**Science Platforms: Interoperability and other Missing Pieces**

**Gregory Dubois-Felsmann**
**Caltech/IPAC-LSST**

**7 December 2018**

# Comments in the session

- Being able to transport a notebook and associated software from one to another, and use it on the different datasets available at different sites

- Seamlessly performing cross-mission queries

- Need for metadata and documentation for the platforms with which we're not familiar.  Need to understand the quality and resolution of the data

- Seamless access to internationally distributed datasets; require information about the science quality of the data you are getting back

- Interoperable authentication - no need for multiple accounts

- Common batch job submission mechanism

- Need a lot of different datasets; avoid as much as possible the need for users to have to understand the meaning of the tables & columns from different datasets; interoperable data models (part of the traditional IVOA remit) -- metadata, UCD's

# Comments in the session (2)

- Worried about imposing too many constraints on data services, archives, etc.; too much normalization limits flexibility (and discovery)

- Discussion of CAOM2, UCDs, etc.; standardization of observational metadata vs. catalogs

- Interoperability for science platforms should be seen to be growing out of the existing body of work in the IVOA, and not a wholly new, huge, VAO-like enterprise

- Container-level interoperability

- If I move a container to a different site, do the local data services at the destination speak the same language?

- Could users define a personal environment that could be replicated at multiple sites?

- Very important to communicate exactly what (processing?) was successful, what can be trusted, what work units were completed?

# What Interoperability?

Clearly there are many flavors of interoperability that are relevant to Science Platforms

- Interoperability of data services over the network (classic IVOA topic)
- Interoperability of source code that accesses data
  - I want to use the same (or minimally modified) code to access (say) imaging survey data from many sources
  - I want to use the same (or minimally modified) code to access data whether it's local or remote (orthogonalize user knowledge of how to code and user knowledge of where the data are)
- Interoperability of scientifically similar data models from multiple sources
- Interoperability of user code at the level of a "unit of analysis"
  - Notebooks
  - Containers
  - Software packages (3rd-party or user's personal software)
- Interoperability of resource access (e.g., to parallel processing, batch systems)

# What Interoperability?

**Clearly there are many flavors of interoperability that are relevant to Science Platforms**

- Interoperability of data services over the network (classic IVOA topic)
- Interoperability of source code that accesses data
  - I want to use the same (or minimally modified) code to access (say) imaging survey data from many sources
  - I want to use the same (or minimally modified) code to access data whether it's local or remote (orthogonalize user knowledge of how to code and user knowledge of where the data are)
- Interoperability of scientifically similar data models from multiple sources
- Interoperability of user code at the level of a "unit of analysis"
  - Notebooks
  - Containers
  - Software packages (3[rd]-party or user's personal software)
- Interoperability of resource access (e.g., to parallel processing, batch systems)

# Evolution of a key set of VO Standards

| Standard | Purpose | Adopted | Proposed |
|---|---|---|---|
| TAP | Catalog data access, table metadata access | 1.0, 2010 | 1.1, 2018 |
| ADQL | Catalog query language | 2.0, 2008 | 2.1, 2018 |
| ObsCore/ObsTAP | Basic observation data model and service definition | 1.1, 2017 | |
| SIA | Simplified image metadata queries | 2.0, 2015 | |
| SODA | On-demand operations like image cutouts | 1.0, 2017 | |
| HiPS | All-sky images and catalogs | 1.0, 2017 | |
| MOC | All-sky Boolean coverage maps | 1.0, 2014 | |
| VOSpace | Remote access to file-oriented data | 2.1, 2018 | |
| **Supporting Standards** | | | |
| VOTable | Tabular data transport with metadata | 1.3, 2013 | |
| DataLink | Two protocols supporting connecting datasets | 1.0, 2015 | |
| DALI | Framework for service definitions | 1.1, 2017 | |
| VOSI | Service self-description standards | 1.1, 2017 | |
| UWS | Protocol for interacting with asynchronous services | 1.1, 2016 | |
| VO-DML | Data model language | 1.0, 2017 | |
| UCD1+ | Controlled vocabulary for semantics of individual data items | 1.3, 2018 | |

# Unfinished business in the existing IVOA ambit

- Fulfilling the promise of the registry and of interoperable discoverability
  - Universal and uniform access to data discovery
  - In practice the quality of registry data is still highly variable, as is the standards-compliance of the registered services

- Interoperable data models that capture a large fraction of the public astronomical data space

- Fully exploiting the standards we have

# Unfinished business in the existing IVOA ambit

- Fulfilling the promise of the registry and of interoperable discoverability

- Interoperable data models that capture a large fraction of the public astronomical data space
  - ObsCore and CAOM2 (non-IVOA but a de-facto standard) do a good job for observation metadata and are being used more and more widely
    - ObsTAP and SIAv2 provide protocols for accessing ObsCore-structured data
    - TAP services can support access to CAOM2, but query patterns are complex and difficult for new users to reproduce (see language bindings discussion below)
  - Standardization of a core data model for object/source/detection catalogs has not gotten above threshold for wide usability
    - VO-DML, approved this year, provides a language for this
    - Still trying to work out standards for annotating data with data models
    - No widely acceptable core Source data model yet exists
    - IMO: Need to focus on building this up incrementally (capture 75% of the value with 25% of the effort?)

