

# Heterogeneous data fusion and data mining in astronomy

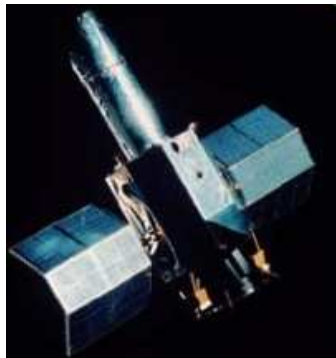


Françoise Genova, Aspera Workshop, 08/10/2010



# Sharing astronomical data

- The astronomy model: competitive AOs to get observation time on facilities
- Re-using data for scientific objectives different from the original ones, i.e. optimize the science return of large ground- and space-based instruments and of large surveys



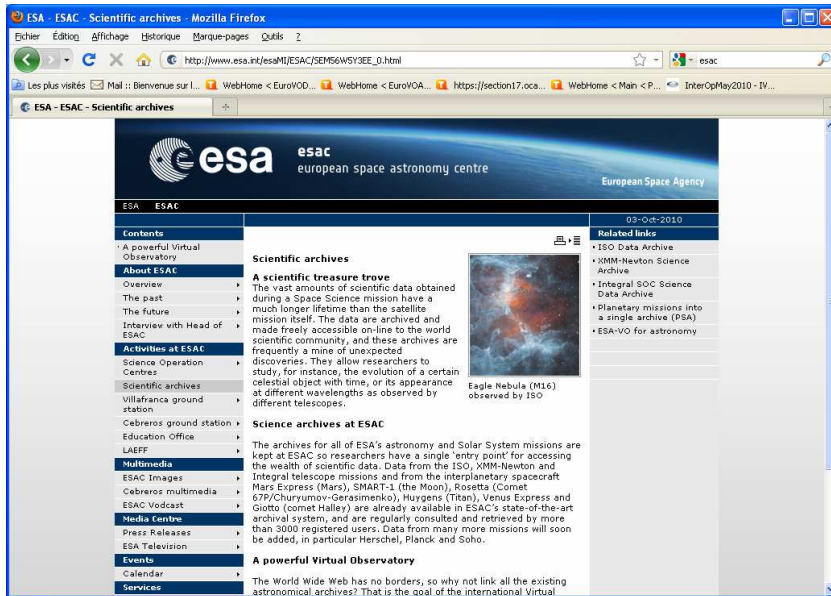
**IUE (1978-1996):** five times more publications from data retrieved in the archive than from the selected observing teams (Wamsteker, Griffin, 1995) – **a major precursor**

# Access to IUE data



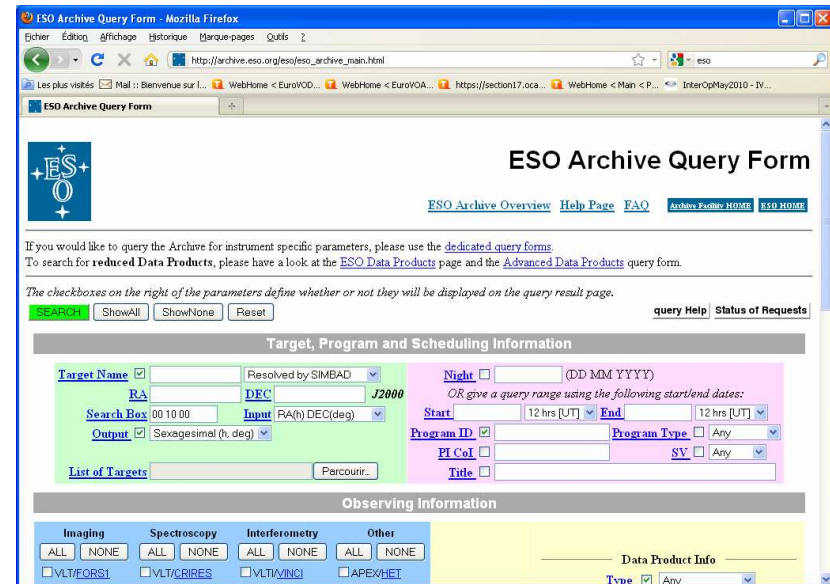
Françoise Genova, Aspera Workshop, 08/10/2010





# European Agency data archives

And MANY others



Françoise Genova, Aspera Workshop, 08/10/2010



## DCA/AIDA Census of European data Centres

- Inclusive definition
- Keywords: provide a service to the community, added-value, sustainability, quality, national and/or international role
- Characterizes the population of European data centres (70 answers from 16 different countries, > 230 forms filled describing a resource)
- Lively snapshot of a very diverse landscape
  - Covers all disciplines of astronomy
  - Large variety of approaches and sizes, from large data centres maintained by European or National Agencies to small teams maintaining a specific service



The screenshot shows a Mozilla Firefox browser window with the following details:

- Browser Title:** DataCentresinVO < EuroVODCA < TWiki - Mozilla Firefox
- Address Bar:** http://cds.u-strasbg.fr/twikiDCA/bin/view/EuroVODCA/DataCentresinVO
- Page Content:**
  - Header:** DATA CENTRE ALLIANCE
  - Navigation:** EuroVODCA, TWiki > EuroVODCA Web > DataCentresinVO (03 May 2007, ChristopheArviset), Edit, Attach
  - Section:** Data Centres in the Virtual Observatory context
  - Text:**

Data Centres are an essential component of the Virtual Observatory, publishing data, metadata and services, and providing the physical storage and computational fabrics. The VO development is a strong incentive to share data and knowledge, and many teams are willing to provide data and services in their domains of expertise. 'Classical' data centres, such as ground-and space-based observatory archives, and generalist data centres, are key providers of added-value services and tools. More and more teams are willing to join with value-added services and tools in specific domains, and VO 'data centres' work in very different contexts - national or international Agencies, scientific laboratories - , and are highly diverse in size and objectives, from small and specific to large and general. Common keywords are the willingness to *provide a service to the community* , provision of *added-value* built on expertise, some kind of *sustainability* , and concern for *quality* . Lessons learnt from the long term history of astronomical data centres show that when beginning these activities, critical parameters are in particular *having a critical mass* adapted to the goals, and ensuring *medium-term sustainability* , which requires at least a strong support from the local authorities. An important factor to win community support, which is indispensable to secure funding, is to find a *national and/or international niche* .

Many types of contribution are possible: data archives, with a particular emphasis put on 'science ready' data; added-value data bases, services; tools, software suites and algorithms, for instance for data visualisation, data analysis and data mining; thematic services to help solving a well-defined science question; full data analysis or reasearch environments. New types of services are emerging, with in particular theoretical services, providing modelling results, or matching models with observations.
- Left Sidebar:**
  - EuroVODCA**
  - Log In or Register
  - EuroVODCA Web**
    - Create New Topic
    - Index
    - Search
    - Changes
    - Notifications
    - Statistics
    - Preferences
    - Help
  - Webs**
    - EuroVODCA
    - Main
    - Sandbox
    - TWiki
  - Links**
    - Euro-VO
    - Euro-VO AIDA
    - IVOA
    - VOtech
    - Communication Network
- Footer:** http://cds.u-strasbg.fr/twikiDCA/bin/view/EuroVODCA/WebHome



Françoise Genova, Aspera Workshop, 08/10/2010



# A multipolar world

- Many facilities and authorities
- Large diversity
- Common VObs framework

# Registry of Resources

- Key element: the ‘yellow pages’
- Compliant with OAI-PMH (allows interoperability with digital libraries)
- Dublin core + disciplinary extensions
- Not a unique registry, but
  - Several harvestable registries (+ publishing registries)
  - + A Registry of Registry



# Euro-VO Registry : Main Page

**EURO VO AIDA** Astronomical Infrastructure for Data Access

The Euro-VO projects: VOTECH EuroVO-DCA EuroVO-AIDA

EURO-VO Registry

Search Resources

Insert Resources

Update Resources

Member of

Powered by

## EURO-VO Full Harvestable VO Resource Registry

Welcome to the EURO-VO Full Harvestable VO Resource Registry.

