## LSDMA Topical Meeting 6 Data Intensive Computing and Applications

#### Lambda

# Integration of Data Repository System and Data Intensive Computing



<u>T. Jejkal</u>, S. Chandna, R. Dapp, V. Hartmann, A. Prabhune, F. Rindone, R. Stotzka, D. Tonne, A. Vondrous, X. Yang



## **Repository System**

#### Repository

Managed location/destination/directory/bucket where digital data objects are

- Registered
- Permanently stored
- Made accessible and retrievable
- Curated

#### Digital data object (dataset) consists of

- Data
- Description for re-use





#### **3** 01.06.2015

## What happens to Digital Data Objects?

#### Access and Retrieval

Query, download, sharing

#### **Preservation**

Bit- and content preservation, curation

#### **Analytics**

Data visualization and -mining, Online Analytical Processing (OLAP), requires "Extract, Transform, Load" (ETL)

> http://www.dwreview.com/Images/DataWarehouse\_Overview.gif http://nssdc.gsfc.nasa.gov/nssdc\_news/dec00/oais\_fig3.gif







## What is different for Research Data?

#### It starts with raw data

- Raw data is the most valuable data
- Processing is reproducible, DAQ not

#### Research data can be huge

- Typical datasets in range of GB/TB
- Local/on-the-fly processing not applicable
- Even management tasks get time consuming

#### Trying to understand research data can be challenging

- Complex processing chain of novel algorithms
- Many datasets (re-)processed with same algorithms
- Provenance information to tweak/compare algorithms and parameters

https://emergentmath.files.wordpress.com/2013/03/untitled2.png





## What is needed?

#### A research data repository system

- Support large scale research data
- Flexibility to be enhanced by DIC

#### **Description by metadata**

- High-level integration and flexibility
- Allows reproducibility and provenance information tracking
- Support for comparing different configurations

#### **Seamless integration**

- Computing should integrate into repository workflows
- No manual data movement
- Should be usable for system- and user-applications



## The Research Data Repository System





- Datasets registered in repository system identified by ObjectId
- Data available in repository storage accessible by repository system

## The Research Data Repository System





- Datasets registered in repository system identified by ObjectId
- Data available in repository storage accessible by repository system
- Configurable access points (AP) for data access/provisioning

## **Description by Metadata**





 Description of task including arguments, version and binary package

## **Description by Metadata**





- Description of task including arguments, version and binary package
- Basic description of execution environment, e.g. repository AccessPoint, custom properties and handler implementation

## **Description by Metadata**





- Description of task including arguments, version and binary package
- Basic description of execution environment, e.g. usable AccessPoint, custom properties and handler implementation
- Capabilities to allow basic decision on taskenvironment assignment

#### → Application testing still necessary!

## **Seamless Integration**





- Task instance linked to one or more OIDs
- Task execution component of repository system takes care of capability matching, data staging, submission, monitoring and result registration



## **Seamless Integration**





- Task instance linked to one or more OIDs
- Task execution component of repository system takes care of capability matching, data staging, submission, monitoring and result registration

Execution **Environment** 



## Conclusions



- First integration of DIC capabilities into repository system
- Metadata description of tasks including arguments, execution environment and capabilities
- Flexible execution environment handler interface
- Used applications must be well tested (focus on recurring tasks)
- Currently first evaluation, afterwards integration of Apache Flink
- Support for provenance standards OPM and PROV

