

Workflows in the VO

Grid and Web Services Session

Jose Enrique Ruiz
IAA-CSIC

October 19th 2011
2011 IVOA Fall Interop Meeting - Pune

WMF

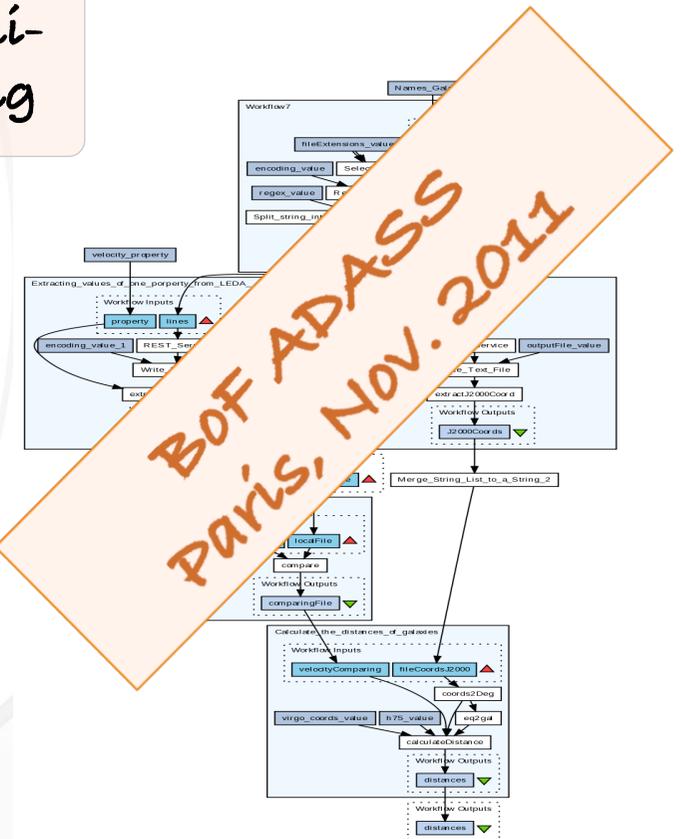
What are workflows?

Combination of **data** and **processes** into a configurable and structured set of steps that implement semi-automated computational solutions in problem solving

Types of workflows in Astronomy

- Personal script-based recipes
Python, IDL, Software..
- Multi-archive **VO** recipes
- Internal group developments
GRID, Clusters..
- Processing pipelines
Provide Data, Computing Infrastructure, Tools..

Scientifically exploitable results vs. **scientific insight**
Easily **accessible** and **reproducible** (Shared)



