# Virtual Observatory

Jaroslav Vazny

Department of Theoretical Physics and Astrophysics

Masarykova Univerzita

A thesis submitted for the degree of

*Master*

Yet to be decided

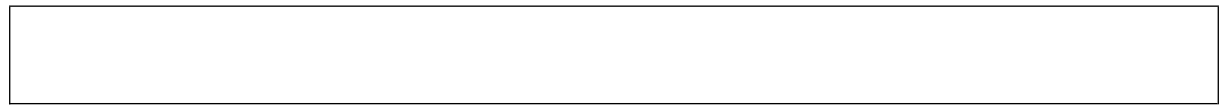I would like to dedicate this thesis to my loving parents ...

# Acknowledgements

And I would like to acknowledge ...

# Abstract

This is where you write your abstract ...

# Contents

# List of Figures

# Introduction

From the dawn of existence astronomy has always been starved for data, but in the last few decades the situation has changed and now we are facing the data flood of biblical proportions. The data are not just increasing in size but in complexity and dimensionality. Ball and Schade [2010] Astroinformatics is the new field of science which has emerged from this technology driven progress. Virtual Observatory, Machine learning, Data Mining, Grid computing are just few examples of new tools available to scientist.
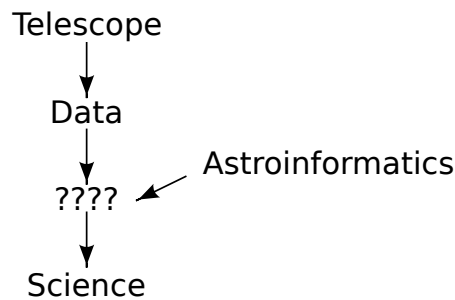
Telescope
↓
Data
↓
???? ← Astroinformatics
↓
Science

Figure 1: Astroinformatics in the context of astronomy Ball and Schade [2010]

Of course astronomers are not alone and particle physics, biology and other sciences are also in the vanguard of the data intensive science. This is great opportunity for interdisciplinary collaboration.

This work deals with the problem of semi-automatic procedures for finding Be stars candidates in the astronomy surveys. More than straight forward process it's trail and error approach probing new possibilities with rather interesting that useful results.

The aim of this work is to be introductory to the technologies of Virtual Observatory and Data Mining and for this reason it is intended to have following properties:

- Main Chapters starts with questions answered in the text and diagram to ease orietation,

- is full of examples,

- is non-linear in nature,

- is meant to be compact and consistent,

- is far from complete.

Chapter one is an introduction to the technologies related to Virtual Observatory. The motivation behind the concept is given without paying too much attention to historical details. Main standards and protocols are discussed and explained. Important aspect are demonstrated on numerous examples. Chapter two is an introduction to Machine Learning and Data Mining in the context of astrophysics. Only methods used in practical part of this work are described in detail: Decision Trees and Support Vector Machines. Examples of several classifications are demonstrated. Third chapter introduces problematic of Be stars. Chapter Four is practical application of previously described technologies and methods. Training data of confirmed Be stars from Ondrejov are correlated with others catalogues to obtain color indexes and spectra. Results are processed by Data Mining algorithms using several libraries and tools. In the last chapter achieved results are critically discussed.

Activities related to this work go beyond this text. Wiki pages were created to present the results and discuss related topic with supervisor as well as with others scientist around the world. Several programs were created to analyze and process acquired data. Source codes were maintained by GIT version system allowing easy sharing. All software used and produced are open source.

# Chapter 1

# Virtual Observatory (VO)

1 What is the motivation behind Virtual Observatory? Is data avalanche
2 problem only in astronomy? What is IVOA? What is Virtual Observatory
3 architecture, standards and protocols

## 1.1 Data avalanche: Opportunity or disaster?

There are two important trends in current astronomy surveys:

- Size: The cumulative compressed data holdings of the ESO archive will reach 1 PetaByte by 2012 Hanisch and Quinn [2010]. Projects like Large Synoptic Survey Telescope (LSST) will produce about 30 TB per night, leading to a total database over the ten years of operations of 60 PB for the raw data.