If you want to know more details about what a "VO Registry" is, you can have a look at the [IVOA Resource Registry](#) specifications.

At the left panel, you can find different utilities to handle Registry data:

- **Search Resources**  
Allows to search among the different resource types in the Registry. For example, clicking on the "Simple Image Access" will display all the available Resources of type "Simple Image Access" in the registries around the world.  
Due to the fact that the Tabular Sky Service Resource type allows for one entry per table, and that CDS contains thousands of tables, and in order to not clobber the access, we have separated the Tabular Sky Service in CDS and non-CDS searches for commodity.
- **Insert Resources**  
Allows the insertion of a new Resource in the EURO-VO Registry.
- **Update Resources**  
Allows the edition of a Resource that resides in the EURO-VO Registry. Resources can only be updated in the Registry where they have been introduced, and not in registry that harvest them from other places. Consult the registry specification for more details (see above URL).
- **Registry Quick Search**  
Allows for a quick search on a string. The string introduced is searched in the following Resource Registry fields:
  - Title, ShortName, Identifier, ResourceType
  - Content -> Description, Subject, Type, ContentLevel

If you have any question regarding the EURO-VO Registry, please send a note to our [Registry manager](#).

last updated: 16-Dec-2008

co-funded project

Demo

<http://registry.euro-vo.org/>

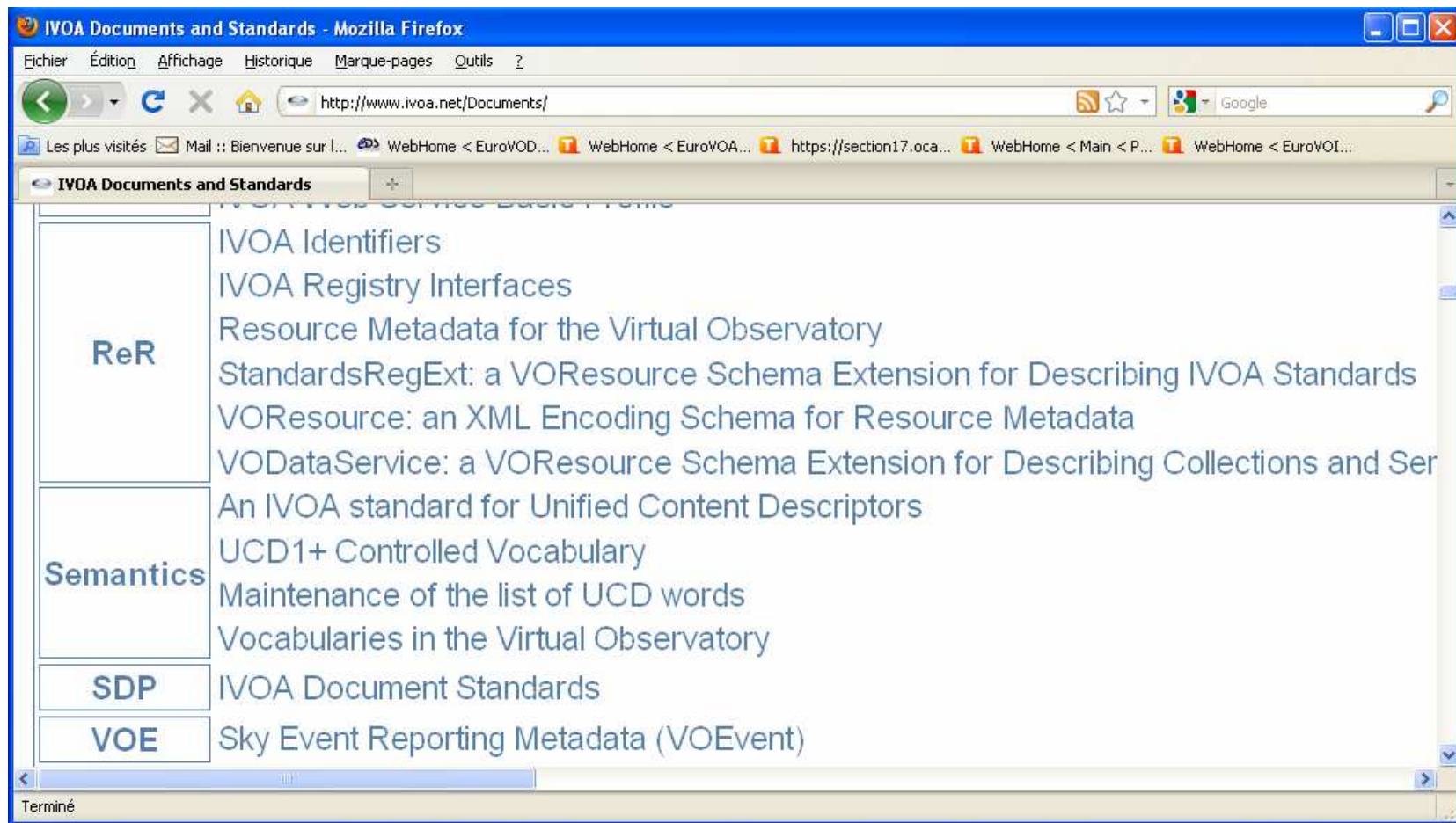


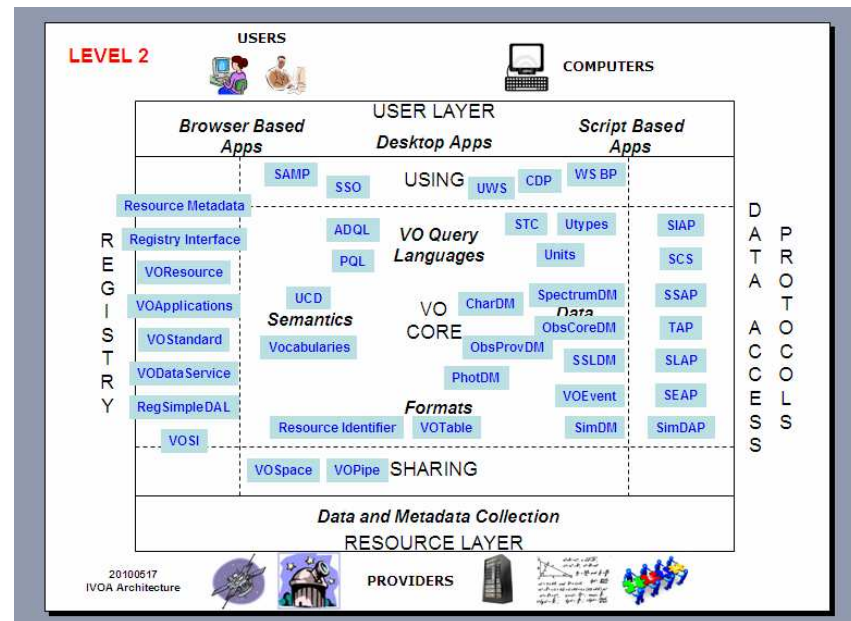
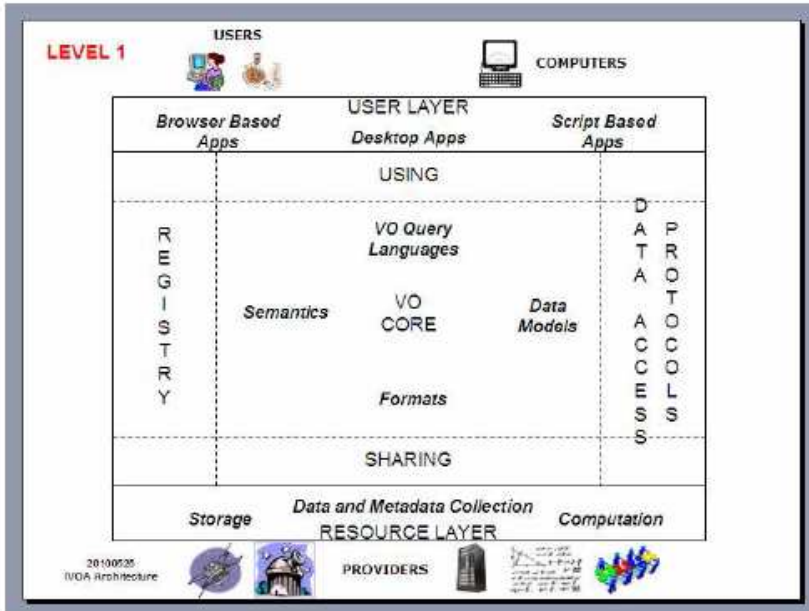
Françoise Genova, Aspera Workshop, 08/10/2010



# Current issues on registry

- Validation of entries (rules, tools)  
    which implies to validate the validators!!
- International effort
- Add extensions to describes the different types of services





Identity metadata

Title Sloan Digital Sky Survey  
ShortName SDSS  
Identifier ivo://stsci.edu/mast/sdss

Curation metadata

Publisher Space Telescope Science Institute/MAST  
PublisherID ivo://stsci.edu/mast  
Creator Sloan Digital Sky Survey Consortium  
Creator.Logo [http://archive.stsci.edu/images/sdss\\_logo.gif](http://archive.stsci.edu/images/sdss_logo.gif)  
Contributor Sloan Digital Sky Survey Consortium  
Date 2003-02-01  
Version SDSS EDR  
Contact.Name Archive Branch, Space Telescope Science Institute  
Contact.Address 3700 San Martin Drive, Baltimore, MD 21218 USA  
