# From astronomical knowledge bases to astronomical ontologies

## Andrea Preite-Martinez
## Sebastien Derriere

# Astronomical Knowledge Bases

- Astronomers already have built various knowledge bases:
    - IAU thesaurus
    - UCDs
    - SIMBAD object types
    - Keywords from journals
- How to build ontologies on top of these ?
- Possible use cases

# IAU thesaurus

- Compilation by astronomy librarians for on-line access to information
- Elaborated in the early 1990s
- Hierarchy of ~1500 terms
  - ~3300 expressions
  - covers various topics
- See Hernandez and Mothe's talk

# IAU thesaurus

# UCDs

- Description of astronomical quantities, measurements
- Hierarchical structure, ~450 words
- Originate from large collection of actual astronomical catalogue descriptions
- Revised by the IVOA for Virtual Observatories

# UCDs



Standardized description of what is measured in tables.

# SIMBAD object types

- Hierarchical classification of astronomical objects in SIMBAD
- ~150 terms
- Link with the dictionary of nomenclature

# SIMBAD object types

```
14.09.08.0:      SN            SN*     SuperNova
14.09.09.0:      Symbiotic*    Sy*     Symbiotic Star
14.14.00.0:    Sub-stellar     su*     Sub-stellar object
14.14.02.0:      Planet?       Pl?     Extra-solar Planet Candidate
15.00.00.0: Galaxy             G       Galaxy
15.01.00.0:   PartofG          PoG     Part of a Galaxy
15.02.00.0:   GinCl            GiC     Galaxy in Cluster of Galaxies
15.03.00.0:   GinGroup         GiG     Galaxy in Group of Galaxies
15.04.00.0:   GinPair          GiP     Galaxy in Pair of Galaxies
15.05.00.0:   High_z_G         HzG     Galaxy with high redshift
15.06.00.0:   AbsLineSystem    ALS     Absorption Line system
15.06.01.0:     Ly-alpha_ALS   LyA     Ly alpha Absorption Line system
15.06.02.0:     DLy-alpha_ALS  DLA     Damped Ly-alpha Absorption Line system
15.06.03.0:     metal_ALS      mAL     metallic Absorption Line system
15.06.05.0:     Ly-limit_ALS   LLS     Lyman limit system
15.06.08.0:     Broad_ALS      BAL     Broad Absorption Line system
15.07.00.0:   RadioG           rG      Radio Galaxy
15.08.00.0:   HII_G            H2G     HII Galaxy
15.09.00.0:   LSB_G            LSB     Low Surface Brightness Galaxy
15.10.00.0:   AGN_Candidate    AG?     Possible Active Galaxy Nucleus
15.10.07.0:     QSO_Candidate  Q?      Possible Quasar
15.11.00.0:   EmG              EmG     Emission-line galaxy
```

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop http://vizier.u-strasbg.fr/cgi-bin/Dic Search Print

Home Bookmarks mozilla.org Latest Builds Google

**Dictionary of Nomenclature of Celestial Objects**

CDS · Simbad · VizieR · Aladin · Catalogues · Nomenclature · Biblio · Tutorial · Developer's corner

**(Last CDS update: 09-Sep-2005)**

Result of query: info cati M$

| Acronym | Use Format | Year | 1st Author | Obj. Type |
|---------|-----------|------|-----------|-----------|
| (LMC M) | [MYM2001] NNN LMC M HHMMm+DDMM | 2001 | MIZUNO N.+ | CO Cloud |
| M | M NNN | 1850 | MESSIER C. | (Opt) |
| M | M 1-NN 2-NN 3-NN 4-NN | 1946 | MINKOWSKI R. | PN |
| M | M NNN | 1975 | MAFFEI P. | V* |
| (M) | GCM +LL.ll+BB.bb | 1981 | GUSTEN R.+ | MCld |
| (M) | [GVC73] {M} R.N | 1973 | GIOVANELLI R.+ | Concentration |
| (M) | [H68] {M} {M} R | 1968 | HULSBOSCH A.N.M. | HVC |
| (M) | [M59] NN | 1959 | MANOVA G.A. | Em* |
| (M) | [M61a] NN | 1961 | MINKOWSKI R. | G in ClG |
| (M) | [MAG95] NNN | 1995 | MINNITI D.+ | GCl |
| (M) | Mills HH+{D}A | 1952 | MILLS B.Y. | (Rad) |
| (M) | [MLV92] NNNNNN | 1992 | MAGNIER E.A.+ | * |
| (M) | MM NN | 1965 | MORAN M. | (Rad) |
| (M) | MSH HH+D-NN | 1958 | MILLS B.Y.+ | (Rad) |