- Complexity: Modern surveys will cover the sky in different wavebands, from gamma- and X-rays, optical, infrared, through to radio. The ability to crosscorelate these observations toghether may lead to the new understanding of physical phenomenas.

Such amount of data is not possible transfer over the network. It imply they are heterogenous, distributed and decentralized in nature.

There is an interesting analogy with the problem (and the solution) which had scientist during LEP project at CERN. Their problem was too many documents in different formats. Tim Berners-Lee [1] designed set of protocols (URIs, HTTP and HTML) which allowed link and share documents Berners-Lee and Cailliau [1990]. This was recognized as generaly useful and World Wide Web was born. An important role plays the World Wide Web Consortium (W3C) in developing Web standards. [2]

---

[1] Sir Timothy John "Tim" Berners-Lee. British engineer and computer scientist and MIT professor credited with inventing the World Wide Web.

[2] Prior to its creation, incompatible versions of HTML were offered by different vendors, increasing the potential for inconsistency between web pages.

## 1.2   International Virtual Observatory Alliance (IVOA)

What is neccessary is sets of standards and protocols to deal with heterogenous distributed data and the authority which encourages their implementation. Such authority is the International Virtual Observatory Alliance (IVOA). It comprises 19 VO programs from Argentina, Armenia, Australia, Brazil, Canada, China, Europe, France, Germany, Hungary, India, Italy, Japan, Russia, Spain, the United Kingdom, and the United States and intergovernmental organizations (ESA and ESO)Hanisch and Quinn [2010].



Figure 1.1: IVOA members

Standards specifications can be obtained on http://www.ivoa.net/.

## 1.3   Architecture

The Virtual Observatory is the middle layer framework which connects the Resource Layer to the User Layer in a seamless and transparent manner. The objective is to improve and unify access to astronomical data and services.
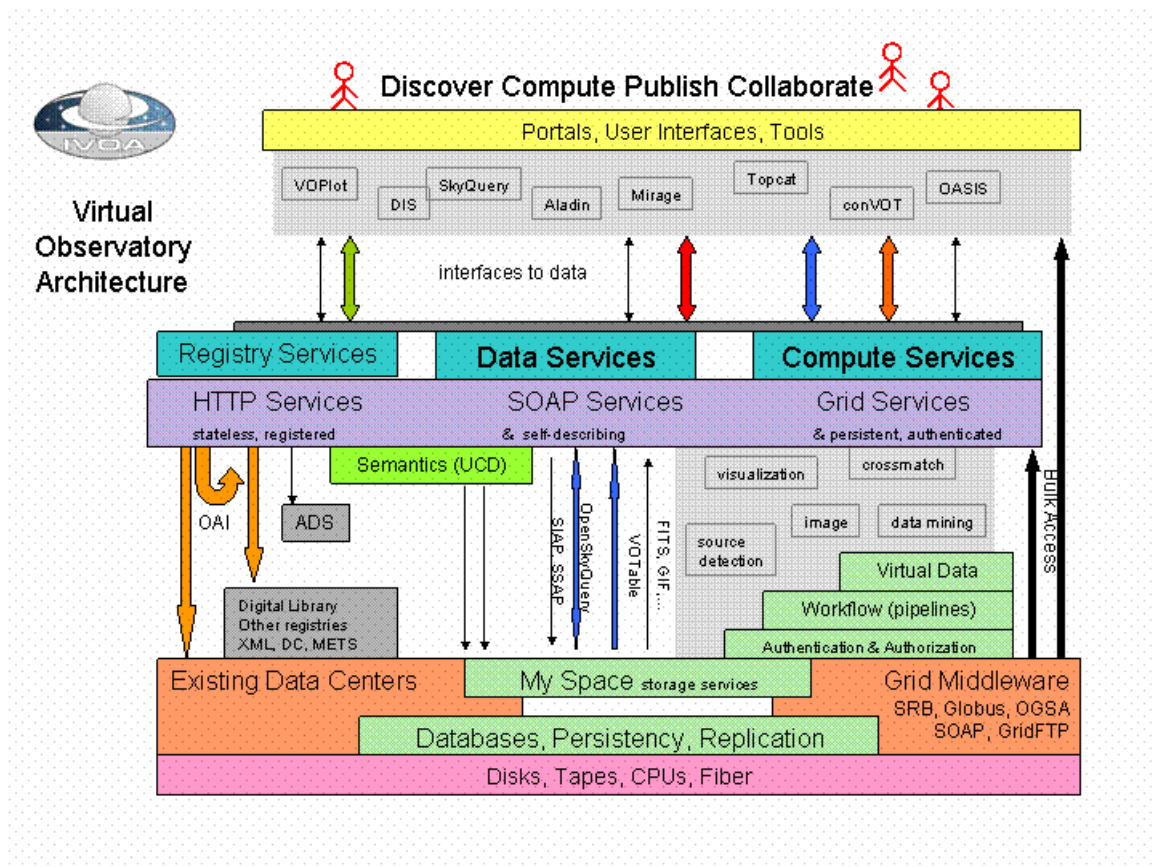
Figure 1.2: VO Architecture

