Adam Ginsburg, Brigitta Sipocz, Brett M. Morris, Thomas P. Robitaille, Leo P. Singer, 4,5 Christoph Deil, Henrik Norman, 7,8 Magnus V. Persson, 9 and the Astroquery collaboration, a subset of the astropy collaboration

THE ASTROQUERY COLLABORATION, A SUBSET OF THE ASTROLT COLLABORATION

ABSTRACT

Astroquery is a collection of tools for requesting data from databases hosted on the internet, particularly those with web pages but without formal application program interfaces (APIs). These tools are based on the Python requests module, which is used to make HTTP requests, and astropy, which provides most of the data parsing functionality. Astroquery has received significant contributions from the broader astronomical community, including several significant contributions from telescope archives. ^{a)}

Corresponding author: Adam Ginsburg aginsbur@nrao.edu; adam.g.ginsburg@gmail.com

¹ Jansky fellow of the National Radio Astronomy Observatory, 1003 Lopezville Rd, Socorro, NM 87801 USA

² Astronomy Department, University of Washington, Seattle, WA 98195, USA

³ Aperio Software Ltd., Headingley Enterprise and Arts Centre, Bennett Road, Leeds, LS6 3HN, United Kingdom

⁴ Astroparticle Physics Laboratory, NASA Goddard Space Flight Center, Mail Code 661, Greenbelt, MD 20771, USA

⁵ Joint Space-Science Institute, University of Maryland, College Park, MD 20742, USA

⁶ Max-Planck-Institut für Kernphysik, Heidelberg, Germany

⁷ Winter Way, Uppsala, Sweden

⁸ESAC Science Data Centre, European Space Agency, Madrid, Spain

⁹ Department of Space, Earth and Environment, Chalmers University of Technology, Onsala Space Observatory, 439 92, Onsala, Sweden

a) The repository associated with this paper is: https://github.com/adamginsburg/astroquery-paper

1. INTRODUCTION

Sharing data is a critical component of astronomical research. Astronomy has historically been a leading field in data sharing, motivated at least in part by questions that cannot be answered with single instruments. In the past few decades, blind surveys have played a huge role in advancing our understanding of the universe.

Data sharing has taken on a variety of forms. The most prominent are the major observatory archives: MAST, NOAO, ESO, IPAC, CADC, CDS (hosting Vizier and SIMBAD), NRAO, CXC, HEASARC, and ESA are the main organizations hosting raw and processed data from ground and space based telescopes. These data archives also serve as the primary means for serving data to users when the data are taken in queue mode, i.e., when the data are taken while the observer is not on-site.

In addition to observatories and telescopes, individual surveys often share their full data sets. In some cases, these data sets are shared via the observatory that acquired them - for example, the all-sky data acquired with Planck, WMAP, and COBE delivered a variety of data products as part of the mission. Other surveys, particularly ground-based surveys, serve their own data. Examples include SDSS, UKIDSS, and likely many more in the near future.

Individual teams and small groups will often share their data. These services do not follow any particular standard and can be widely varied in the type and amount of data shared. Sometimes these data are shared via the archive systems (e.g., IRSA at IPAC hosts many individual survey data sets), while others use their own web hosting systems (e.g., MAGPIS).

Finally, there are other data types relevant to astronomy that are not served by the typical astronomical databases. Examples include molecular and atomic properties, such as those provided by Splatalogue and NIST.

Astroquery arose from a desire to access these databases from the command line in a scriptable fashion. Script-based data access provides astronomers with the ability to make reproducible analysis scripts in which the data are acquired and processed into scientifically relevant results with minimal overhead.

In this paper, we provide an overview of the astroquery package. Section 2 describes the basic layout of the software and the shared API concept underlying all modules. Section 3 describes the development model.

2. THE SOFTWARE

Astroquery consists of a collection of modules that mostly share a similar interface, but are meant to be used independently. They are primarily based on a common framework that uses the Python requests package to perform HTTP requests to communicate with web services.

For new development, there is a template that lays out the basic framework of any new module. All modules are based on having a single core class that will have some number of query_* methods. The most common query methods are query_region, which usually provide a "cone search" functionality, i.e., they search for data within a circularly symmetric region projected on the sky.

