# Workflows
# Access and Massage VO Data

**José Enrique Ruiz**

on behalf of the Wf4Ever Team

# Wf4Ever
# Advanced Workflow Preservation Technologies for Enhanced Science
## 2011 - 2013

1. Intelligent Software Components (ISOCO, Spain)
2. University of Manchester (UNIMAN, UK)
3. Universidad Politécnica de Madrid (UPM, Spain)
4. Poznan Supercomputing and Networking Centre (Po
5. University of Oxford and OeRC (OXF, UK)

6. Instituto Astrofísica Andalucía (IAA-CSI
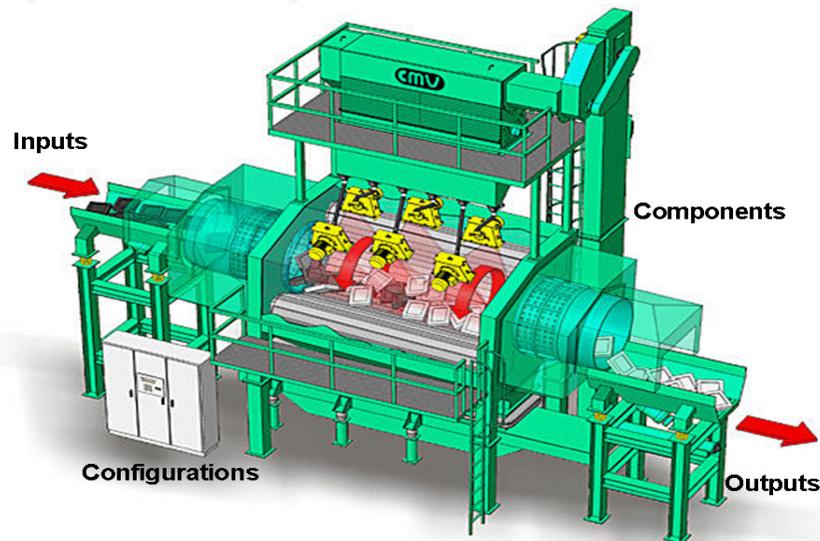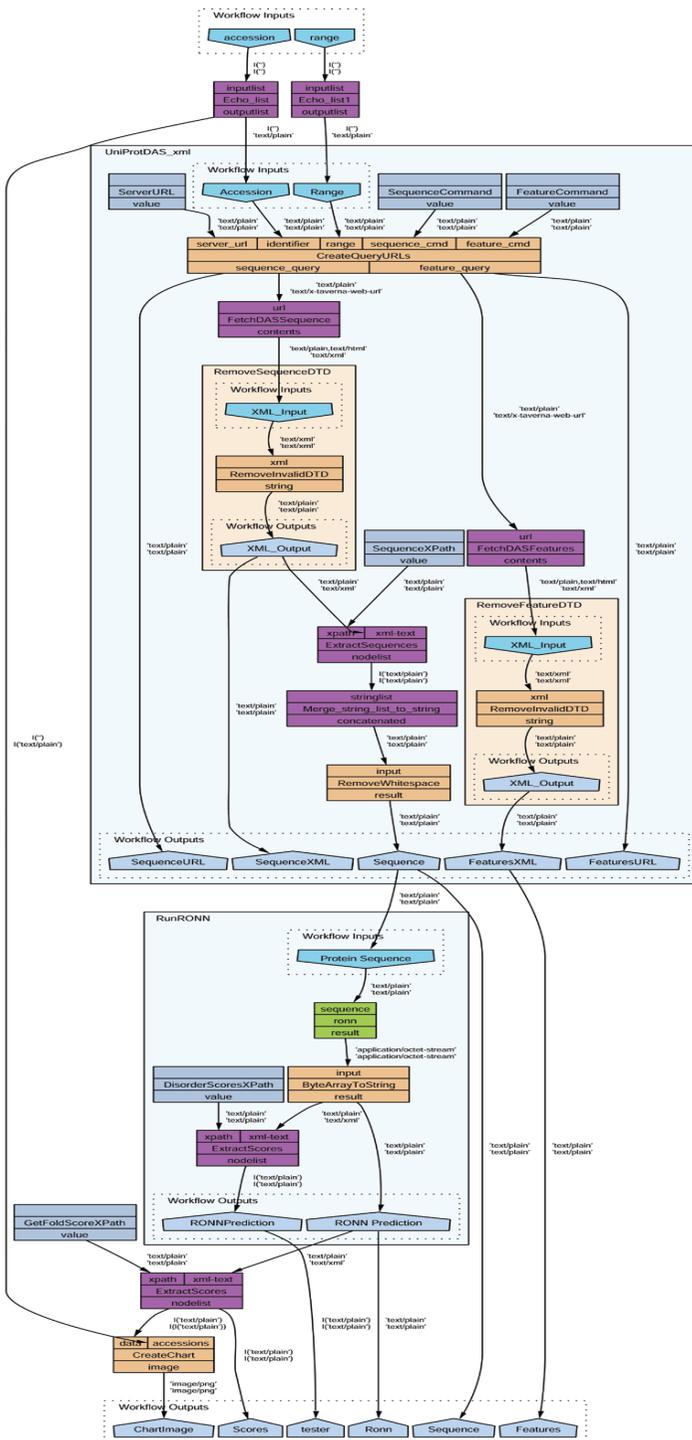7. Leiden University Medical Centre (L

Reproducible Science

# What is a Scientific Workflow?



» A mechanism for coordinating the execution of services and codes, and linking together resources.