The Architecture is depicted on the figure 1.2.The level of abstraction goes from top to bottom. Starting with iterfaces, used by people or applicatian to discover resources. Next level is the service layer implemented by standard protocols, followed by the hardware level where actual data are stored. This onion like structure hide the complexity of the lower layer and provide data and metadata to the higher layer. This concept is similar to TCP/IP [1] protocol.

The essence of VO architecture is service orientation. Each service is autonomous with well defined defined boundaries. Very important aspect of VO implementation is the adoption of formats and protocols used in astronomy (FITS) and computers science (XML [2] , Web service [3] SOAP [4]) for many years. In other words VO does not reinvent the wheell but it's stands on the shoulders of giants.

## 1.4 VOResources

A resource is a general term referring to a VO element that can be described in terms of who curates or maintains it and which can be given a name and a unique identifier. Just about anything can be a resource: it can be an abstract idea, such as sky coverage or an instrumental setup, or it can be fairly concrete, like an organization or a data collection. Benson et al. [2009]

UML [5] diagram of the resource in on the figure 1.3. Resource can be a generalization of organization, data collection, apllication or service. Organization can be linked together with other organization. The same is true for data collection. Organization ia a generalization of and/or provider which can own zero to N services. Publisher can have zero to N resources.

---

[1]TCP/IP (Transmission Control Protocol/Internet Protocol). The basic communication language or protocol of the Internet.

[2]Extensible Markup Language (XML) is a set of rules for encoding documents in machine-readable form.

[3]method of communication between two electronic devices over a network.

[4]Simple Object Access Protocol, is a protocol specification for exchanging structured information in the implementation of Web Services in computer networks.

[5]Unified Modeling Language. Standardized general-purpose modeling language in the field of object-oriented software engineering.
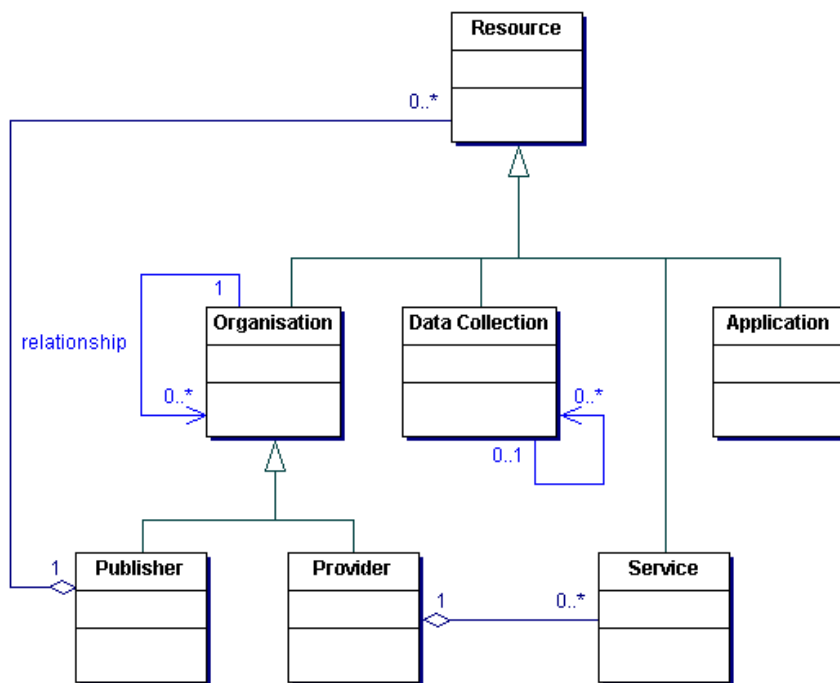
Figure 1.3: UML diagram of VOResource