Contact.Email [archive@stsci.edu](mailto:archive@stsci.edu)  
Contact.Telephone +1-410-338-4547

General content metadata

Subject galaxies, quasars, stars, CCD photometry,  
spectroscopy, redshift, sky surveys  
Description The Sloan Digital Sky Survey is using a dedicated  
2.5 m telescope and a large format CCD camera to obtain im-  
ages of over 10,000 square degrees of high Galactic latitude sky  
in five broad bands (u', g', r', i' and z', centered at 3540, 4770,  
6230, 7630, and 9130 Å, respectively). Medium resolution spec-



Françoise Genova, Aspera Workshop, 08/10/2010



Source [SDSS Survey 10 months of obs.](#)  
 2002AJ....123..485S  
 ReferenceURL <http://archive.stsci.edu/sdss/index.html>  
 Type Survey, Catalog, EPOResource  
 ContentLevel Research  
 Relationship mirror-of  
 RelationshipID <ivo://sdss.org/sdss/edr>

Collection and service content metadata

Facility Apache Point Observatory, Sloan 2.5-m Telescope  
 Instrument Five-band clocked CCD camera

Coverage.Spatial PositionInterval FK5 145.17 -1.25 235.9 1.25 PositionInterval  
 FK5 250.71 52.15 267.0 66.29 PositionInterval FK5 350.43  
 -1.25 359.99 1.17 PositionInterval 0.0 -1.25 56.37 1.17

Coverage.RegionOfRegard 0.0001  
 Coverage.Spectral Optical  
 Coverage.Spectral.Bandpass u', g', r', i', z'  
 Coverage.Spectral.MinimumWavelength 400.e-9  
 Coverage.Spectral.MaximumWavelength 850.e-9  
 Coverage.Temporal.StartTime 1999-12-25  
 Coverage.Temporal.StopTime 2001-07-15  
 Coverage.Depth 3 e 6



Françoise Genova, Aspera Workshop, 08/10/2010





#### Data quality metadata

DataQuality	A
ResourceValidationLevel	4 [provided by registry curator]
ResourceValidatedBy	ivo:/us-vo.org/registry
Uncertainty.Photometric	3.e-7
Uncertainty.Spatial	0.00003
Uncertainty.Spectral	1.e-11
Uncertainty.Temporal	0.1

#### Service metadata

Service.AccessURL	<a href="http://archive.stsci.edu/cgi-bin/sdss/catalog">http://archive.stsci.edu/cgi-bin/sdss/catalog</a>
Service.InterfaceURL	<a href="http://archive.stsci.edu/sdss/catalog.html">http://archive.stsci.edu/sdss/catalog.html</a>
Service.BaseURL	<a href="http://archive.stsci.edu/cgi-bin/sdss/catalog">http://archive.stsci.edu/cgi-bin/sdss/catalog</a>
Service.HTTPResultsMIMEType	text/xml
Service.StandardID	ivo://ivoa.net/Services/ConeSearch
Service.MaxSearchRadius	0.2
Service.MaxReturnRecords	5000
Service.MaxReturnSize	5.e8

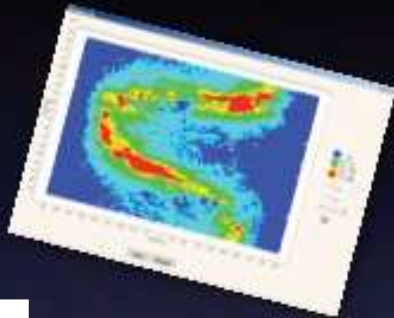
# SAMP

- The Simple Application Messaging Protocol
- Enables software tools to interoperate and communicate