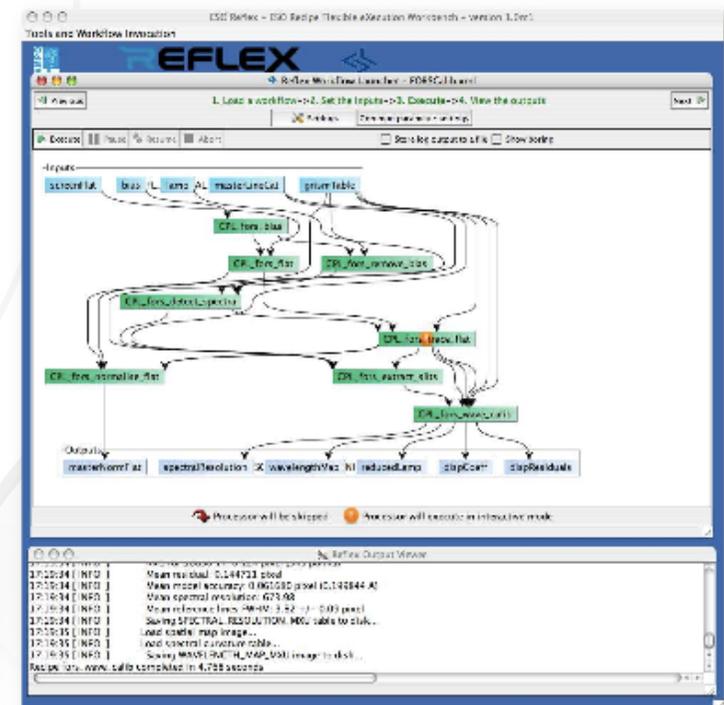
ESO Reflex

Finland's in-kind contribution to ESO

- Prototype/feasibility study
- Initially based on Taverna 1

Current implementation based on Kepler

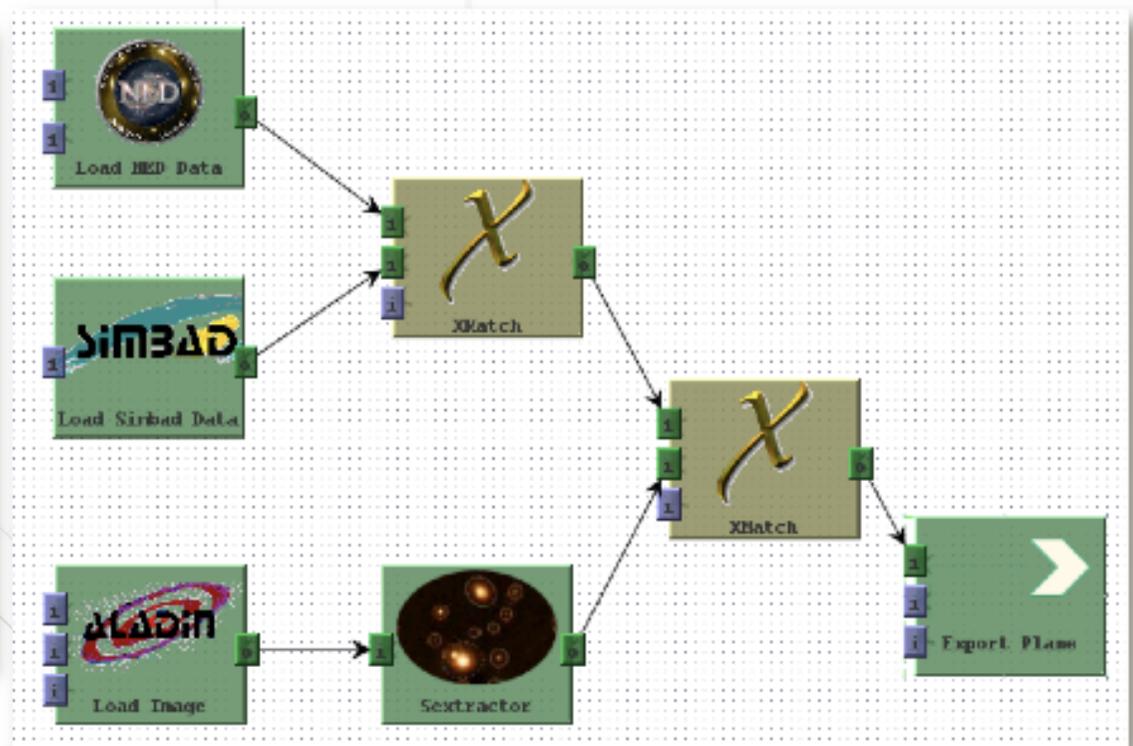
- Main intent: replace CLI (pipeline orchestrator)



