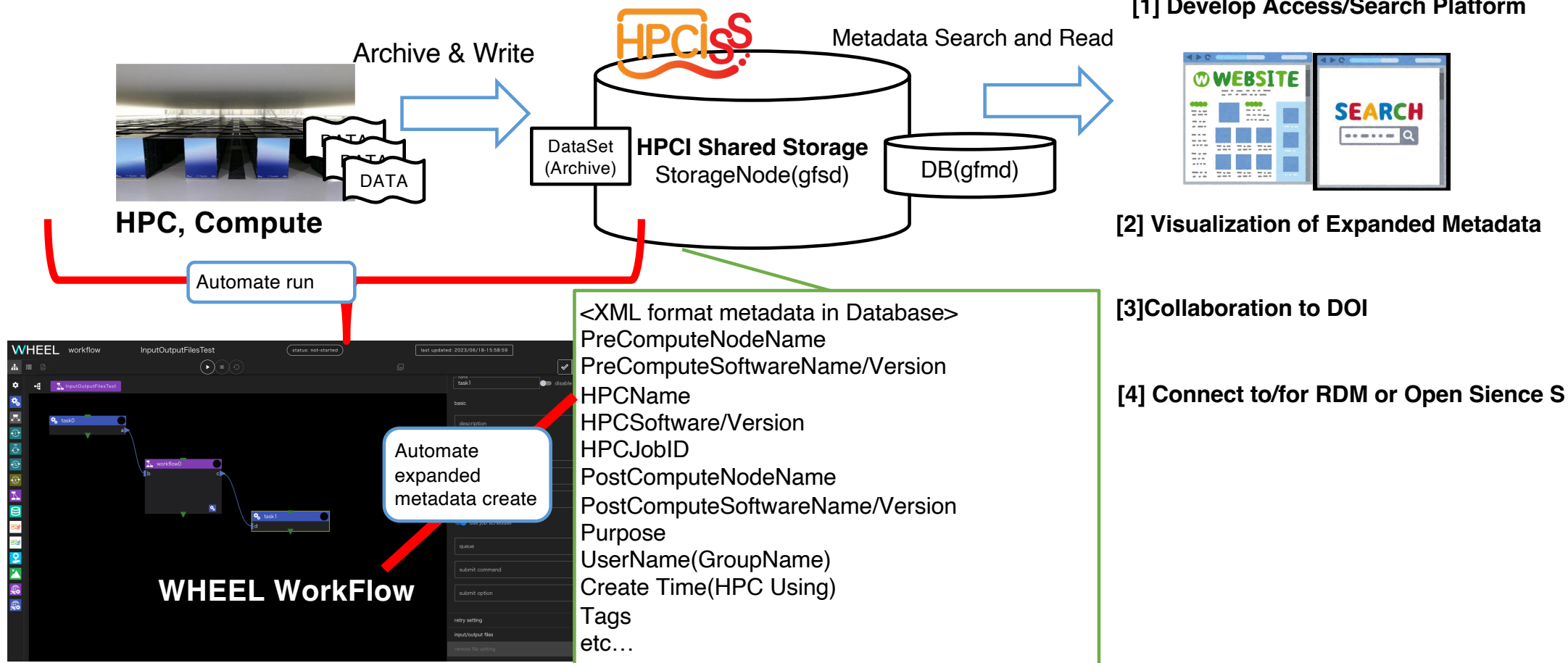
- Long-term goal: integrate with DOI assignment

Progress

- FY2024: Enabled WHEEL I/O to HPCI-SS via Gfarm API

- FY2025: Developing automated metadata collection during HPC compute & pre/post-processing



Metadata

SuperComputer Name
HPC Software/Version
JobID
Tags
Purpose
UserName(GroupName)
Create Time(HPC Using)
Etc…

https://sca25.sc-asia.org/wp-content/uploads/2025/04/pos121s1.pdf

# Next Work

**Future challenges**

**[1] Develop Access/Search Platform**

Archive & Write

Metadata Search and Read

**HPCI Shared Storage**
**StorageNode(gfsd)**

DataSet (Archive)

DB(gfmd)

**HPC, Compute**

DATA
DATA
DATA

Automate run

**[2] Visualization of Expanded Metadata**

**[3]Collaboration to DOI**

**[4] Connect to/for RDM or Open Sience S**

Automate expanded metadata create

**WHEEL WorkFlow**

<XML format metadata in Database>
PreComputeNodeName
PreComputeSoftwareName/Version
HPCName
HPCSoftware/Version
HPCJobID
PostComputeNodeName
PostComputeSoftwareName/Version
Purpose
UserName(GroupName)
Create Time(HPC Using)
Tags
etc…

https://sca25.sc-asia.org/wp-content/uploads/2025/04/pos121s1.pdf

9

# SCA/HPCAsia 2026: Call for Submissions

**SCA 2026** Supercomputing Asia — Gathering the **Best of HPC** in Asia

**HPC Asia 2o26**

- Event Overview:
  - Date: January 26-29, 2026
  - Venue: Osaka International Convention Center (Osaka, Japan)
  - Theme: "Everything with HPC –AI, Cloud, QC, and Future Society"

- Call for Submissions: Papers, Posters, Workshops, BoFs, and Tutorials

| Papers | Posters | Workshops | Birds of a Feather | Tutorials |
|---|---|---|---|---|
| Paper abstracts: 29 Aug 2025 Submissions close: 5 Sep 2025 Result notification: 20 Oct 2025 | Submissions close: 27 Oct 2025 Result notification: 14 Nov 2025 | Submissions close: 30 Jun 2025 Result notification: 31 Jul 2025 | Submissions close: 1 Sep 2025 Result notification: 1 Oct 2025 | Submissions close: 11 Jul 2025 Result notification: 15 Aug 2025 |

**For more details, please visit our website:**

https://www.sca-hpcasia2026.jp/

Thank you