Example of resources

```
1  stilts regquery query="shortName like 'AIASCR'"
2  regurl=http://registry.euro-vo.org/services/RegistrySearch
3  ofmt=votable-tabledata > resourceExample.vot
```

```
1   <?xml version='1.0'?>
2   <VOTABLE version="1.1"
3    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
4    xsi:schemaLocation="http://www.ivoa.net/xml/VOTable/v1.1 http://www.ivoa.net/
         xml/VOTable/v1.1"
5    xmlns="http://www.ivoa.net/xml/VOTable/v1.1">
6   <!--
7    !  VOTable written by STIL version 3.0 (uk.ac.starlink.votable.VOTableWriter)
8    !  at 2011-03-24T00:45:59
9    !-->
10  <RESOURCE>
11  <TABLE nrows="1">
12  <LINK title="Registry Location" href="http://registry.euro-vo.org/services/
         RegistrySearch"/>
13  <PARAM arraysize="23" datatype="char" name="Registry Query" value="shortName
         like 'AIASCR'">
14  <DESCRIPTION>Text of query made to the registry</DESCRIPTION>
15  </PARAM>
16  .
17  .
18  .
19  <DATA>
20  <TABLEDATA>
21    <TR>
22      <TD>ivo://asu.cas.cz</TD>
23      <TD>AIASCR</TD>
24      <TD>Astronomical Institute of the Academy of Sciences of the Czech Republic
             Naming Authority</TD>
25      <TD>Astronomical Institute of the Academy of Sciences of the Czech Republic
         </TD>
26      <TD>http://stelweb.asu.cas.cz/web/index/index-en.php</TD>
27      <TD>Petr Skoda &lt;skoda@sunstel.asu.cas.cz&gt;</TD>
28    </TR>
29  </TABLEDATA>
30  </DATA>
31  </TABLE>
```

```
32   </RESOURCE>
33   </VOTABLE>
```

## 1.5   Data Access Protocols

### 1.5.1   Cone Search Protocol

### 1.5.2   Simple Image Access Protocol

### 1.5.3   Simple Spectra Access Protocol

## 1.6   Data Formats

### 1.6.1   VOTable

### 1.6.2   FITS

**Motivation**

"An archival format must be utterly portable and self-describing, on the assumption that, apart from the transcription device, neither the software nor the hardware that wrote the data will be available when the data are read." Council [1995]

FITS (Flexible Image Transport System) was originaly created for exchange of radio astronomy images between WSRT [1] and the VLA [2] Schlesinger [1997]. It is now used as a file format to store, transmit, and manipulate scientific data and it is (thanks to it's revolutionary design) de facto standard in astronomy.

**Structure**

One file can contain several HDUs (Header Data Units).The first part of each HDU is the header, composed of ASCII card images containing keyword=value statements that describe the size, format, and structure of the data that follow.

- Primary header and data unit (HDU).

- Conforming Extensions (optional).

- Other special records (optional, restricted).

Standards and documents related to FITS are maintaned by IAUFWG [3] and aviable at http://fits.gsfc.nasa.gov.

