



STScI | SPACE TELESCOPE
SCIENCE INSTITUTE

EXPANDING THE FRONTIERS OF SPACE
ASTRONOMY

The Roman SOC Data Management System

Henry Ferguson STScI

Science Operations and the Data Management System

STScI Roman Science Operation Center responsibilities include:

- Planning & scheduling all observations
- Calibration and support of the Wide Field imaging
- The archive (MAST) for all mission data
 - Most Roman science will be archival due to the survey nature of the mission

NASA Astrophysics *Big Data*:

- Data accumulated per week likely to be $\gg 100x$ *Hubble*
- Both catalogs and pixel-level data sets provide unique science opportunities
- Downloading and processing exceeds resources typically available



Barbara A.
MIKULSKI ARCHIVE FOR
SPACE TELESCOPES

Science data products from multiple mission partners

- Calibrated and mosaiced images, extracted spectra, catalogs, etc.
- Staged in the cloud and co-located with significant computational resources
- Open source, modular imaging pipeline facilitates custom reprocessing

Data storage & processing

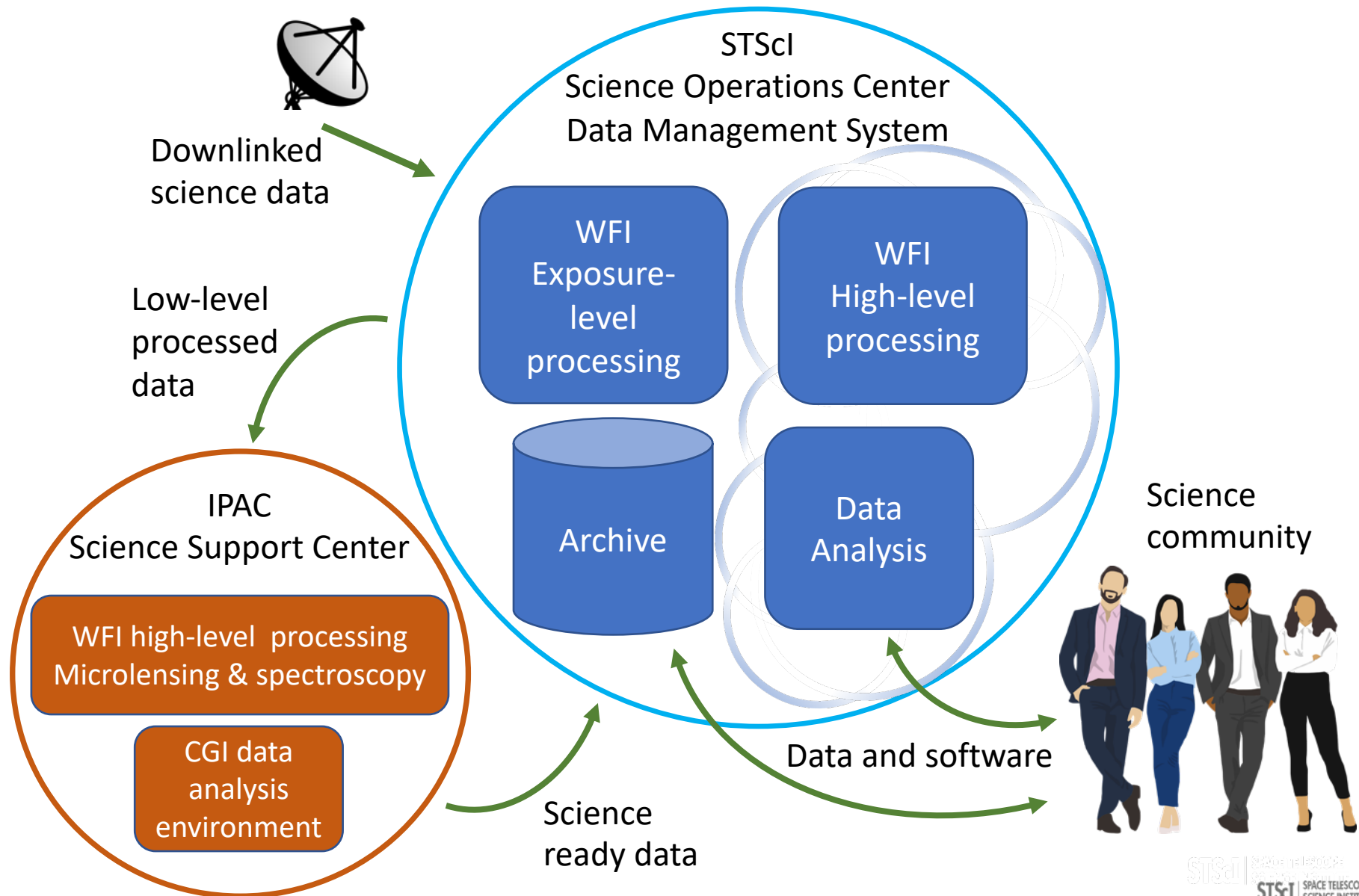
- Cloud-based high-level data processing brings software to the data
- Jupyter Lab environments ease access, sharing and repeatability
- Software environment for the community in sync with mission data processing