# Unfinished business in the existing IVOA ambit

− Fulfilling the promise of the registry and of interoperable discoverability

− Interoperable data models that capture a large fraction of the public astronomical data space

− Fully exploiting the standards we have
  - High-quality, regularly validated services
  - Clean registries
  - Extensive use of DataLink to allow navigating from data to related data

## Client-side access to IVOA services

- Greatest focus has been on data-center to tool-developer relationship

- The development of language bindings to IVOA network protocols and data serializations, or even their standardization, *are not part of the IVOA's remit*
  - The IVOA has no mandate to formally favor any particular language (and doing so would be a form of anti-interoperability)

# Client-side access to IVOA services

– Science-user-centric language bindings have been around for years, but not a focus of effort until recently

  – Tension:
  the people who are good at designing and implementing language bindings appropriate to a given language (e.g., Python)

  are generally not the same as

  the people with in-depth understanding of the standards they may be implementing against, and even less so the ones understanding the "big picture" of the architecture: how the IVOA standards are envisioned to work together

  – A language binding that too-explicitly exposes the existence of the underlying protocols is not necessarily what the community wants
    • Astronomers would naturally like to think about domain objects: observations, sources, spectra, and not about formal protocols

# Python existing practice: Compare Astroquery and PyVO

- Astroquery is organized around a template for query interfaces, which *may* be used to implement modules that provide access to astronomical data services
  - Interoperability is at the Python level, assuming the template is followed, and makes no assumptions about standardization of the underlying data services
  - Availability of standardized (IVOA) data services does allow many Astroquery modules to share implementation code

- PyVO is explicitly oriented around the existence of the IVOA protocols and provides 1:1 bindings to them
  - Using PyVO requires understanding the existence and purpose of the protocols

- Vision: build out an Astroquery-like interface that hides protocols but takes advantage underneath of high-quality, robust bindings to the protocols
  - Especially: one that encourages the development of a DataLink-based ecosystem

## The promise of data modeling

- Progress in the availability and use of rigorous common data models for (subsets of) published astronomical datasets (for observations, catalogs) *enables exposing that data through common language-specific models* (e.g., Python classes)

- Baby steps: recreating Astropy "SkyCoord" objects in the Python bindings to IVOA services that return spherical coordinates
  - How far can this go?
  - Is it in conflict with efficient in-memory data models for tabular data?
  - How much does it help astronomers to have rows in a source database representable as Python objects?

# What Interoperability?

## Clearly there are many flavors of interoperability that are relevant to Science Platforms

- Interoperability of data services over the network (classic IVOA topic)
- Interoperability of source code that accesses data
  - I want to use the same (or minimally modified) code to access (say) imaging survey data from many sources
  - I want to use the same (or minimally modified) code to access data whether it's local or remote (orthogonalize user knowledge of how to code and user knowledge of where the data are)
- Interoperability of scientifically similar data models from multiple sources
- Interoperability of user code at the level of a "unit of analysis"
  - Notebooks
  - Containers
  - Software packages (3rd-party or user's personal software)
- Interoperability of resource access (e.g., to parallel processing, batch systems)

# Open-source data analysis codes

- Open-source software efforts are producing robust, mature, and immediately applicable layered libraries that provide out-of-the box solutions for many computational and visualization tasks.

- They are becoming a common language that is known to young people entering our disciplines.

# Open-source infrastructure tooling

− Containerization frameworks and container-orchestration frameworks such as Docker and Kubernetes are greatly simplifying system administration and application deployment.

# Open-source server-side interactive application delivery

− The Web-based notebook interface (here, JupyterLab) is providing a straightforward means of delivering
interactive,
server-side,
next-to-the-data
computing to
our users.

# So what's missing?

In implementing the "Science Platform" vision of offering computing resources close to the data, the following pieces seem not to be as mature or as seamlessly easy for our users to apply:

- Frameworks for efficient multi-tenant application of parallel tools such as Dask
- Cross-site APIs for managing batch jobs and workflows of batch jobs
- Finer-grained transportability of packaged user analysis code across data centers
  - Including from agency-managed assets to commercial cloud providers – support multiple models!
- Tools for automatically matching code to data across multiple data centers
  - Seems to require additional metadata standardization
  - Need widely-adopted solutions to security concerns about container frameworks
- Interoperable authentication and authorization
- Frameworks for reusability, reproducibility, human-friendly version control, collaborative editing, and citeability that work smoothly with the now-popular notebook-oriented environments
  - Be aware of developing concerns/criticism of a too-pervasive shift to notebooks without this

There is a need for R&D, development of experience bases and lessons-learned, and finally (where appropriate) standardization in these areas.