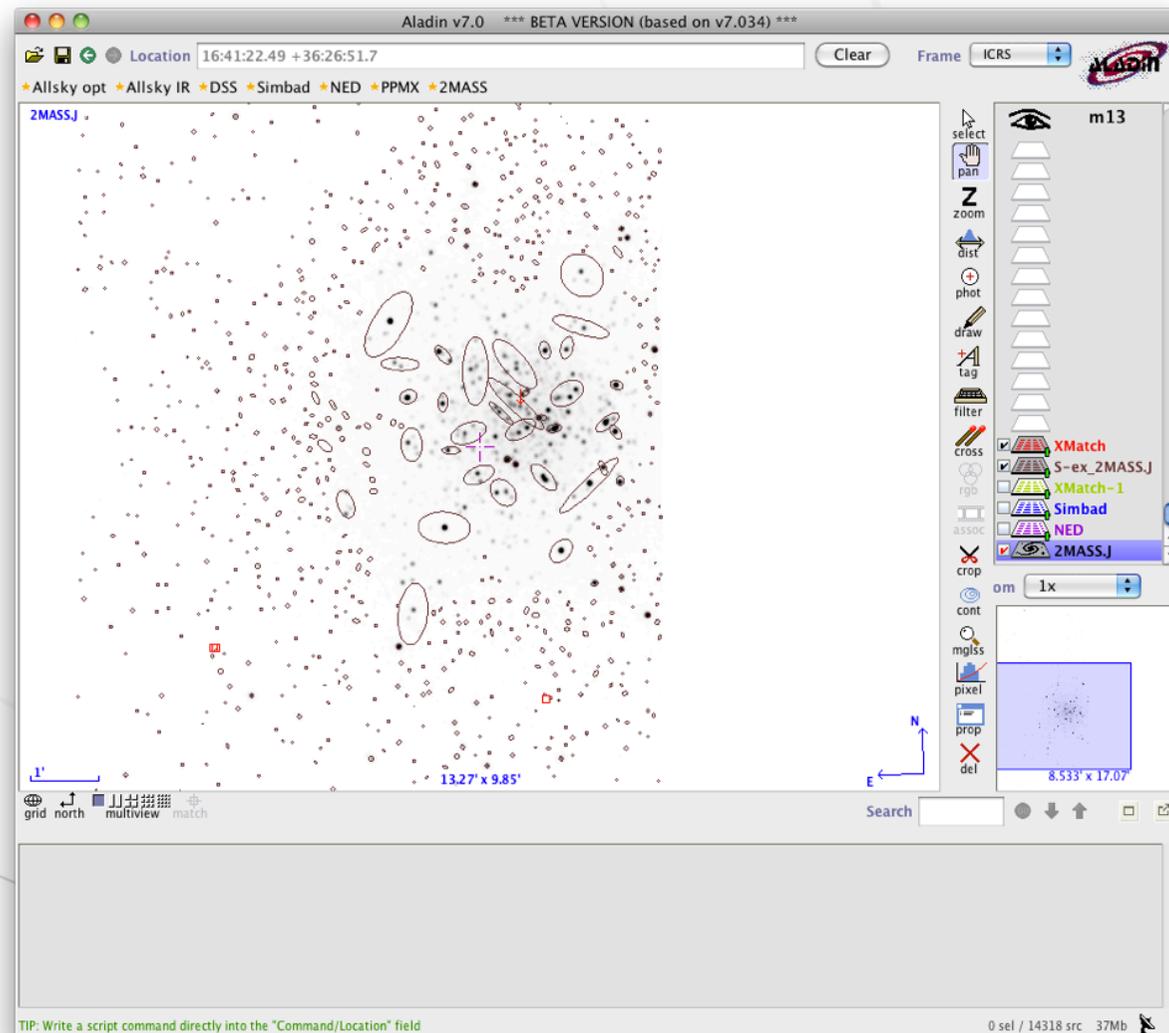
Aladin JLOW Plugin

Aladin plugin API permits graphical replacement of Aladin tools

```
#AJS
get NED m13 14'
get Simbad m13 14'
sync
xmatch NED Simbad 4
sync
get Aladin(2MASS,J) m13
sync
get SExtractor(2MASS.J)
sync
xmatch "XMatch*1" "SExtractor*"
sync
export -votable "XMatch*2"
```



Aladin JLOW Plugin



Wf4Ever

2011 – 2013 EU funded FP7 STREP Project
Preservation of Workflows

- Preservation of workflows and associated material
- Archival, classification, indexing in semantic repositories
- Provide advanced access and recommendation capabilities
- Collaborative working platform
- Sharing, re-use, re-purpose
- Digital Libraries

Jose Enrique Ruiz
Tuesday DCP Session

Cyber-SKA

Provide infrastructure that will be required to address the needs of future radio telescopes such as the Square Kilometre Array

web based workflow builder

- Image segmentation
- Image mosaicking (Montage)
- Spatial reprojection
- Plane extraction from data cubes

The screenshot displays a web-based workflow builder interface. At the top, there are three buttons: 'Create Pipeline', 'Execute Pipelines', and 'Clear All Pipelines'. Below these are five tabs: 'Segment', 'Mosaic', 'Plane Extract', 'Compress', and 'Stage'. The interface shows two pipeline configurations, Pipeline Number: 0 and Pipeline Number: 1. Pipeline 0 consists of four stages: 'file list' (with an 'Add files...' button), 'segment' (with parameters: bbox swl: 60d, bbox swr: -0.5d, bbox nwl: 160d, bbox ner: 0d, and a 'remove' button), 'mosaic' (with parameter: Background correction: 0, and a 'remove' button), and 'stage' (with parameter: Directory prefix: fourth, and a 'remove' button). Pipeline 1 consists of three stages: 'file list' (with an 'Add files...' button), 'planeextract' (with parameters: Plane start: 0, Plane end: 0, and a 'remove' button), and 'stage' (with parameter: Directory prefix: [empty], and a 'remove' button). Blue arrows indicate the flow between stages in each pipeline.