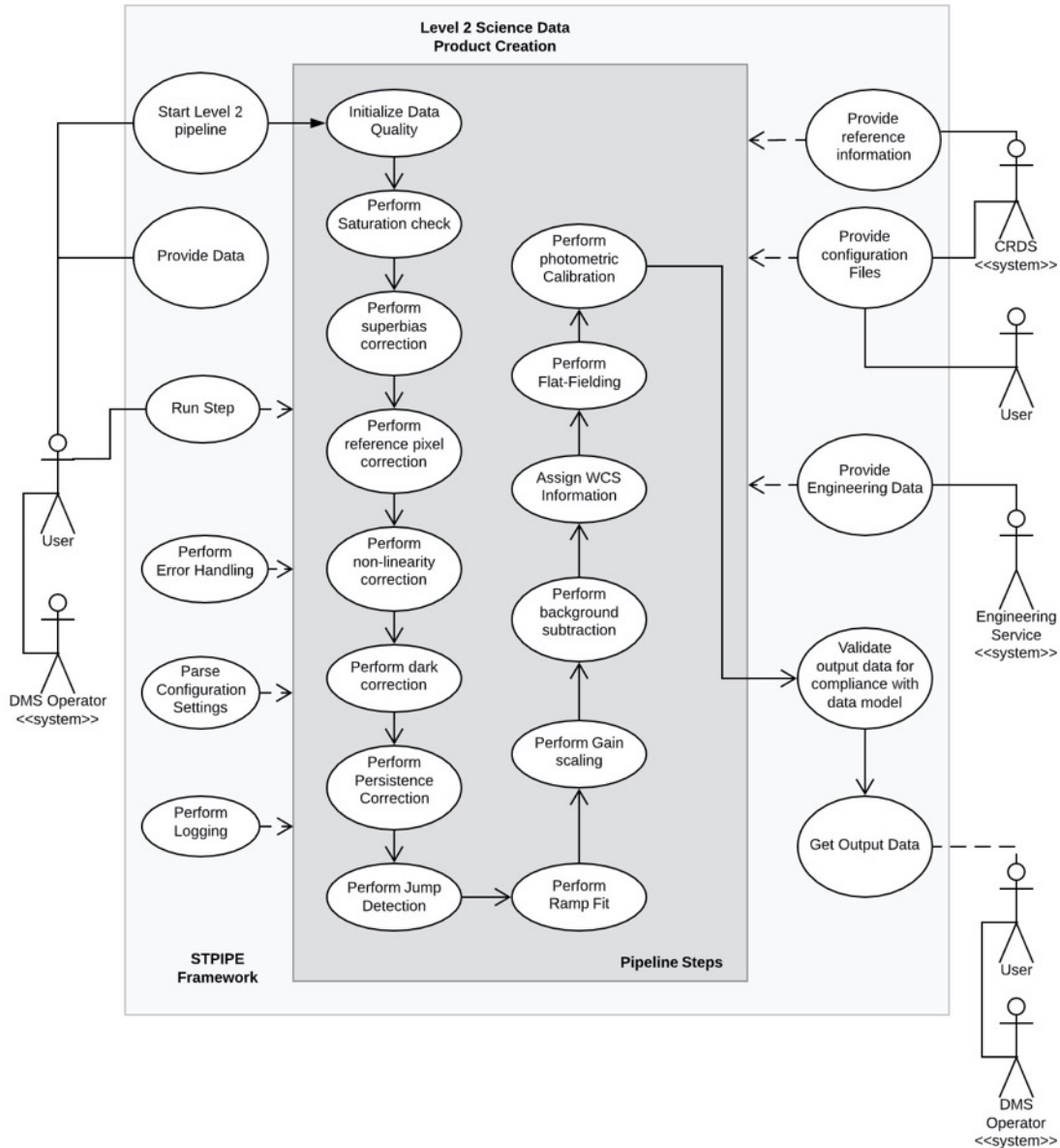


- System approved at PDR had the following attributes:
 1. The JWST science calibration pipeline with minor adaptations
 2. High-level WFI processing using software provided by external science teams
 3. A science platform (HLPP) that allows users to interact with the data and high-level processing software in the cloud
 4. An Archive with HST/JWST/MAST like functionalities
 - Including science data from the SOC, the SSC and high-level community products.
 - Archiving selected WFI and CGI ground test data
 - Storage of all Roman mission data products

Roman Data Management



Exposure Level Processing Flow