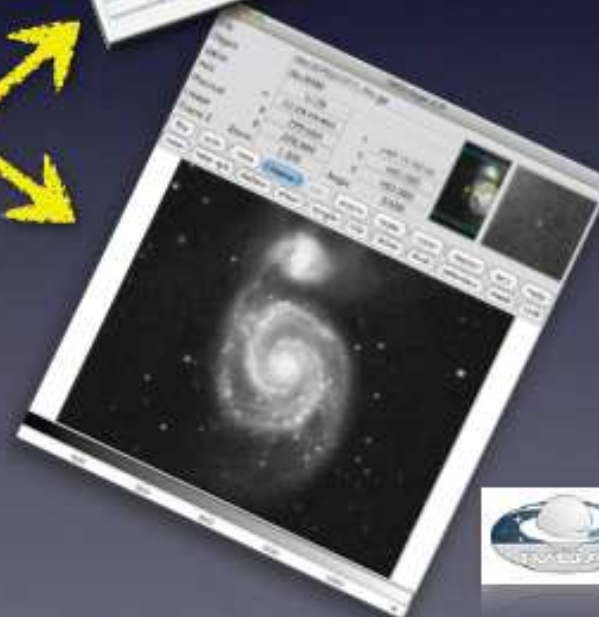
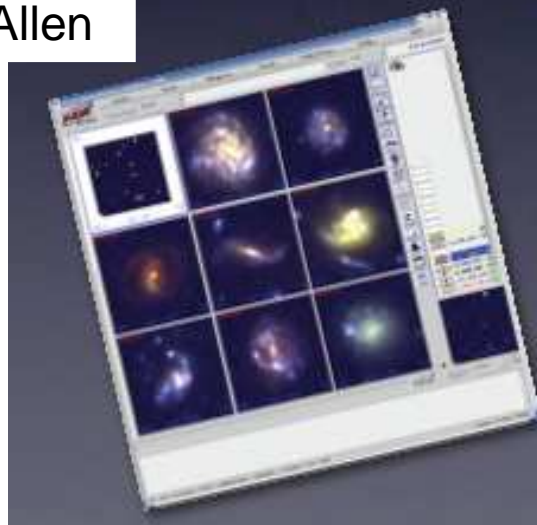
*IVOA members have recognised that building a monolithic tool that attempts to fulfil all the requirements of all users is impractical, and it is a better use of our limited resources to enable individual tools to work together better. One element of this is defining common file formats for the exchange of data between different applications. Another important component is a messaging system that enables the applications to share data and take advantage of each other's functionality.*

# Interoperability of VO Tools

- interaction between tools

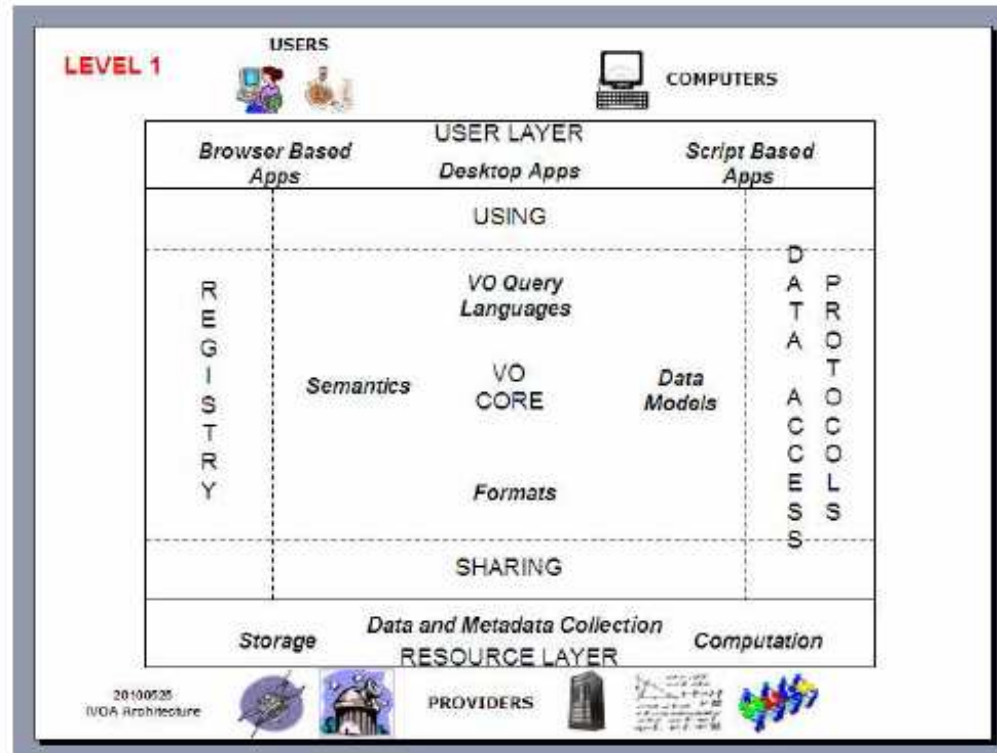


NAME	RA	DEC	TYPE	CLASS	CLASS	CLASS	CLASS	CLASS	CLASS
NGC 147	00:52:02.1	41:53:04.3	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy
NGC 148	00:52:02.1	41:53:04.3	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy
NGC 149	00:52:02.1	41:53:04.3	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy
NGC 150	00:52:02.1	41:53:04.3	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy
NGC 151	00:52:02.1	41:53:04.3	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy	Galaxy



Courtesy of  
M.G. Allen

# Semantics



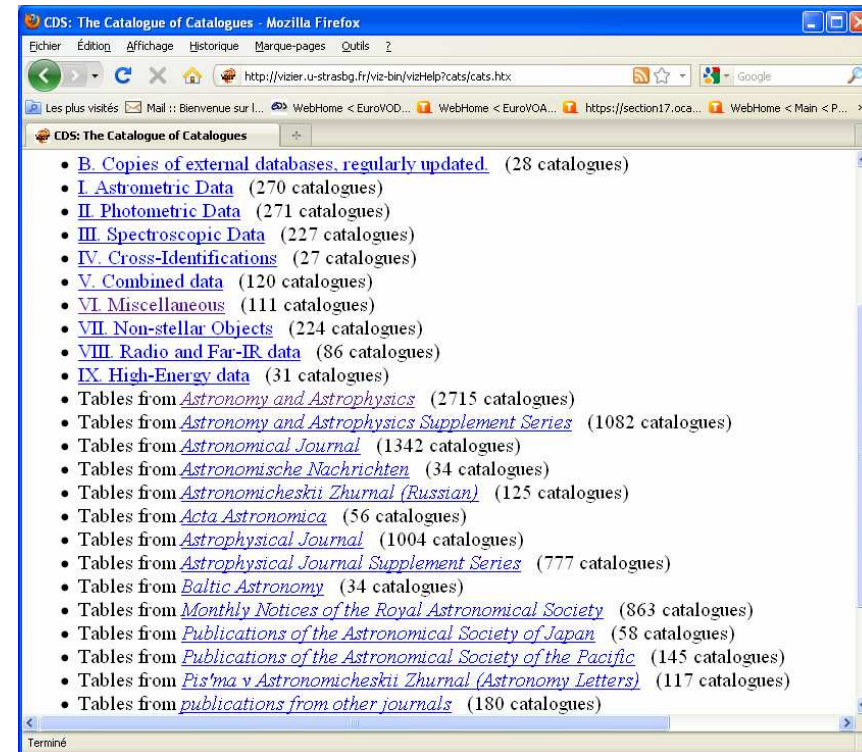
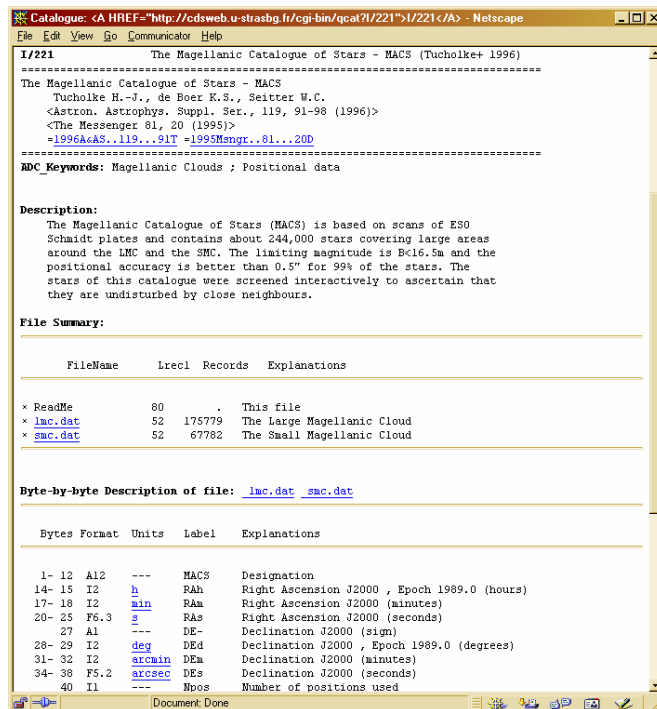
# Catalogues and published tables

