# R&D @ CDS
## and other developments

André Schaaff on behalf of the CDS Team

**CDS Scientific Council 2016**

---

## □ Why R&D ?

- The CDS team has always spent time on R&D activities to follow the technological evolutions
- These evolutions are now very fast and in various fields (interactions, visualization, mobility, components, Big Data & Open data, Clouds, etc.) with a lot of actors in both the commercial and the Open Source domains
- It is becoming hard to test and evaluate everything in addition to the everyday work
- The R&D activity is now well identifed at CDS, it involves several persons of the staff with the help of interns and short contracts
- It provides also topics to present and discuss during the Infusion meetings

# Internship programme

- We have now an internship programme to hire IT students (11 in 2016) to work with us on several topics, R&D and other developments
- During the three last years
  - 31 interns, 358 weeks / 6.9 years
- + short contracts following a few internships
  - To push the work to the production side
  - To work on short developments during the Summer
- Possible hiring on projects

# Internships in 2016

- Aladin fisheye visualization of HiPS surveys
- A monitoring system to track and visualize the number of queries of CDS services
- XObsCoreFits, a user interface to provide data to be ingested into CDS VizieR (images, spectra, time-series)
- Development of Python code for spectral analysis of CALIFA IFU data cubes
- Analysis of photometric flux conservation in the use of HiPS for images and cubes
- "A la découpe" Sky progressive survey (HiPS) server
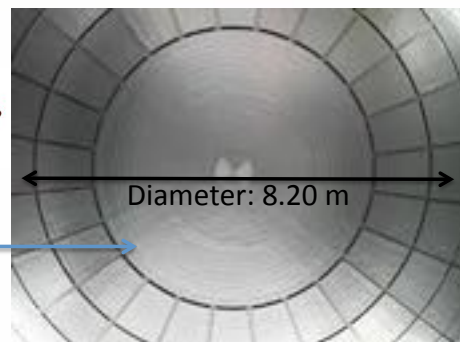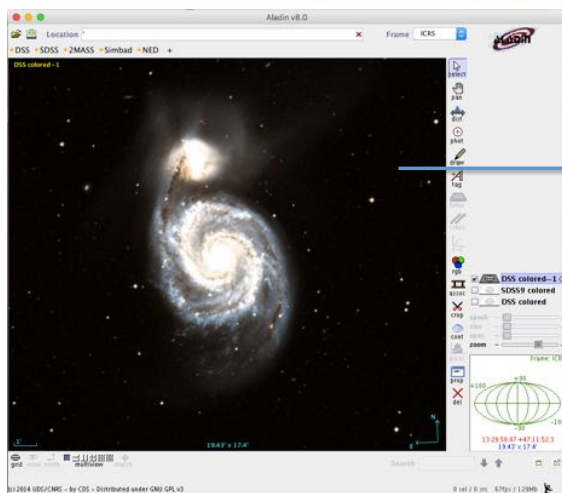
# Internships in 2016 (2)

- Development of a system to discover new Virtual Observatory services (VO Alerts)
- 3D visualization in a Web browser (large datasets, interpretation and immersion), → 3 internships
- Apache Spark and X-Match → on going
- "DevOps" at CDS, for the deployment of services in containers (Docker) → on going
- + Collaboration around "Binding Database Metadata with Scientific Papers" with L. Michel → on going
- + Participation to R&D project tutoring at ENSIIE Strasbourg

# Aladin fisheye visualization of HiPS surveys



Diameter: 8.20 m

F. Bonnarel, S. Derriere, P. Fernique, A. Schaaff, M. Wendling (JdS), B. Rota (JdS)
Intern: Arnaud Steinmetz (ENSIIE Strasbourg)