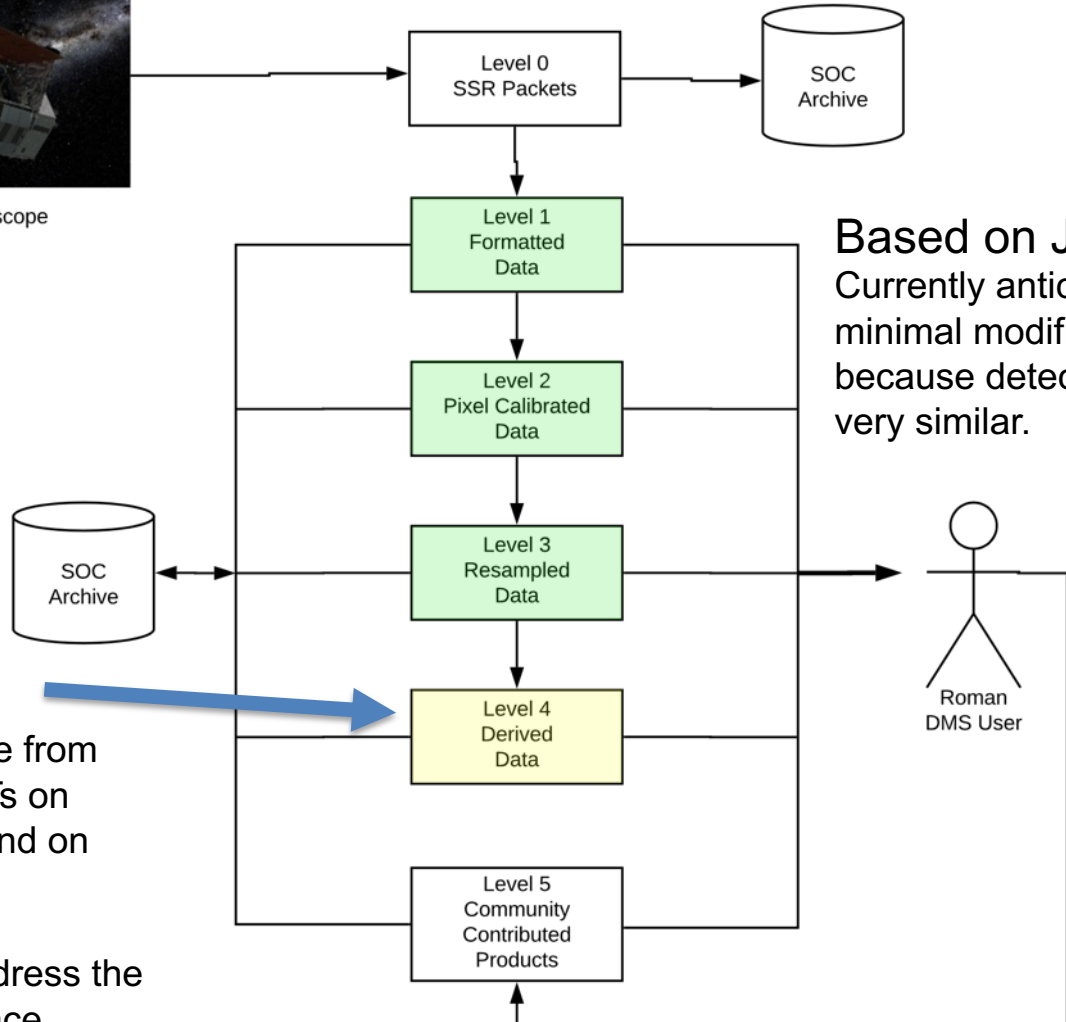
- Roman science data are public
- Users will be able to retrieve science data from MAST.
 - including data from the SOC, the SSC as well as I&T data and high-level community products.
- Expect archive services to evolve
- Currently incorporating Jupyter analysis+visualization tools into the archive for JWST.
 - Improving access to high-level products with services like z.mast and exo.mast.
- Higher-level products (level 2 and beyond) will be available in the cloud as well
 - SOC is currently scoped for cloud hosting of SOC data, not SSC, CGI or community products (although they will be in MAST)



