# Efficient metadata management with the AMI ecosystem

CHEP, Kraków (19 - 25 Oct 2024)

F. Lambert, P.-A. Delsart, J. Fulachier, J. Odier
ami@lpsc.in2p3.fr

# Outline

1. Metadata

2. The AMI ecosystem

# 1. Metadata

# What are metadata?

"Metadata is data that provides information about other data. It describes the characteristics, content, and context of the data, making it easier to understand, organize, find, and manage."

# Metadata are essential for science

- A thought-provoking article from Nature (vol. 533, 2016):
  - **>70%** of researchers have failed to reproduce another scientist's experiments.
  - **>50%** have failed to reproduce their own experiments.
- Metadata must help make data **FAIR** for reproducible science:
  - **F**indable: The first step in (re)using data is to find them.
  - **A**ccessible: Long-term preservation and easy access to data.
  - **I**nteroperable: Open, widely shared languages and formats to combine metadata.
  - **R**eusable: Metadata must provide information about the origins of the data and the conditions for its reuse.

# AMI and metadata challenges

- AMI (**A**TLAS **M**etadata **I**nterface) is a generic ecosystem dedicated to scientific metadata.
  - Over 24 years of experience within the ATLAS collaboration at CERN
  - Several years with smaller collaborations like NIKA2, n2EDM, and others
- This experience has shaped our vision of how metadata challenges should be addressed:
  - How can physicists efficiently select the data they need?
  - How to deal with heterogeneous sources of metadata?
  - How can metadata help ensure that data can be reused long after the experiment ends (in 5, 20, or even more years)?
- The primary goal of AMI is to help physicists identify the data that will be most useful to them.

# 2. The AMI ecosystem

# AMI ecosystem in a nutshell

- **Front-end**: AWF (**A**MI **W**eb **F**ramework) - Modern JavaScript
  - Controls for building Web applications to select and display data
  - Fully configurable default "search" application

- **Back-end**: AMI Core - Java
  - Microservices providing interoperable outputs like XML, JSON, CSV, etc.
  - Interaction with any kind of datasource (auto-detection of DB structure)

- **Task Server**: A distributed super-CRON
  - Extracting metadata from primary sources (pull mode)
  - (Re)processing and storing metadata in AMI

- **Clients**: Python, C++, Java, JavaScript, etc.

- **Query Language**: MQL (**M**etadata **Q**uery **L**anguage)
  - Designed for non-database experts
  - No need for knowledge of the underlying DB schema

# AMI typical usage



**Figure 1:** A typical data and metadata workflow.

# End-users and AMI

- Varied profiles of end-users with diverse needs:
  - Some users need to write scripts to access data.
  - Others only require simple Web applications or command-line tools.

- Varied levels of expertise among end-users:
  - Some are scientists with development skills, while others are not.
  - Most users are not experts in SQL.
  - They often do not know the structure of the various metadata databases.

- AMI offers solutions tailored to these needs and expertise levels:
  - User-friendly Web applications (point-and-click).
  - Scriptable clients for more advanced users.
  - MQL, a high-level metadata-oriented language, designed for all users.

# MQL: A Metadata-Oriented Language

- MQL handles metadata entities; a dataset is defined by its characteristics.

```
SELECT *
WHERE
  DATASET.STATUS = 'VALID' AND DATASET_KEYWORDS.KEYWORD = '1jet'
```

**Figure 2:** MQL query on dataset entity.

- SQL manages database objects; a dataset is a "table" with fields.

```
SELECT *
FROM DATASET, DATASET_KEYWORDS
WHERE
  DATASET.STATUS = 'VALID' AND DATASET_KEYWORDS.KEYWORD = '1jet'
  AND
  DATASET_KEYWORDS.DATASETFK = DATASET.IDENTIFIER
```

**Figure 3:** SQL query.

# AWF: Search by Criteria Interface

Designed for point-and-click users.



**Figure 4:** Interface for searching by criteria.

# AWF: Search Result Interface

Also designed for point-and-click users.