## CDS VizieR service

A single standard

# Catalogues

A single view of the tables published in journals;  
collab. journals + data centres  
Lists of observations



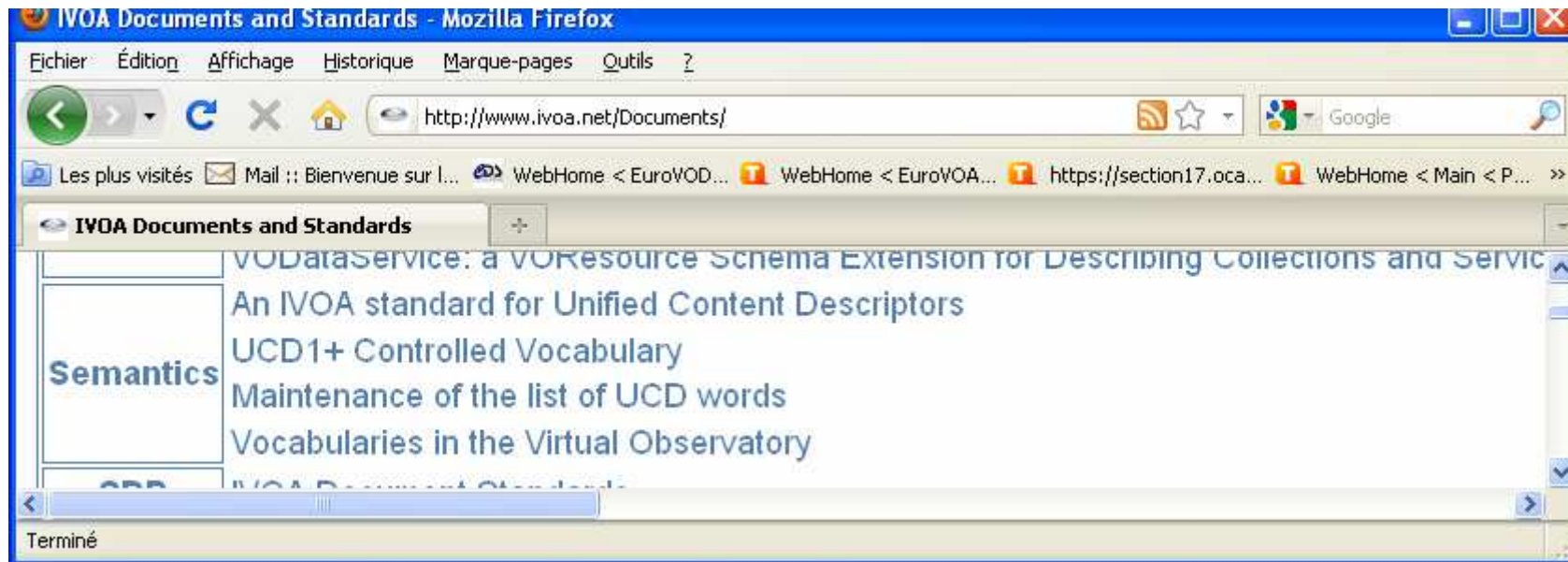
Françoise Genova, Aspera Workshop, 08/10/2010