- Decided at mission level
 - Division of responsibilities between MOC, SOC and SSC
- Past SOC studies
 - Cloud service providers \Rightarrow AWS
 - Mission Data Formats \Rightarrow ASDF
 - Science Platform \Rightarrow JupyterHub
- Underway
 - Database technologies
 - Extent of cloud integration for MOC, SOC and SSC
 - Evolution of cloud vs. on-premises for all missions
 - Management & policies for community use of the cloud resources (later in this talk)



Roman Space Telescope



Based on JWST
Currently anticipate minimal modifications because detectors are very similar.

Now in scope for SOC:
Original plan was to integrate software from SITs. New plan is to work with the SITs on features and algorithms, but not depend on SITs for software deliveries.

Worked closely with the Project to address the highest priority risks for meeting science requirements.

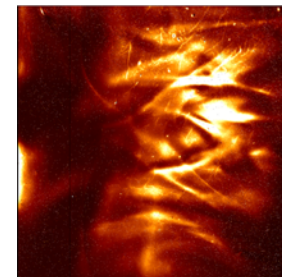
- Simulations

- Package (based on JWST Mirage) to simulate WFI level-1 data
 - Including the most important instrument signatures
 - Hugely beneficial for testing pipeline instrument-signature removal
- Idealized simulations for Level 2
 - Enables artificial source injection in the pipeline
 - This is critical for quantifying systematic uncertainties for many different science topics

- Point spread functions

- Empirical calibrations vs. position and time
- Queryable library
- Programmatic access via API

