

Implementing, with Python and PostgreSQL, Virtual Observatory services for publishing survey data from the Observatorio Astrofísico de Javalambre

Javier Hernandez¹

Tamara Civera¹

¹*CEFCA - Centro de Estudios de Física del Cosmos de Aragón, Teruel, Aragón, Spain*

Abstract. The Centro de Estudios de Física del Cosmos de Aragón (CEFCA) is carrying out from the Observatorio Astrofísico de Javalambre (OAJ, Teruel, Spain) two large area multiband photometric sky surveys, J-PLUS and J-PAS, covering the entire optical spectrum using narrow and broad band filters. As an effort to make public the data, we offer Virtual Observatory (VO) compliant services to make more versatile the access to the data through the multiple VO compliant existent tools. For example, catalogues are offered through TAP protocol and images can be searched and downloaded using SIAP protocol. This contribution summarizes why we decided to make our own implementation, the process followed to choose which services to implement according to the kind of data generated by the survey, why Python and PostgreSQL were chosen and, finally, the lessons learned.

1. J-PLUS and J-PAS Data Publication

As large area multiband photometric sky surveys, J-PLUS¹ and J-PAS² have images, catalogue data, and photo-redshift computations as its main output artifacts. The final data release will contain hundreds of thousands images (about 4000 pointings in 59 filters for J-PAS) with tens of thousands of objects detected in each image. Besides single image catalogues, the main catalogue is obtained detecting objects in a reference filter image (R filter) and measuring the object's fluxes in all of the other filters (using *Sextractor dual* mode), the result is a *photo-spectra* of each object detected in the reference filter.

The obvious decision to publish survey data was to create a custom web portal (Civera 2020), offering advanced tools for data search, visualization, and download. The portal includes a powerful sky navigator (where we can select an object and see its properties and photo-spectra), image search, cone search, object list search, and image catalogue download. But we realized that the tools were not enough appropriate for bulk data processing, automation, and user customization (for example catalogue files

¹<http://www.j-plus.es/>

²<http://j-pas.org/>

just include a fraction of the data, because including all the properties will bloat the files). And also did not offer a standardised way of accessing the data.

Due to the success of IVOA³ Virtual Observatory standards, and availability of tools supporting them, we decided to offer support for some of the services so astronomers can use existing tools and create scripts to operate with our data. There were already existing open source software to publish data as IVOA services, like the offered by GAVO⁴ and the Spanish VO⁵, but they did not fulfil all our needs. For example we wanted integration with our portal, secure access for certain catalogue versions, and support for a few custom extensions like array data types in TAP. We saw that IVOA services are web based, and in general they are simple, so we decided to create our own implementation of the few services needed, and integrate them in our portal.

Our database included information of images, object catalogues, and computed redshifts. We decided to implement on top of it the Simple Image Access (SIA) service for image search, that returns links to our download service of full fits images, *cutouts*, or colour images. We also included Cone Search for each object table that included indexed coordinates, and Table Access Protocol (TAP) to allow execute ADQL queries on the tables storing information of images, filters, objects and different photo-redshifts computed using several tools.

In our web pages we included also certain IVOA functionality, we included support for Simple Application Messaging Protocol (SAMP) so for example the image search page include a button that can send a found image to a open window of the Aladin⁶ tool (which supports SAMP). We also included for users that do not want to use a external tools a web page for executing TAP asynchronous queries.

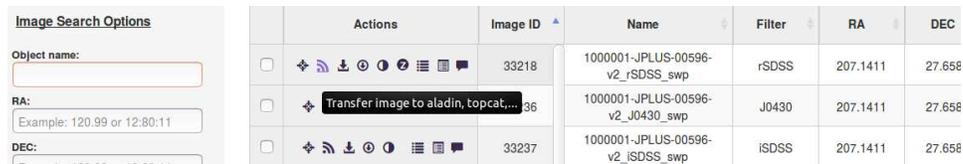


Figure 1. SAMP image transfer from web portal.

2. Implementation Tools

For the implementation the main parts to address were the storage for the catalogue data and the engine for web applications.

2.1. The Database Layer

J-PLUS and J-PAS Catalogue final data is expected to occupy several terabytes so one our priorities was to have a solid database server, known to manage terabyte sized