---

[1]Westerbork Synthesis Radio Telescope
[2]Very Large Array
[3]International Astronomical Union FITS

### 1.6.2.1 Examples

There are many libraries for working with FITS files. The official list is aviable at
http://fits.gsfc.nasa.gov/fits_libraries.html. PyFITS, library for Python programming language was used for following examples. PyFITS is a development project of the Science Software Branch at the Space Telescope Science Institute http://www.stsci.edu/resources/software_hardware/pyfits.

Reading FITS headers.

```
In [1]: import pyfits
In [2]: hdulist = pyfits.open('spSpec-53237-1886-248.fit')
In [3]: hdulist.info()
Filename: spSpec-53237-1886-248.fit
No.    Name        Type      Cards   Dimensions  Format
0    PRIMARY    PrimaryHDU   213  (3874, 5)   float32
1               BinTableHDU   54  6R x 23C    [1E, 1E, ...
2               BinTableHDU   54  44R x 23C   [1E, 1E, ...
3               BinTableHDU   18  1R x 5C     [1E, 1E, ...
4               BinTableHDU   32  53R x 12C   [1J, 1J, ...
5               BinTableHDU   26  36R x 9C    [19A, 1E, ...
6               BinTableHDU   14  3874R x 3C  [1J, 1J, 1E]
```

Printing primary HDU.

```
In [4]: print hdulist[0].header
-------> print(hdulist[0].header)
DATE-OBS= '2004-08-20'       / 1st row - TAI date
TAIHMS = '10:36:18.11'       / 1st row - TAI time (HH:MM:SS.SS) (TAI-UT = appr
TIMESYS = 'tai   '           / TAI, not UTC
TAI-BEG =       4599713999.00 / Exposure Start Time
TAI-END =       4599717089.00 / Exposure End Time
MJD     =              53237 / MJD of observation
MJDLIST = '53237 '           /
VERSION = 'v3_140_0'         / version of IOP
CAMVER = 'SPEC1 v4_8'        / Camera code version
OBSERVER= 'prn   '
OBSCOMM = 'science '
TELESCOP= 'SDSS 2.5-M'       / Sloan Digital Sky Survey
```

Updating FITS file.

```
In [16]: prihdr = hdulist[0].header
In [17]: prihdr.update('observer', 'Astar')
In [18]: prihdr.add_history('I updated this file 3/27/11')
```

Example from program pf (plot fits) created for purposes of this work to plot $H\alpha$ emission in the spectra.

```python
def read(file):
    """ Read fits file. Convert wavelength to angstroms """
    data = pyfits.getdata(file)
    w = lambda x : 10.0**(3.5796 + x*10.0**(-4))
    x = np.arange(1,data[0].size + 1)
    xx  = w(x) # convert to actual wavelenght
    return np.asarray([xx, data[0]])

def plot(file,xdata,ydata,spLine):
    fig = plt.figure()
    ax = fig.add_subplot(111)
    graph = ax.plot(xdata,ydata, 'r')
    ax.set_title(file)
    ax.set_xlabel("$Wavelenght [\\AA]$")
    ax.set_ylabel("$Energy [10^{-17} erg/s/cm^2/\\AA]$")
    ax.axvline(x=spLine, color = 'g', ls ='--')
```

I would also like to add an extra bookmark in acroread like so ...

# References

N.M. Ball and D. Schade. ASTROINFORMATICS IN CANADA. 2010. v, 1

K. Benson, R. Plante, E. Auden, et al. IVOA Registry Interfaces. *IVOA Working Draft*, 2009. 6

T. Berners-Lee and R. Cailliau. WorldWideWeb: Proposal for a HyperText project. *European Particle Physics Laboratory (CERN)*, 1990. 3

National Research Council. Preserving Scientic Data on our Physical Universe. 1995. 9

Norman Gray. An RDF version of the VO Registry. *IVOA Note (V1.0, 2007 September 20)*, 2007.

RJ Hanisch and PJ Quinn. The international virtual observatory. *Retrieved from http://www. ivoa. net/pub/info/TheIVOA. pdf on*, 24, 2010. 3, 4

B. Schlesinger. A Users Guide for the Flexible Image Transport System. 1997. 9