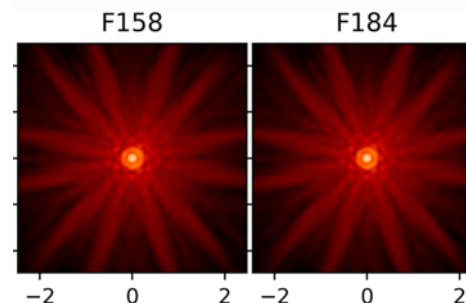
Instrument signatures



Idealized

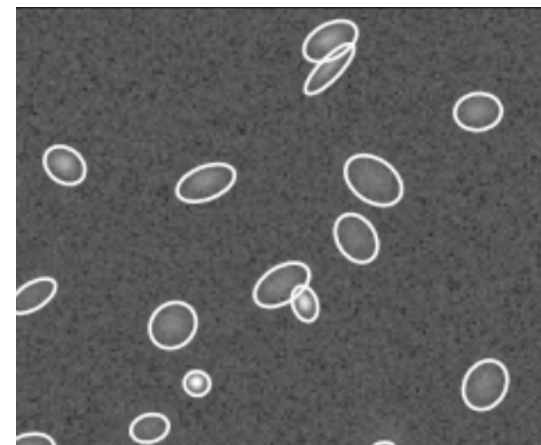
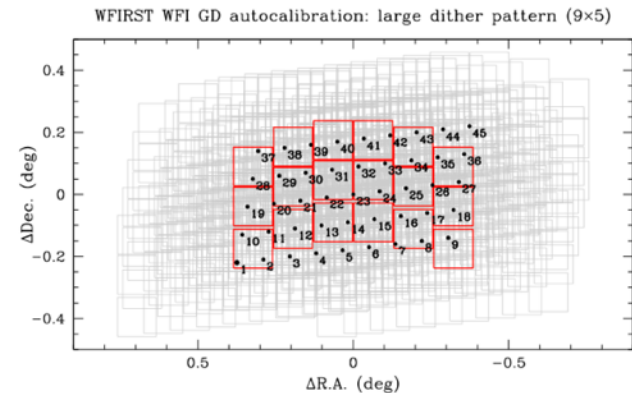


Point spread functions





- Astrometric calibration
 - Referenced to GAIA
 - Consistency between image metadata and catalogs
- Catalogs
 - Catalogs of static sources
 - Matched-PSF multi-band photometry
 - Including photometric redshifts (based only on Roman data)
 - Variable sources
 - Multiple epochs of catalogs
 - Catalogs of difference images
 - Spectroscopic extractions & redshifts (SSC)
 - Matched to sources in imaging catalog



Integrating high-level science data



- SSC responsible for microlensing data & spectroscopy
- SOC and SSC pipelines diverge at the highest levels, but:
 - Data will be distributed through common archive and (ideally) common formats with meta-data
 - HLSS and HLIS data will be integrated and matched per science requirements
 - SOC Science calibration pipeline will be public so the community can use SOC and SSC modules for other applications
 - The goal is for the archive to appear seamless to outside users



- Public data products contributed by the science community are likely to be widely used. Examples include:
 - Joint photometry with complementary data sets
 - Photometric redshifts that use complementary data sets
 - Value-added catalogs of derived properties (e.g. from SED fitting)
 - Hybrid spectroscopic and photometric catalogs
 - Survey-level calibrations
 - Improved astrometry & photometry after constraining for consistency across the full survey
 - Window functions, masks, PSF kernels, etc.
 - Transient-free template images
- Details & cadence to be defined through future community engagement and opportunities