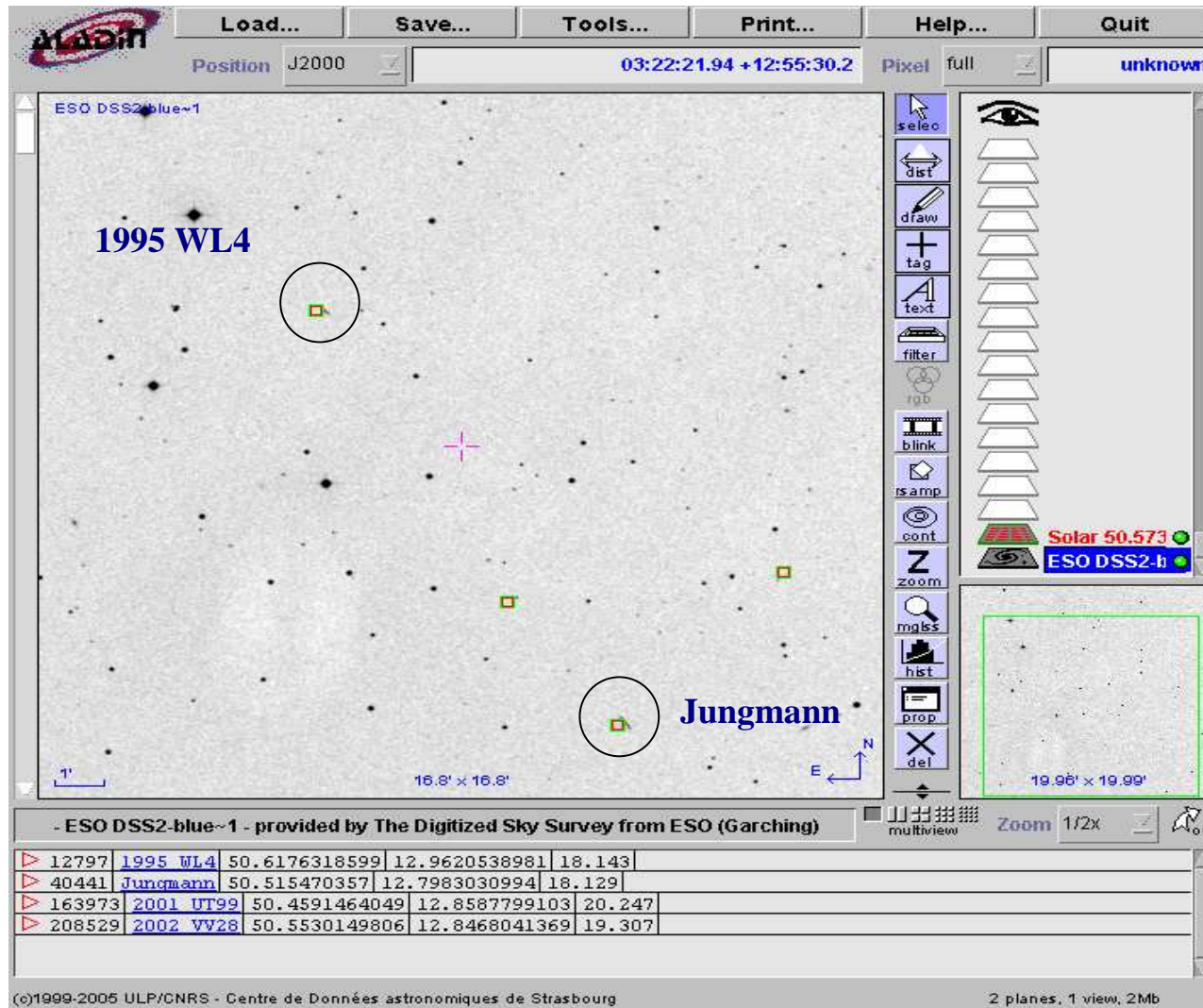
# The UCDs

- Describes quantities in astronomy
  - Starting point: the CDS catalogue service VizieR
  - 100.000 columns in VizieR
- > IVOA standard Unified Content Descriptor





- Maintenance of UCDs: a procedure to add new ones
- Standard for vocabularies: W3C (SKOS, RDF)
- Notes on ontologies (ontology of object types, use case)



SkyBot  
 (IMCCE,  
 Paris  
 Observatory)  
 +  
 Aladin (CDS,  
 Strasbourg  
 Observatory)  
 +  
 VO standard  
 (VOTable)

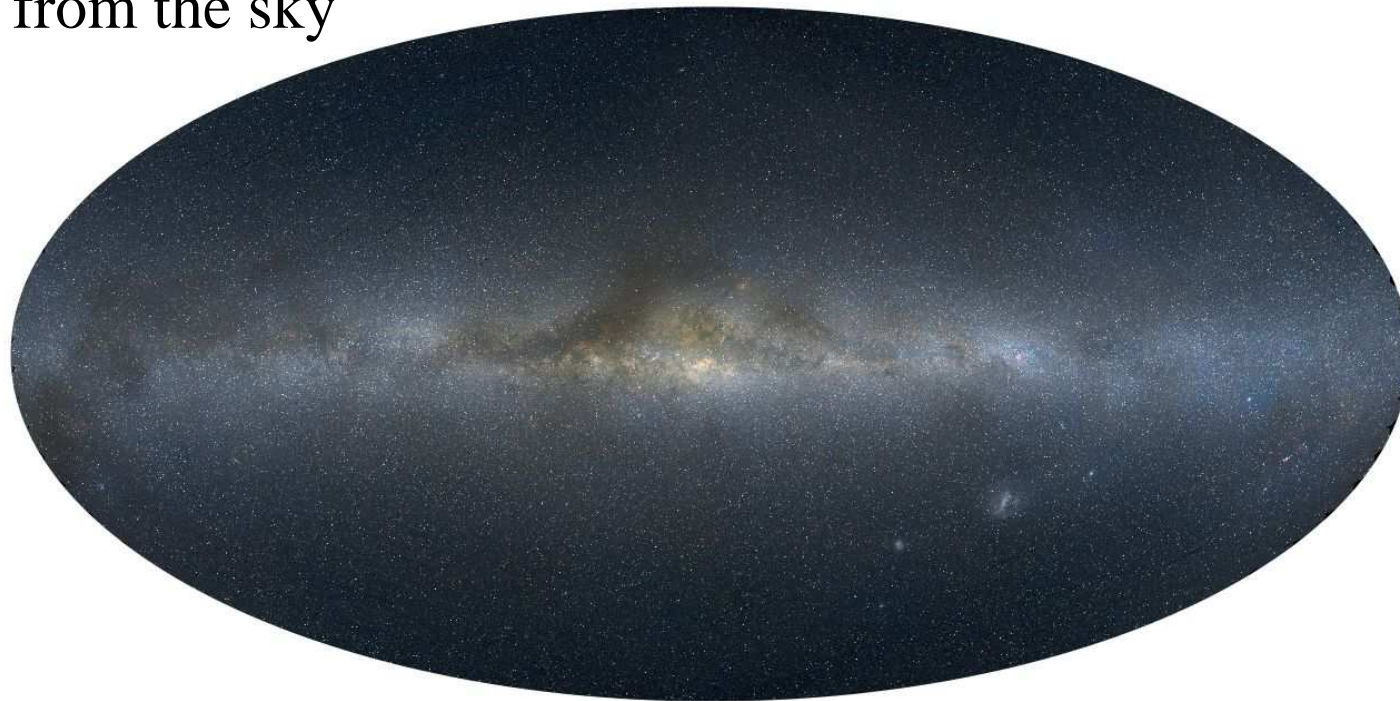
*Fast solar system  
 moving targets*

- Also, e.g., VOEvent – LSST! (Ray's talk)