» The combination of data and processes into a configurable, modular, structured set of steps that implement semi-automated computational solutions in scientific problem-solving.

» The implementation of a scientific method.



Inputs

Components

Configurations

Outputs

» **IVOA Note Definition**

*These are networks of analytical steps that may involve, e.g., database access and querying steps, data analysis and mining steps, and many other steps including* <mark>*computationally intensive jobs on high performance cluster computers.*</mark>

» **Wf Software**

› Taverna

› Kepler

› Pegasus

› Triana

› ESO Reflex

**Related Initiatives**

› ER-Flow

› VAMDC

› Helio-VO

› Cyber-SKA

› IceCore

› Montage

› Astro-WISE

› AstroGrid

**In the VO**

› GWS WG

› VO France WF WG

› VAMDC

› AstroGrid

Capturing Actions
Reproducibility

## AstroTaverna Workflows

Retrieving and Manipulating VO Data  **+ Catalog        TML Pages**

- ConeSearch
- SIA
- SSA
- TAP coming soon…                                                    **Web Services**

- Tabular Data (VOTables)                                    **ccess to JDBC databases**
- Images, but not yet Spe

- Crossmatching, Fi        solving, Coordinates and reference system transfor        ssage.. (**STILTS**)
- Overplotting so        s on Images and filtering, overplot circles, ellipses, etc. as a        of physical magnitude. Resampling, crops, blinks, mosaics, mov    blinks, RGBs, fusion, diff.. (**ALADIN**)
- **SAMP** for final inspection

**+ Advanced Analysis using Scripts**

*No interactive actions and decisions based on visual inspection*

## VOData Access: VO Services Discovery



🌐 http://amiga.iaa.es/p/290-astrotaverna.htm

## VOData Massage: VOTables, STILTS, Aladin, TerminalSim



http://amiga.iaa.es/p/290-astrotaverna.htm

## Massage of Tabular Data

**Calculation of Luminosity Profiles for a Sample of Galaxies extracted from SDSS DR8**



**GALFIT**

**IRAF**
*Image Reduction and Analysis Facility*

**90 galaxies observed in 3 bands**

SExtractor

## Aladin Scripts and Macro executing in GUI/noGUI mode

## VO compliant data from pipelines

Traditional data processing pipelines, e.g., instrumental or survey data processing pipelines, which produce higher, level data products. At present there are many variants of these and they have little or no direct connection to VO, aside from possibly producing VO-compliant data or being optionally driven from VO.

It is not clear how much VO mechanisms are needed at this level (VO compliant data and metadata, modelling provenance, etc.)

## Driving Data Processing Pipelines from the VO

In this case we have a traditional data processing pipeline and the remote user or client software invokes a job to do some pipeline reprocessing, e.g., to custom reprocess an instrumental dataset to produce a new image, cube, etc. The "workflow" in this case runs at a single site, and VO is used to drive the job remotely (SSO, UWS) and manage the results (VOSpace, VO data services).

We could think on integrating the traditional data processing pipelines we already have with VO, to allow VO users to do on-the-fly reprocessing to generate data products which can be analysed with VO (custom reprocessing of observatory data for example)

Some attempts to integrate general processing applications have been made with CEA and UWS.

## Distributed Data Analysis Workflows

In this case a user or a client defines and executes a distributed workflow, which invokes services on multiple remote sites via the VO infrastructure. The workflow would be entirely in VO-space, driving simpler services at the individual sites.

The AstroTaverna developments provide a graphical tool for the composition and design of workflows based on VO services and data from different archives and facilities.

**Self Descriptive Web Services**: S3, SimDAL, PDL, DataLink

## Much wider FoV and spectral coverage

» Large volumes for an observed datacube

» Subproducts are **Virtual Data** generated on-the-fly

| | Low Res | | High Res | | Extreme Res | |
|---|---|---|---|---|---|---|
| Number | 4 Bytes | 4B | 4 Bytes | 4B | 4 Bytes | 4B |
| Resolution | 2,048 x 2,048 | 16MB | 8,192 x 8,192 | 268MB | 12,288 x 12,288 | 603MB |
| Channels | 16,384 | 0.27TB | 16,384 | 4.39TB | 16,384 | 9.8TB |
| Stokes & Weighting | 1 | 0.27TB | 1 | 4.39TB | 4 + 1 | 49.5TB |

ASKAP Cubes
Prof.  Kevin Vinsen

## Automated surveys

» Huge amounts of tabular data

» Services for KDD



Extraction of scientifically relevant information from a multidimensional parameter space

» Exploration services

» Anomaly detection

» Cross-matching data

» Dimensionality reduction

» A cloud of Web Services

> Archives should evolve from data providers into
>
> » Virtual data  providers
>
> » Software tasks providers

» Archives speaking Web Services

> Astronomy of multi archives/facilities/wavelength
>
> Interconnected and interoperable archives
>
> » Data -> Virtual Observatory
>
> » Software Tasks

Preservation

**Process should benefit of the same privileges acquired by data**

Preserving the method ensures replication of final results at any moment