*Presented at IVOA Trieste*
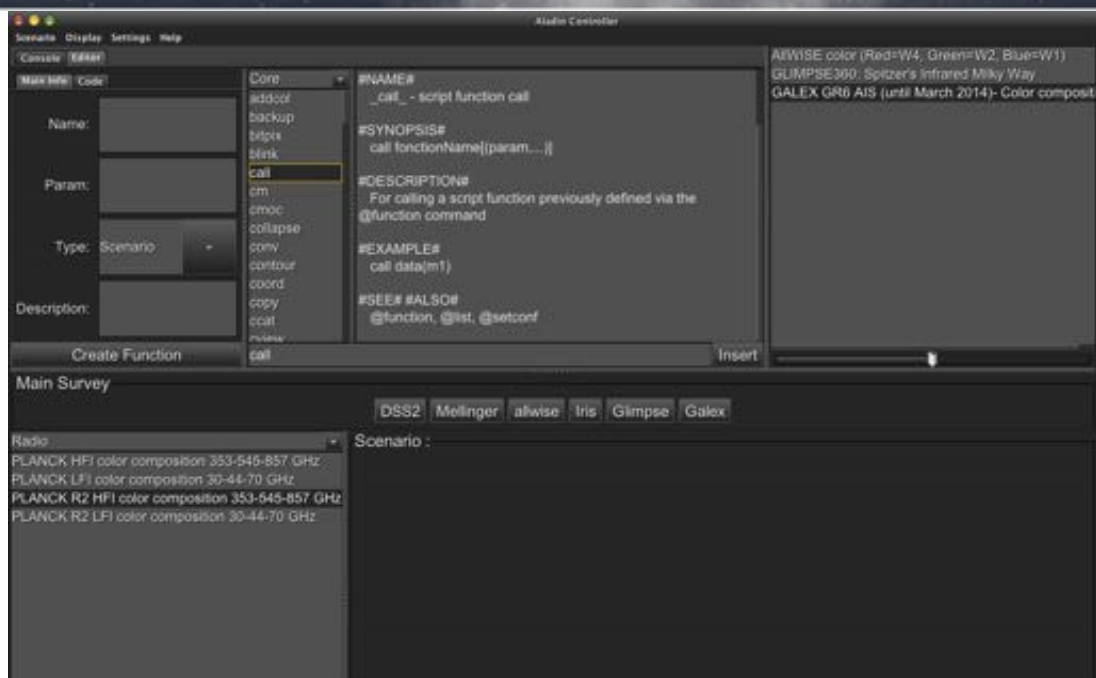
# From the screen to the dome

- How to display Aladin HiPS surveys on the dome ?
  - The "fisheye" projection we use is actually the ARC (zenithal equidistant) projection
  - In this projection angular distances to the zenith are conserved
  - Well adapted to projection on the sphere
- How to manipulate ?
  - A plugin was the best solution

# Planetarium plugin

# Illustrations



Credits: IRIS



Credits: GALEXGR6
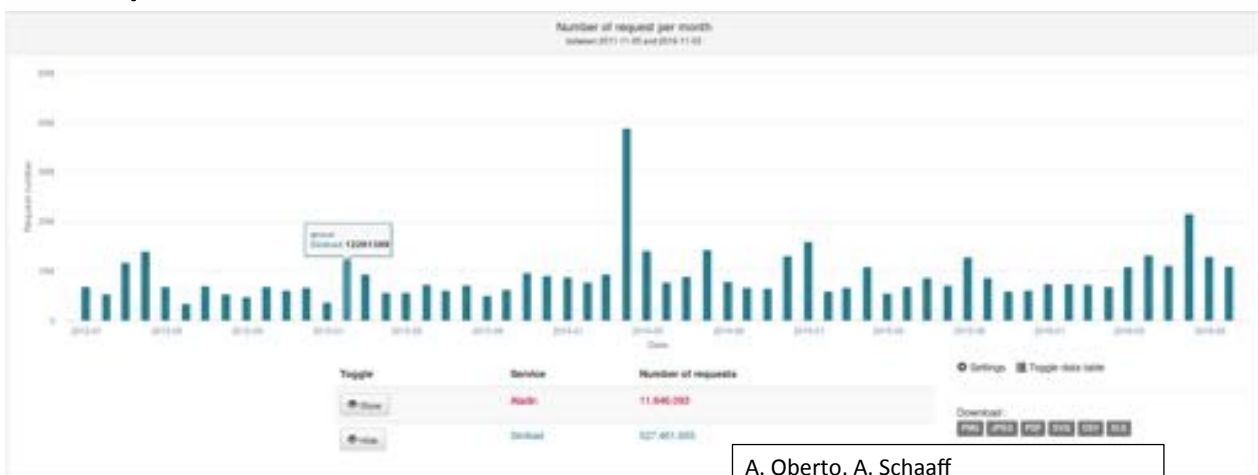
---

# A monitoring system to track and visualize the number of queries of CDS services

- Grouping the service logs in one "store" to provide on-demand statistics



A. Oberto, A. Schaaff
Intern: Thierry Lacoste (IUT Champagne)

# Slide 1

- World use of a service…



Future R&D aspects: log mining

# Slide 2

# XObsCoreFits



- Mapping between FITS documents (images, spectra, time-series) and IVOA ObsCore
- Documentalists complete a pre-processing phase done through the Saada API
- VizieR data can also be included