An example using the SIMBAD interface is shown below (see http://astroquery.readthedocs.io/en/latest/simbad/simbad.html):

from astroquery.simbad import Simbad
result_table = Simbad.query_region("m81")

In this example, Simbad is an instance of the astroquery.simbad.St class. The returned result, stored in the variable result_table, is an astropy table.

While there is a common suggested API described in the template module, individual packages are not required to support this API because, for some, it is not possible. For example, the atomic and molecular databases refer to physical data that is not related to positions on the sky and therefore their Astroquery modules cannot include query_region methods.

2.1. The API

The common API has a few features defined in the template module. Each service is expected to provide the following interfaces, assuming they are applicable:

- query_region A function that accepts an astropy SkyCoord object representing a point on the sky plus a specification of the radius around which to search. The returned object is an astropy Table.
- query_object A function that accepts the name of an object. This method relies on the service to resolve the object name. The returned object is an astropy Table.
- get_images For services that provide image data, this function accepts an astropy SkyCoord and a radius to search for data that cover the specified target. The returned object should be a list of astropy.io.fits.HDUList objects.

Beyond these basic functions, there is a series of functions with the same names, but with the additional suffix _async. The query_*_async functions return a requests.Response object from the accessed website or a list of such objects, providing developers with the ability to access the data in a stream or access only the response metadata. The <code>get_images_async</code> method returns <code>FileContainer</code> objects that similarly provide 'lazy' access to the data, but specifically for FITS files. Developers need only implement these <code>_async</code> functions because a wrapper tool exists to convert <code>_async</code> functions into their corresponding non-asynchronous versions.

2.2. Testing

Astroquery testing is somewhat different from most other packages in the Python ecosystem. While the tests are based on the astropy package-template and use pytest to run and check the outputs, the astroquery tests are split into remote and non-remote. The remote tests exactly replicate what a user would enter at the command line, but they are dependent on the stability of the remote services. Historically, we have found that it is quite rare for all of the Astroquery-supported services to be accessible simultaneously.

We therefore require that each module provide some tests that do not rely on having an internet connection. These tests rely on *monkeypatching* to replace the remote requests, instead using locally available files to test the query mechanisms and the data parsers. Monkeypatching in the context of pytest results in code that is generally more difficult to understand than typical Python code, but a set of tests independent of the remote services is necessary.

The non-remote tests are run as part of the continuous integration for the project with each commit. The remote tests are run as part of a regularly-scheduled cron job. Running the remote tests infrequently also helps reduce the burden on the remote services.

 $2.3.\ Other\ utilities$

There are several general-use utilities implemented as part of Astroquery, such as a bulk FITS file downloader and renamer and a download progressbar. There is also a schema system implemented to allow user-side parameter validation. So far, at the time of writing, this tool is only implemented in the ESO and Vizier modules, but it could be expanded to other modules to reduce the number of doomed-to-fail queries sent through Astroquery.

3. THE DEVELOPMENT MODEL

Astroquery is an astropy affiliated package and is a core part of the astropy ecosystem (Astropy Collaboration et al. 2013). It is a standalone project and will always remain independent of the astropy core package.

Astroquery has received contributions from 53 people as of June 2017. While the primary maintenance burden is shouldered by 2-3 people at any given time, most individual modules have been implemented independently by interested volunteers.

Some contributions have come as direct institutional support. The ESA GAIA and ESASky modules were provided by developers working for ESA. Meanwhile, the MAST and VO Cone Search query tools were added by developers at STSCI, with the latter moved over from astropy.vo.

Astroquery has received support from the Google Summer of Code program, with two students (co-authors Madhura Parikh and Simon Liedtke) from 2013-2017.

Anyone can contribute to Astroquery. The maintainers are committed to helping developers make new modules that meet the requirements of Astroquery.

4. DOCUMENTATION AND REFERENCES

Several authors have independently described how to use various **astroquery** modules, which is a helpful practice we encourage.

- https://www.cosmosim.org/cms/news/cosmosim-package-
- https://arxiv.org/abs/1408.7026

REFERENCES

Astropy Collaboration, Robitaille, T. P., Tollerud, E. J., et al. 2013, A&A, 558, A33