---

# Keywords from journals

- Normalized list of keywords for the main astronomical journals

- Classification of bibliographical references

PHYSICAL DATA AND PROCESSES

acceleration of particles
accretion, accretion disks
astrobiology
astrochemistry
atomic data
atomic processes
black hole physics
conduction
convection
dense matter
diffusion
elementary particles
equation of state
gravitation
gravitational lensing
gravitational waves
hydrodynamics
instabilities
line: formation
line: identification
line: profiles

# Challenges

- Combine these knowledge bases to perform complex queries
- Combine several sources of (meta)data : catalogues, journals, ...
- Help astronomers to find resources in a growing list of heterogeneous sources
    – Virtual Observatory's Registry
    – Advanced information retrieval, metadata mining

# Towards Ontologies

- What technologies ?
    – Semantic Web
    – Ontologies
- How they connect to the VO
    – Registry, UCD, ADQL
- What components can be built?
- Suggested implementations

# What technologies?

- Ontologies and Semantic Web are very active research domains

- XML-based

- Strong support by W3C – definition of open standards (approved after years of discussions!)

- Interest of many commercial companies in these techniques

# Ontologies

- An ontology is a formal description of a set of concepts and their relationships to each other

- Why develop one ?  (Ontology 101)
  - to share common understanding of the structure of information among people or software agents
  - to enable reuse of domain knowledge
  - to make domain assumptions explicit
  - to analyze domain knowledge

# Ontologies

- Ontologies rely on Description Logics
- Reasoners can make inferences
  - defined concepts vs primitive concepts
  - necessary (sufficient) conditions
- Building an ontology is an iterative, collaborative process:
  - domain ontologies
  - task ontologies

# Ontologies

- Ontologies can be edited/stored with different tools/formats:
  - DAML + OIL (Ontology Inference Layer)
  - OWL (Web Ontology Language)
- Quite mature editors:
  - OILed
  - Protégé

# Ontologies



**W3C Recommendation**

## OWL Web Ontology Language Guide

W3C Recommendation 10 February 2004

**This version:**
http://www.w3.org/TR/2004/REC-owl-guide-20040210/
**Latest version:**
http://www.w3.org/TR/owl-guide/
**Previous version:**
http://www.w3.org/TR/2003/PR-owl-guide-20031215/
**Editors:**
Michael K. Smith, Electronic Data Systems, michael.smith@eds.com

Chris Welty, IBM Research, chris.welty@us.ibm.com

Deborah L. McGuinness, Stanford University, dlm@ksl.stanford.edu

Please refer to the **errata** for this document, which may include some normative corrections.

See also translations.

Copyright © 2004 W3C® (MIT, ERCIM , Keio), All Rights Reserved. W3C liability, trademark, document use and software licensing rules apply.

## Abstract

The World Wide Web as it is currently constituted resembles a poorly mapped geography. Our insight into the documents and capabilities available are based on keyword searches, abetted by clever use of document connectivity and usage patterns. The sheer mass of this data is unmanageable without powerful tool support. In order to map this terrain more precisely, computational agents require machine-readable descriptions of the content and capabilities of Web accessible resources. These descriptions

---

# Semantic Web



**W3C®** Technology and Society **domain** | Semantic Web **Activity**

## Semantic Web

The **Semantic Web** provides a common framework that allows **data** to be shared and reused across application, enterprise, and community boundaries. It is a collaborative effort led by W3C with participation from a large number of researchers and industrial partners. It is based on the Resource Description Framework (RDF), which integrates a variety of applications using XML for syntax and URIs for naming.