³<http://www.ivoa.net/>

⁴<https://www.g-vo.org/>

⁵<http://svo.cab.inta-csic.es/main/index.php>

⁶<http://aladin.u-strasbg.fr/aladin.gml>

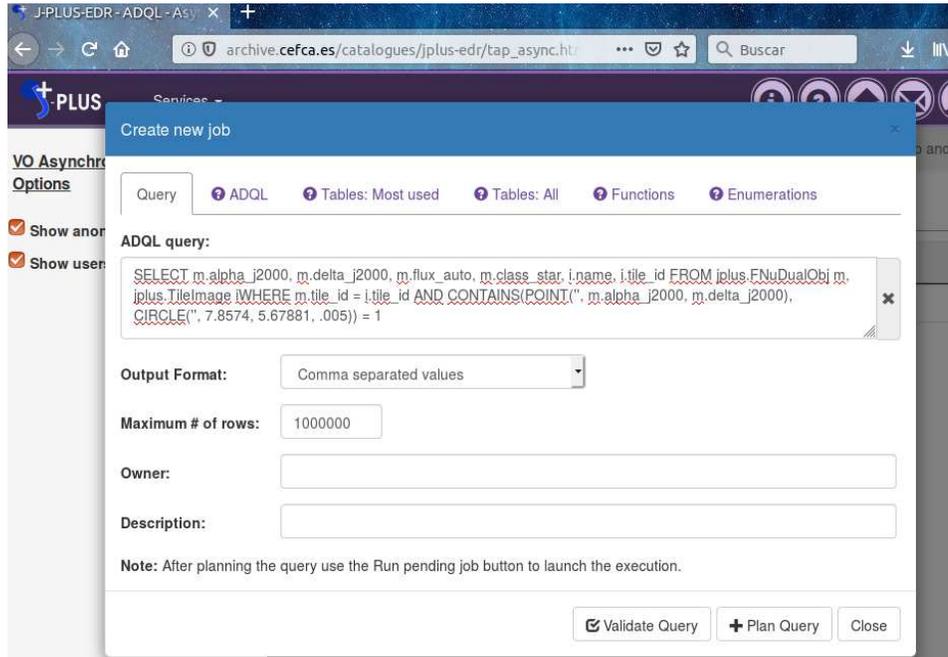


Figure 2. Asynchronous TAP page.

databases, with also rich functionality like possibility to add new functions and new data types. We also wanted a Open Source Relational database so we decided to use PostgreSQL. PostgreSQL has another very important feature for us, support for *array* data types. With flux and error measures for about 16 apertures in 59 filters that give us, just for photometry and its error, about 1888 columns of data in the object table, a nightmare. We found that using a column of array type and assigning a position for each filter in the array we simplified the problem a lot, now users can obtain easily the whole photo-spectra or using an index extract certain band.

The ability to add new data types and complex functions in different program languages was also very important for the implementation of TAP that includes a lot of math and geometric functions.

For cone search, and the use of similar functions in TAP, having a positional index in the database is very important. Although PostgreSQL has some support for spatial indexes, we decided to use Healpix⁷ for indexing object coordinates. Healpix indexes are integer values that are very compact to store in the database and allows the most reliable and efficient index type in PostgreSQL (B-Tree), additionally it allow users other uses like can be the creation of property maps (example density maps).

2.2. The Web Layer

The Python programming language was in use in the house for the reduction pipelines, so we evaluated its suitability for creating web applications. With a little research we

⁷<https://healpix.jpl.nasa.gov/>

found that the Python WSGI specification with a lot of supporting servers (Apache mod_wsgi, Unicorn, uWSGI, ...) and frameworks (Pyramid, Flask, Django, ...) gave us an excellent support for the task. And also we knew the large amount of astronomy, database access, and miscellaneous libraries available in Python.

We opted to use the Pyramid framework that allows to easily map a Python function to a web URL, to extract web parameters and access to url segments. Also the framework allows different forms to generate the response, like using templates or returning a stored file.

3. Lessons learned

Regarding TAP, our initial idea was to pass directly the queries to the database, but for security reasons and SQL dialect translation, we opted to create a Python parser which needed a significant effort but offered interesting opportunities like table rename and hiding, inserting *limit* expression to set the maximum output rows, or include what we called enumerations that allow to use names when referring the array position of filters in a magnitude/flux array column (for example 'jplus::rSDSS' for the 'r' filter position). Also supporting TAP we found that implementing the geometric functions is very complex due to the rich functionality that is defined, so at the moment only partial support for that functions exists (extra work is needed to support all coordinate frames or all figures operations).

IVOA centres on public data so security access is not standardised. We have some catalogue versions with limited access for a certain period of time so we need registration and authentication, but some VO tools do not offer support for authentication or it is not documented. We finally achieved to support authentication for some tools like Topcat using Basic HTTP authentication.

Performance is always something that at some point you have improve. The large number of objects in the catalogue databases makes that performance of functions in the database very important. We initially implemented in Python the needed functions for ADQL but later we moved to a C implementation because it is ten times faster. A TAP query can take some minutes to execute so executing concurrently them is mandatory. Python threads was our first election, but Python threads have some blocking issues. Fortunately Python *multiprocessing* package has a similar API and lacks that locking problems, so migration was quick.

Acknowledgments. Funding: Fondo de Inversiones de Teruel, PGC2018-097585-B-C21 (MCIU/AEI/FEDER), Aragon Reseach Group E16_17R.

References

Civera, T. 2020, in ADASS XXIX, edited by TBD (San Francisco: ASP), vol. TBD of ASP Conf. Ser., TBD