G. Landais, P. Ocvirk, L. Michel
Intern: Félix Royer (IUT Belfort-Montbéliard)

***Poster at ADASS Trieste***

# Spectral analysis of CALIFA IFU data cubes

- CALIFA, a collection of ~600 galaxy images in hyperspectral data cubes
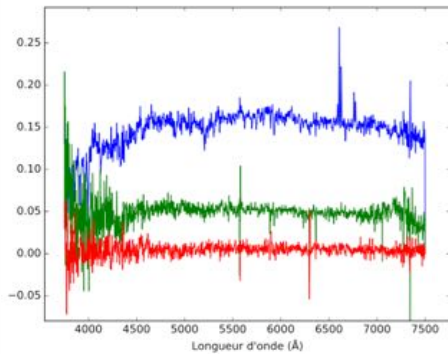- Build of velocity maps and variance



Spectra in a cube at (40,40)(blue), (20,40)(red), (0,40) (green)



Images from a cube at different wavelengths

M. Louys, C. Boily
Intern: Thibaut Buchert (ENSIIE Strasbourg)

# Spectral analysis of CALIFA IFU data cubes (2)

- The aim of was to improve an existing Python application at several levels: speed, test on a large set of cubes, code quality, documentation, etc.
- Rewriting in C of some parts with a computing time gain: 1 minute for 10 CALIFA cubes, 100 times faster



Residuals maps

Output examples of the application

NGC 0001
Fitting with adjacent pixels

# "A la découpe" Sky progressive survey (HiPS) server

- The aim was to develop a server (based on Java Servlets under Tomcat) to generate on the fly images from progressive surveys
- Use of the IVOA SIAv2 standard



Image generation from "pixels", problems

P. Fernique, T. Boch
Intern: Thomas Janowskyj (IUT Charlemagne)

# Photometric conservation in HiPS processing

*How can we assess aperture photometry degradation after HiPS transformation?*
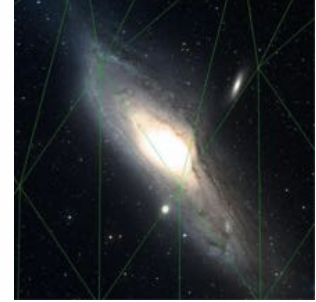
- HiPS is becoming an IVOA standard
- it is based on HealPix tessellation of the sky which defines cells of equal area on all sky and implies transformation of the original data
- We have raised with the help of an intern the question of potential quality degradation in producing the HiPS version



**Test data and procedures:**
- The IRIS data collection is a reprocessing of the IRAS Sky Survey Atlas , recalibrated using DIRBE and made of 12.5 degrees joint images at 12, 25, 60 and 100 μm
- Mean surface brightness is estimated in circular areas on both original images and corresponding areas in the HiPS representation using *Aladin* code controlled via a python module.

***Poster at ADASS Trieste***

M. Louys, F. Bonnarel, C. Bot, P. Fernique
Intern: Damien Teodori (ENSIIE Strasbourg)

**Figure 1: IRIS @100micron  HiPS visualisation within** Aladin: high resolution on the right

# Photometric conservation in HiPS processing (2)



Figure 2: Fractional difference (as a percentage) between the HIPS and original surface brightnesses as a function of the average surface brightness. The histogram of the fractional difference is shown on the right.

**Results and conclusions:**
- Relative photometric difference is extremely small (< 0.3% !!!) typically much smaller than other sources of uncertainties (fig 2)
- No obvious spatial bias in the results but further tests needed.
- Tests scheduled for source photometry on appropriate surveys
- **For IRAS/IRIS data it is fairly equivalent to measure photometry on data distributed in the HiPS format than on original ones**

---

# □ VO Alerts

- A service to notify users of newly-published VO resources relevant for a list of sky targets
  - Cone Searches, Simple Image/Spectrum Access
- Motivation: difficult to know what's new in VO
  - CDS news: all new VizieR catalogues
  - VO Registry: information must be pulled, and tested for relevance
- Similar to SimWatch, but for VO services
  - Monitor newly added references for SIMBAD objects

***Poster at ADASS Trieste***

T. Boch, S. Derriere
Intern: Loïc Gasiorowski (IUT Charlemagne)

# VO Alerts (2)

# 3D Visualization in a Web browser

- Light tool to visualize several kinds of 3D data in a Web browser (based on WebGL)

- On going work around the visualization of large datasets
  - $4096^3$ simulation data cube
  - Too large to be loaded in a Web browser
  - → data on a server + progressive visualization on the client side

A. Schaaff, D. Aubert, N. Deparis, N. Gillet, P. Ocvirk
Interns: Jérôme Desroziers (IUT Charlemagne), Thibault Bouchard et Nicolas Adam (ENSIIE Strasbourg)