"The Semantic Web is an extension of the current web in which information is given well-defined meaning, better enabling computers and people to work in cooperation." -- Tim Berners-Lee, James Hendler, Ora Lassila, The Semantic Web, Scientific American, May 2001

**On this page:** Activity Statement | Specifications | Publications | Presentations | Groups

**Nearby:** Advanced Development | SWAD-Europe | Simile | Semantic Web Coordination | RDF | RDF Core | RDF Data Access | Web Ontology | Best Practices and Deployment | Interest Group | Developer Tools

# Semantic Web

- Data shared and reused among application and community boundaries

- RDF: Resource Description Framework (XML-based, with URIs and namespaces)

- Goal:
    - communication between software agents, users
    - discover valuable applications (resource discovery) and exploitation paths (workflows)

# Application in biology

- BioHaystack (IBM Watson Research)

• Integration of data from heterogeneous databases
  • access protocols
  • data formats
  • softwares
• RDF as underlying model
• Link to myGRID

# Application to the VO

- Astronomy is a user community which is familiar with many issues motivating the Semantic Web research:

• Resource description in the VORegistry, with unique identifiers and XML format: mapping to RDF?
• VOTable, with metadata and data grouped in the same document: metadata sharing is at the core of the semantic web
• UCD and the IAU thesaurus for the semantic description will be very valuable resources for building an astronomical ontology
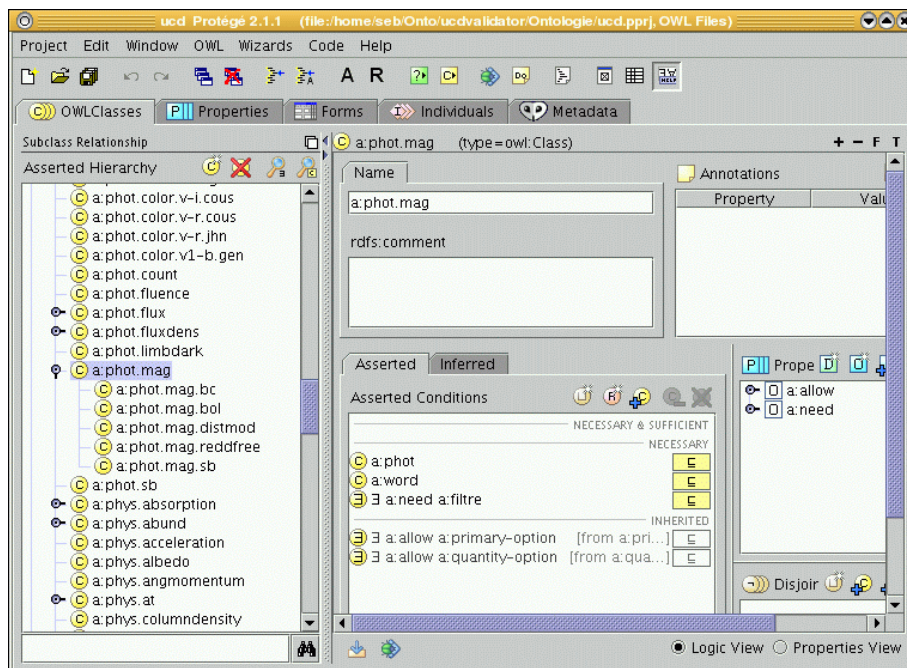
# How to proceed ?