Montage

- FITS Image Mosaicking
- Toolkit for Desktops, Clusters and Grids

Astro-WISE

- Distributed data storage and computing infrastructure
- Track process provenance of final data products
- Calibration and analysis of images

Helio-VO

- Solar physics virtual Observatory
- Enable workflow execution via Taverna Server

Workflows VO France

- Provide use cases mainly oriented VO
- AIDA Workflow System implements FITS validation with ChardM

The next generation of archives

Much wider FOV and spectral coverage

- Huge sized datasets (~100 TB)
- Big Data science highly dependent on I/O data rates
- Subproducts as **virtual data** generated on-the-fly

Automated surveys

- Huge amount of tabular data
- Services for **Knowledge Discovery in Databases**

We are moving into a world where

- computing and storage are cheap
- data movement is death

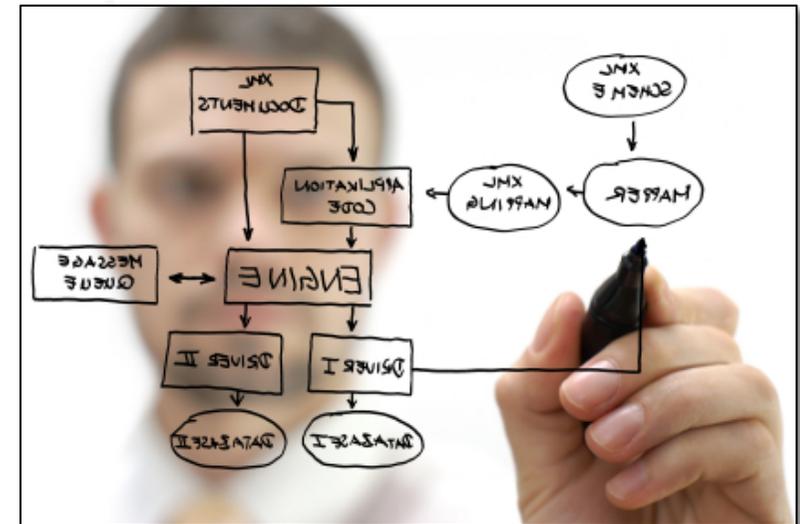
Archives should evolve from data providers into **virtual data** and **services providers**, where web services may help to solve bandwidth issues.

Archives speaking web services

- Smaller virtual data subproducts
- Distributed, multi-archive, multi-wavelength astronomy

web services based workflows as a disruptive working methodology

- Reproducible
- Repeatable results
- Encourage best practices
- Modular nature allows
 - Re-use
 - Re-purpose
- Expose
 - Provenance
 - Scientific method
- Formative
- Foster collaborative work



Distributed data analysis in the VO

- Panchromatic, multi-archive, multi-facility
- Executes in the VO Infrastructure
- Orchestration of simple services

Present processing pipelines

- Produce exploitable data
- Provenance modelling
- VO compliant data

Data processing from the VO

- Provide custom re-processing to VO users
- Virtual data generation through UWS in VOspace

Workflows VO Characterization

- Inputs
- Outputs
- Processes
- Descriptions
- Metadata
- Etc..

IVOA Working Groups

- **Data Modeling**
Characterization, Provenance..
- **Semantics**
Ontologies, Vocabularies, Annotations
- **Data Access Layer**
TAP, self-descriptive Protocols
- **Grid and Web Services**
UWS, VOSpace, SSC
- **Applications**
SAMP
- **IG . KDD**
Knowledge Discovery and Mining
- **IG . Data Curation and Preservation**
Persistent Identification and Curation of VO Resources..
Wf4Ever Project

*Carlo Maria Zwölf
Thursday GWS Session*

*Fabio Pasian
Liason with Open Grid Forum*



IVOA Note

Scientific Workflows in the VO
workflow@ivoa.net