# Illustration



User interface, visualization engine,
several functionalities

# Illustration (2)



Navigation in
the data

# ☐ Apache Spark and X-Match

- Evaluation of Spark in the frame of a use case, the "cross-match" of source catalogues:
  - Improvement of the existing service (one server) ?
  - Up to scale capability ?
  - Which cost (€, manpower, performances) ?
- Spark is a kind of MapReduce done mostly in memory

*A. Schaaff, F.-X. Pineau, Julien Nauroy (Université Paris Sud)*
*Interns: Paul Trehiou (UTBM), (Noémie Wali (UTBM) in 2015)*

***Oral at ADASS Trieste***

***Presented at IVOA Trieste***

---

# ☐ The data…

- Source catalogues (>10,000 available)
- Examples (number of sources):
  - 2MASS[1], 470,992,970
  - SDSS[2] DR9, 469,053,874

**Example of a ReadMe file associated to 2MASS source catalogues available through the VizieR service**



[1]2MASS, Two Micron All Sky Survey,
[2]SDSS, Sloan Digital Sky Survey

# …and the CDS "cross-match" service

- The "cross-match" service does a cross correlation of sources between (very) large catalogues (current size: $10^9$), for the full sky, a cone or a HEALPix cell.

Fuzzy join between 2 tables (A and B) of several hundred millions of data

Credits: http://healpix.jpl.nasa.gov/

Data is not distributed but organised and stored on one server

The sky is cut into diamonds of the same size, pixels, each source or sky object is a numbered pixel.

---

# First experiment result

- Input data (SDSS DR7 (primary sources) and 2MASS): 54GB and 58GB file size; 357 175 411 and 470 992 970 elements
- Output data: 49 208 820 elements

X-Match service reference time was: 10 minutes

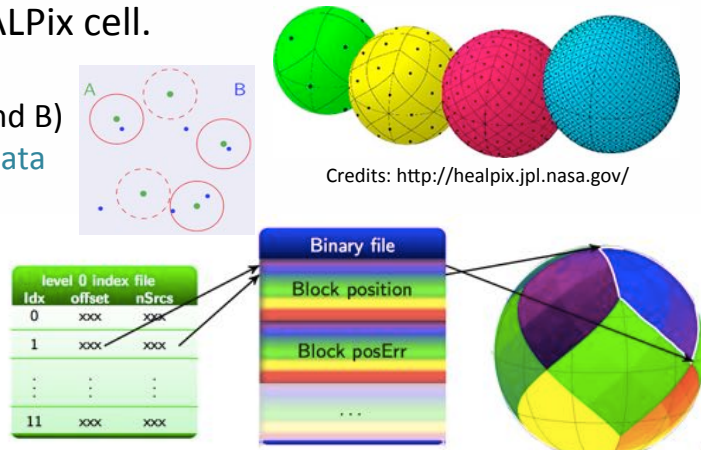| Cross-Match (source duplication done in phase 2 with all the data as output) HDFS block size= 128MB for the input files ; sdss7.csv and t 2mass.csv replicated 2 times | | | | | |
|---|---|---|---|---|---|
| HashPartitioner | 60 partitions | | | | |
| HDFS output files size | 32MB | | | | |
| Number of nodes Spark/HDFS | 5 | 7 | 9 | 10 | 11 |
| Phase 1: prepare | 23,0 | 16,0 | 14,0 | 14,0 | 13,0 |
| mapToPair (sdss7.csv) | 5,1 | 4,9 | 4,9 | 4,8 | 4,7 |
| saveAsHadoopFile (sdss7.bin) | 5,7 | 2,7 | 2,0 | 2,3 | 1,5 |
| mapToPair (2mass.csv) | 5,7 | 5,2 | 5,2 | 5,1 | 5,0 |
| saveAsHadoopFile (2mass.bin) | 6,5 | 3,6 | 1,9 | 1,6 | |
| Phase 2: join | 31,0 | 21,0 | 13,0 | 11,0 | 9,9 |
| mapToPair (sdss7.bin) | 7,2 | 4,7 | 3,5 | 3,0 | |
| flatMapToPair (2mass.bin) | 11,8 | 8,3 | 5,5 | 4,9 | 4,3 |
| saveAsTextFile (crossMatch_D.txt) | 12,0 | 7,6 | 3,4 | 2,4 | 2,3 |
| TOTAL | 54,0 | 37,0 | 27,0 | 25,0 | 22,9 |

# Comments

- On going work: we expect a significant improving of the performances, with a reasonable hardware cost