**Figure 5:** Interface for displaying search results.

# AWF: Search Modeler Interface

A tool for admin users to create search-by-criteria interfaces.



**Figure 6:** Search Modeler interface (admin user).

# Microservices: Interact from the Web

A tool for advanced users to interact using defined "commands".



**Figure 7:** MQL query "command" executed from a Web application.

# Microservices: Interact from a Shell

Advanced users can interact using defined "commands" in a shell.



**Figure 8:** MQL query "command" executed from a shell.

# Microservices: Interact with Python

Programmers can interact using an existing client in a program.



```
########################################################
# IMPORT PYAMI CLIENT MODULE AND API STATIC FUNCTION #
########################################################
import pyAMI.client

########################################################
# INSTANTIATE THE PYAMI CLIENT FOR ATLAS #
########################################################
client = pyAMI.client.Client(['atlas-replica-v2'])

########################################################
#  PRINT RESULT AS TEXT #
########################################################
res = client.execute('''SearchQuery -entity="dataset" -catalog="mc23_001:production"
-mql="SELECT * WHERE totalEvents=100000 AND dataType='EVNT' LIMIT 2 OFFSET 0"''')
print(res)
~
                                                              1,1          All
```

**Figure 9:** MQL query "command" executed with a Python script.

# Task Server: Run any kind of task

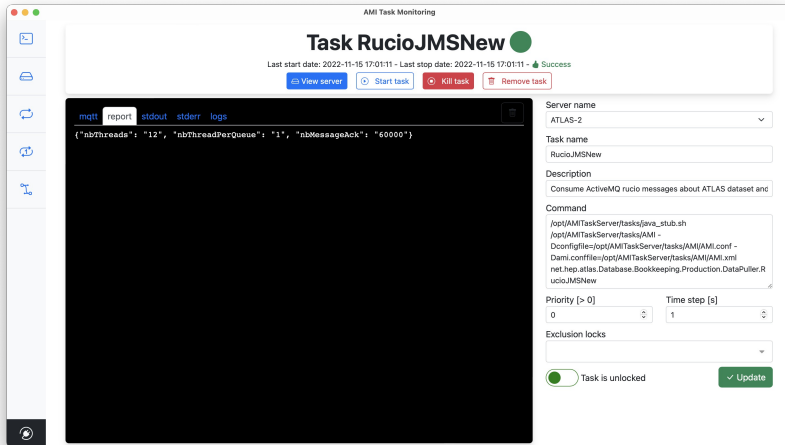The metadata manager can define tasks to aggregate metadata...



**Figure 10:** Configuration of recurrent tasks.

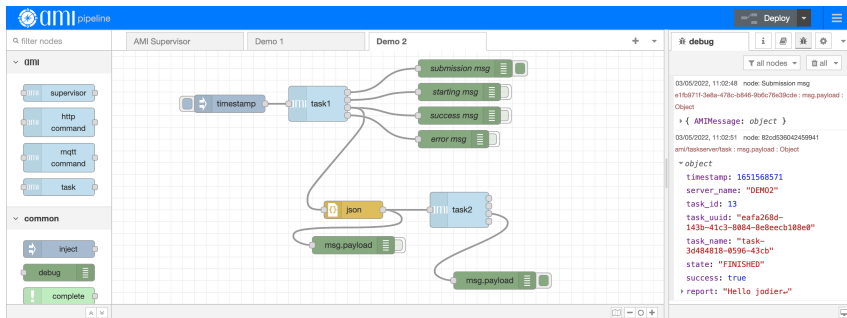# Task server: pipelined tasks

... and configure pipelined tasks.



**Figure 11:** Chained task executions.

# Try AMI

- Official site
  - https://ami-ecosystem.in2p3.fr/
- Docker-compose-based demo
  - Test on your laptop: https://github.com/ami-team/AMIDemo/
  - Test online: http://demo.ami-ecosystem.in2p3.fr:667/
- Documentation (Admin/end-user/developer guides, MQL langage)
  - https://ami-ecosystem.in2p3.fr/doc/
- Contact
  - ami@lpsc.in2p3.fr

# Thank You for Your Attention!