# Long term data accumulation and re-use

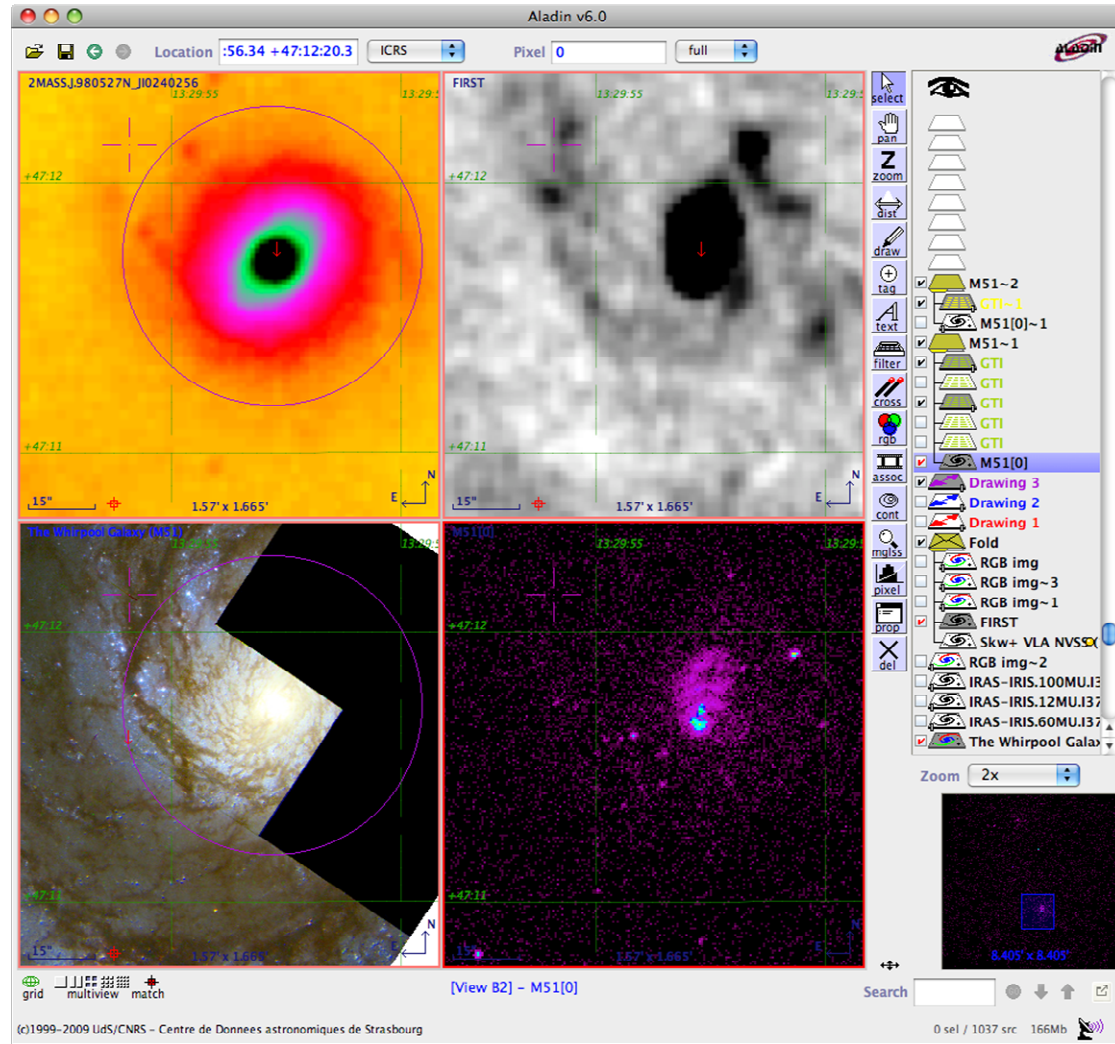
Mellinger colored

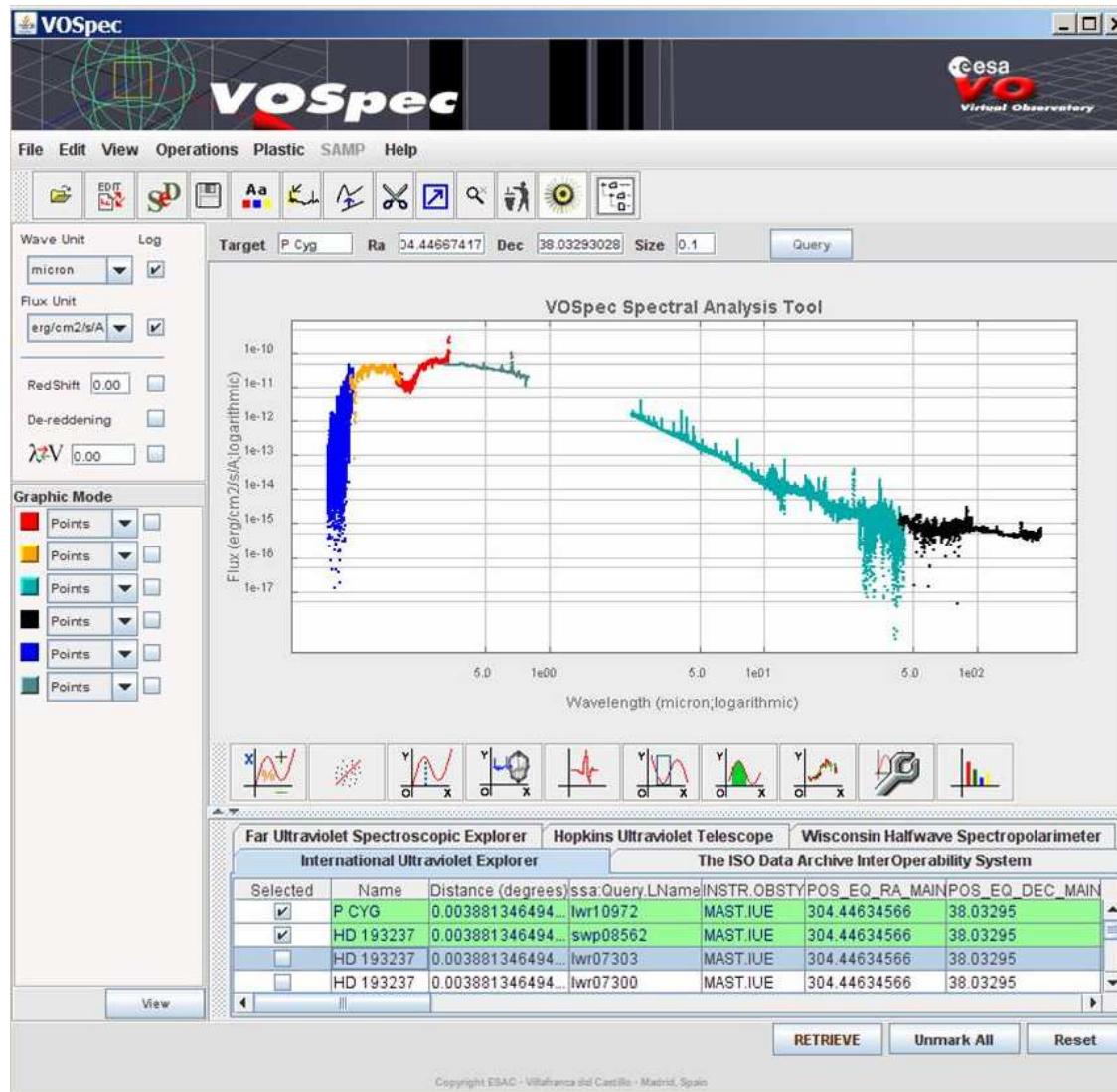
A view from the sky



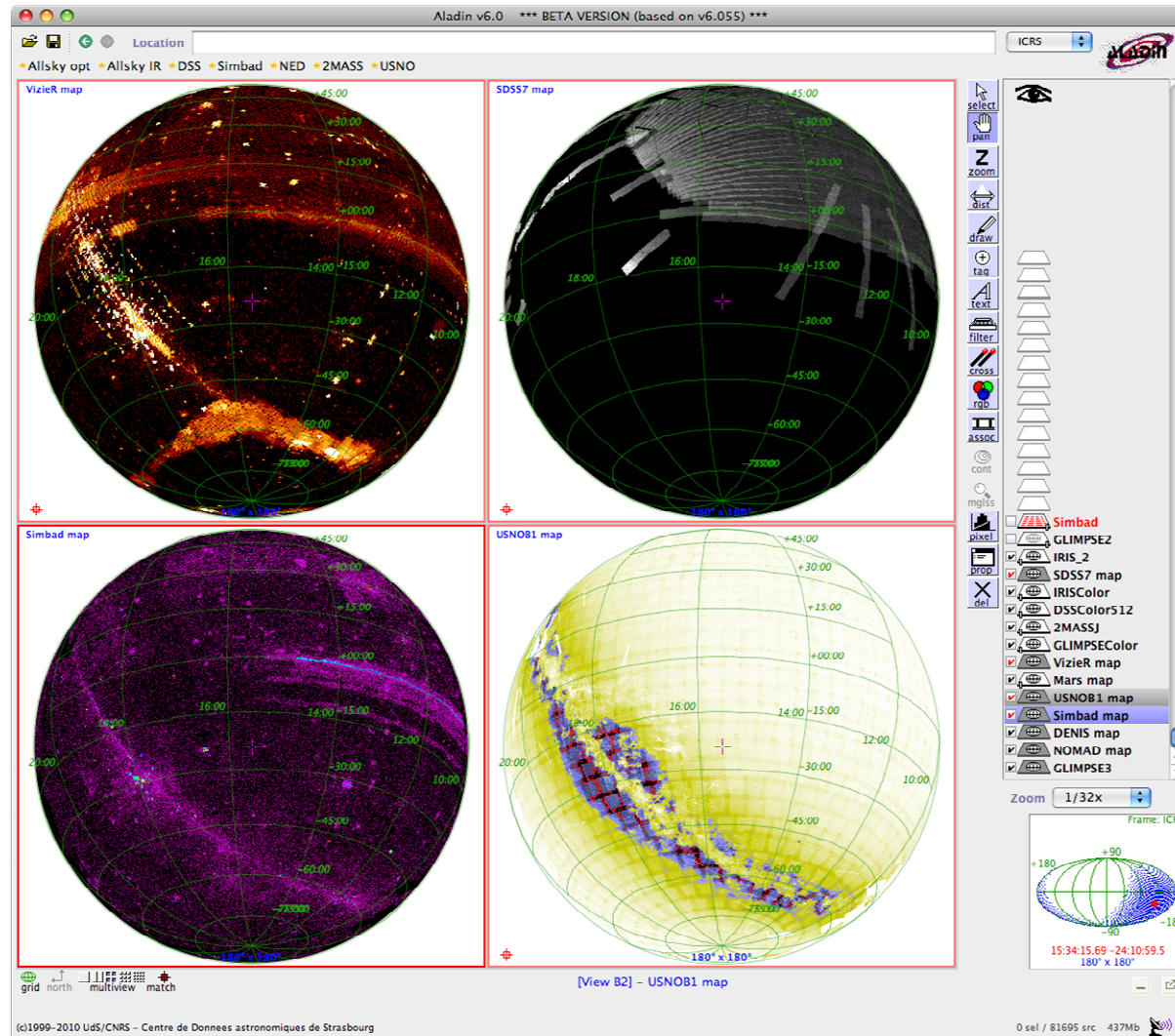
360° x 180°







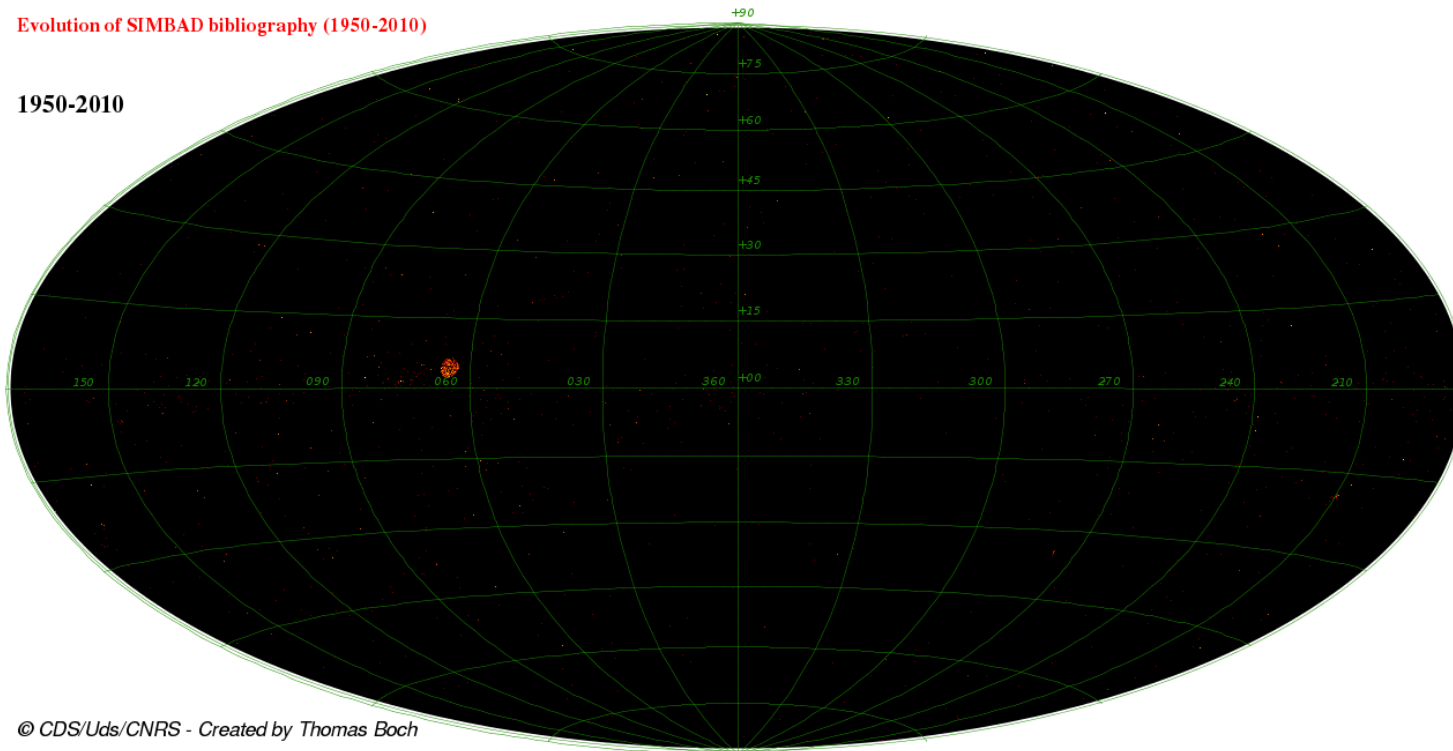




# An example of construction on the long term Constructing knowledge from bibliography (i.e., research results) SIMBAD view

Evolution of SIMBAD bibliography (1950-2010)

1950-2010



© CDS/Uds/CNRS - Created by Thomas Boch

Aladin v6.9 \*\*\* BETA VERSION (based on v6.910) \*\*\*

Position 040.40045 -01.00678 Effacer Référentiel Gal

Allsky opt Allsky IR DSS Simbad NED PPMX 2MASS

Mellinger colored

344.3° x 145.2°

simbad query - MAIN\_ID - Main identifier for an object

...	MAIN ID	OTYPE	RA	DEC	PMRA	PMDEC	PLX VALUE	RV VALUE	GA...	GA...	G...	G. SP T...	MORPH TYPE	FILTER...	U	FLUX ERROR	FLUX BIBCODE	F I F I
1	<a href="#">HESS J1857+026</a>	gamma	284.2958	+02.6667														
3	<a href="#">HESS J1858+020</a>	gamma	284.5833	+02.0900														
2	<a href="#">MGRO J1908+06</a>	gamma	287.18	+06.18														
7	<a href="#">HESS J1912+101</a>	gamma	288.0	+10.1														
2	<a href="#">HESS J1923+141</a>	gamma	290.75	+14.10														
7	<a href="#">QSO B2005-489</a>	BL Lac	302.35579...	-48.83158931	15.80	-2.20		20538				B...						
2	<a href="#">QSO B2356-309</a>	BL Lac	359.783050	-30.628031				45457	0....	0....	90							

(c) 2010 UDS/CNRS - by CDS - Distributed under GPL v3 licence

57 sel / 57 src 92Mo