- In 2016, we were invited to present this work during 2 workshops in Paris (70 attendees) and Clermont Ferrand (55 attendees)
  - Proposals to collaborate and to provide us computing and storage facilities, other invitations

- Presentations and discussions at IVOA Sydney, Cape Town and Trieste

# "DevOps" at CDS

- The use of components (light virtualization) is becoming one of the most popular technologies with Docker as flagship

- We use it to deploy Apache Spark on the clusters

- Study of possible applying to CDS services like VizieR for the mirroring maintenance

- A major topic is to implement a prototype allowing a user "to move his code to the data" through components

A. Schaaff, F.-X. Pineau, G. Landais, L. Michel
Interns: Paul Trehiou (UTBM)

# Binding Database Metadata with Scientific Papers

- Objectives
  - Providing database users with a relevant description of the exposed data.
  - Helping users to locate resources described or mentioned in a text corpus
  - Exploring possible ways to facilitate the job of the documentalists
- Principle
  - Bringing together a text corpus related to the service
  - Indexing the corpus with a text search engine
  - Integrating this text search facility inside this service
- Foreseen Applications
  - XMM-Newton Catalogue: providing users with a scientific description of the catalogue columns
  - Vizier: searching catalogues by querying the corpus of README files as a whole

L. Michel (SSC XMM-Newton Strasbourg) in collaboration with the CDS (engineers and documentalists)
Intern: Sinan Acar (UTBM)

# Binding Database Metadata with Scientific Papers (2)



*Full text containing a selected fragments matching the requested keywords*

*Text fragments matching the requested keywords*

User Interface

Text Indexes

Indexing Engine

Text Corpus
- PDF
- ASCII
- HTML

Dictionary
Column names to Keywords

*Operation workflow*

*Data preparation workflow*

# R&D project tutoring at ENSIIE Strasbourg

- Not an internship, the work is done with the engineer school facilities, integrated in the semester
- Each student has 150 hours to spend on it
- Subjects around Immersive 3D Visualization, multi-Google Cardboard management from a laptop, interactions, etc.
- Students: Nicolas Buecher, Jonathan Chastenet, Mélody Deloffre, Jonathan Hurter, Pauline Kobersi, Raoul Le Perlier, Vincent Stébé, Damien Teodory

# And "after internship" short contract examples (2016)

- Thierry Lacoste: CDS logs to an operational version
- Felix Royer: improvement of XObsCoreFits
- Jérôme Desroziers: new (IVOA) VOTable Javascript parser (used in CDS Portal, AstroDeep, and soon in Aladin Lite), presented at IVOA Trieste
- Joris Vigneron (intern 2015): cleaning in Simbad database

# Future investigation plans 2018-2020

- Interfaces and interactions
  - High resolution (4K and perhaps more...) screens will be common, merging between (OS, multitouch) smartphone / tablet / laptop / desktop is on going and will probably be done during the period
  - Other kinds of interactions are emerging (voice, gesture, ...)
  - Balance between standalone apps / Web apps is moving
  - Continuous R&D in this frame is crucial

# Future investigation plans 2018-2020 (2)

- "Big Data"
  - Providing the data: a continuous R&D effort is needed to provide an added value on the access mechanisms to the data which include both the organisation of the data, the metadata and the technologies
  - Providing the tools: to access and explore all the CDS data (and external data in the context of the interoperability)
    - Knowledge databases, machine learning, etc.
  - Log Mining: capitalize on the centralization of the CDS logs to learn about the user workflow to improve the integration of the services

# Future investigation plans 2018-2020 (3)

- Immersive 3D-Visualization
  - Oculus Rift & similar devices: take into account the 2 previous years (after the public release of these devices) and continue (or not) the studies (and implementations) done previously
  - Continuous R&D around this topic
- Clouds
  - Following previous experiments (HiPS in the clouds), identify and implement in the cloud services (probably mirrors) of the CDS

# Future investigation plans 2018-2020 (4)

- Social networks
  - Probably one of the main evolution of the services
  - Deepest involving of the users at several levels
  - Astronomers, developers, documentalists: the fourth part of the team could be build in this frame
- Connected objects
  - A large variety of devices around 2020
  - Could be seen as gadgets in astronomy
  - (But) Probably many use cases (alerts, news, etc.) which are not yet well identified in the Science field

# Conclusion

- **Large spectrum** from the system/hardware to the Data/ Text mining
- The R&D activity is a way to update and improve the services... and the skills of the permanent staff
  - New technologies
  - Presentation of the work at the end of the internships
  - Etc.
- It is also a "transversal" activity making people work together
- It is not only technical...
- Special thanks to Sandrine Langenbacher, Thomas Keller and Jean-Yves Hangouet