- Don't only try to build one single ontology to describe all astronomy

- Start with well-defined small ontologies, with simple use cases

- Then, build links between these ontologies and allow more complex problems to be solved
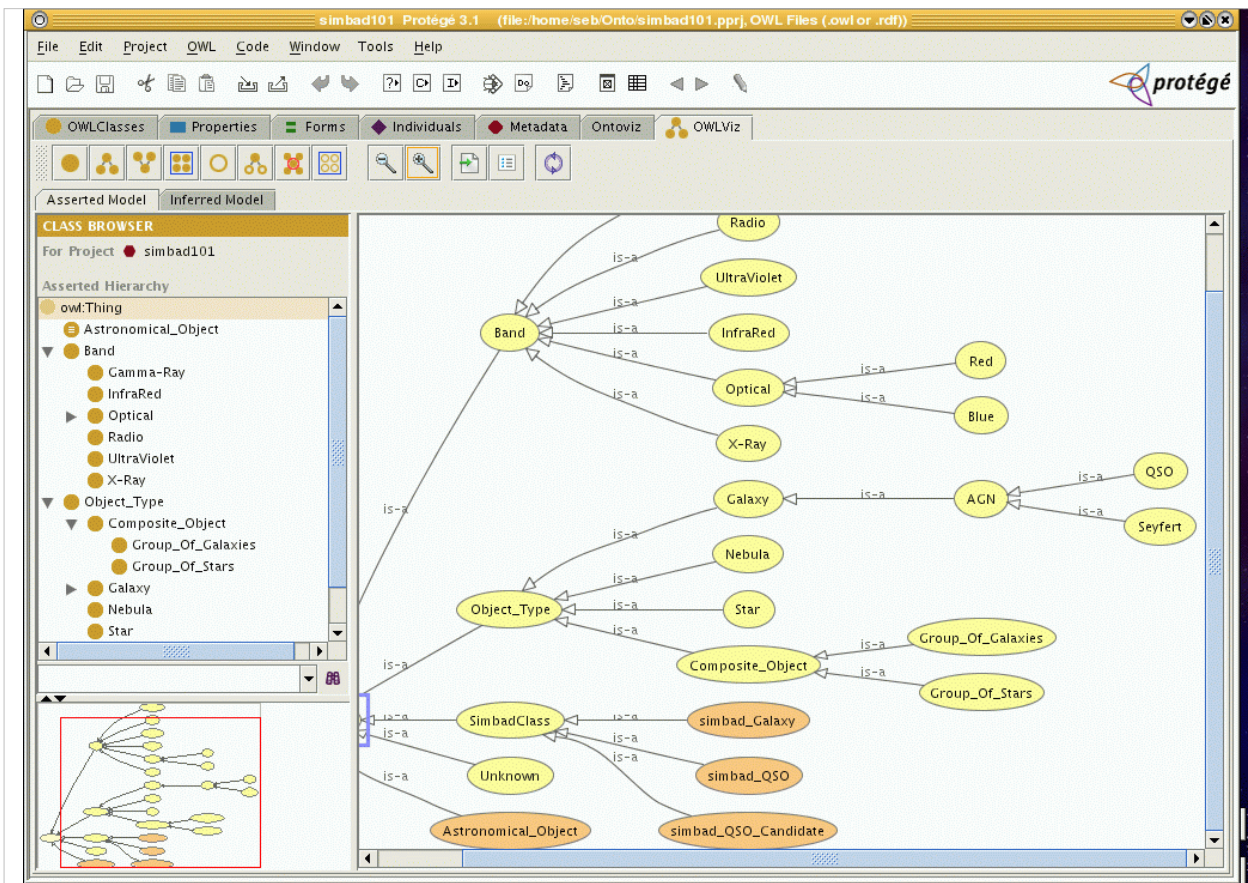
# UCD and ontologies

- Collaboration between CDS and LORIA
- Syntactic rules for building UCDs from words
- UCD validation
- Application to UCD assignment:
  – given a description, which UCD is relevant?

# UCD and ontologies

# Ontology of object types

- Collaboration CDS–INAF in VOTech

- Based on SIMBAD object types

    – improve object hierarchical description

    – link between object type and wl/frequency domain

- Allow multiple classification (possible in SIMBAD4)

- Link with UCD and IAU thesaurus

# Use case 1

- Use ontology to validate entries during registry population (for data curators):
  - interactive keywords refinement
  - check if global keywords consistent with specific measurements (UCDs)

# Use case 2

- Assist users in registry exploration
- Goal: retrieve relevant datasets from the registry, using the ontology
  - search for "young star" would be interpreted by a reasoner, and recursively matched to all sub-categories of "young star"
  - in case of null result, suggest request broadening (TTauri star -> young star)

# Use case 3

- Help automated object names detection in articles
  - analyze the context of the paper
    - keywords
    - object types
    - measurements (UCDs)
  - use dictionary of nomenclature
  - identify object names, and most likely acronym

# Questions ?