Computing & Data Resource Management for the Science Platform

Work in Progress

Evolution

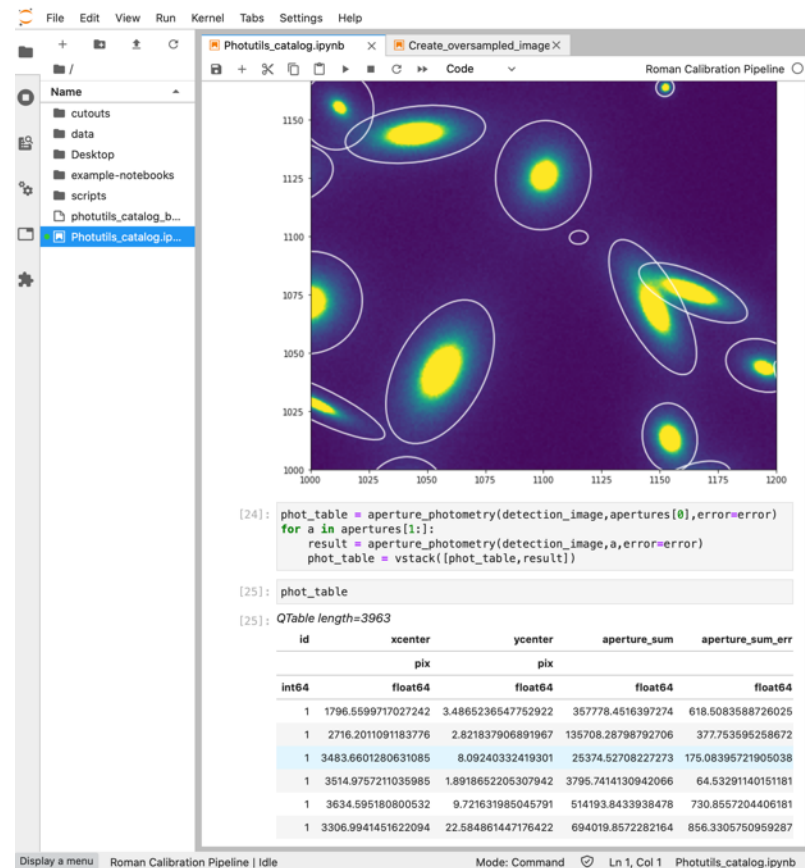
- Major projects are increasingly moving to the cloud
- Funding strategies could/should evolve accordingly
 - Traditional Hubble model is to enable hardware purchases & computing service fees in grants
 - Supercomputing proposals are typically separate from grants
 - NSF XSEDE program; NASA HEC program
 - Pros & cons:
 - May not be efficient to fund cloud allocations in individual small grants
 - Want to avoid double-jeopardy of separate proposals for computing vs. getting observing time or archival funding
 - Want to enable people with funds to be able to use them and work in the same environment
- Cross-institutional collaboration brings opportunities & challenges
 - How to manage access & allocations for non-US astronomers

Why the Cloud?

- **Putting both the computing and the science-ready data in the commercial cloud offers the following benefits:**
 - Convenient scalability for both data volume and computational demands
 - Flexible solutions for specific computing needs (e.g. GPUs or I/O optimized computing)
 - Lower total costs to NASA relative to multiple on-site installations
- **Benefits to the science users include:**
 - Efficient access to the data
 - Computing resources for exploratory work are available with no need to write a grant proposal
 - Local IT and software support costs are greatly diminished
 - Easier collaboration with astronomers across institutions
 - A powerful and stable science software environment



- Log in with your MyST account
- JupyterHub instance
 - Roman science calibration pipeline software installed and configured
 - Full Python + Astropy ecosystem installed and configured
 - Ability to install other packages and your own code
- Flexible, scalable architecture
 - Simple to add CPU & storage
 - High-throughput access to the data
 - Can scale up resources (e.g. GPUs or neural engines) as science needs & technology evolves



Concept for managing Roman platform



- The platform should provide a sufficient base level of resources
 - Most Roman grants wouldn't need a computing line item.
 - Projects needing exceptional resources could still apply for funds
 - Allows much more global optimization for science than case where funds are locked in small grants.
- Tier concept to support most users
 - Very lightweight process for getting access and increasing allocations
 - Will require periodic renewal
 - Will require rather frequent purging of stored data to control costs
- Looking to to enable ~1000 Roman papers/year
 - The Roman-data-intensive work; not necessarily all the computing
 - Not long-term archival storage of intermediate projects

Resources

- Priorities
 - Provide resources for work that needs to access large quantities of Roman data
 - Make it easy for the user community to get access, use & collaborate
 - Cost effectiveness & cost predictability
 - Biggest concern is persistent file-system storage (EFS)
 - Do not *require* use of the platform to do Roman science
- Not part of the concept
 - Providing resources for everything else
 - e.g., simulations & modeling, reducing complementary data
 - Don't plan to micromanage but will set quotas based on Roman data-analysis needs.
 - Providing a dedicated machine-learning environment
 - Could add this later if there is demand
 - Resources for high-school & undergraduate education
 - Using Roman data for coursework
 - Cost sharing with non-US stakeholders
 - CGI data analysis on the same platform
 - TBD
 - Co-location & support for analysis of SSC high-level products on the same platform
 - Co-location & support for community contributed products on the same platform

Relation to the Core Surveys



- Using WFI: weak lensing, supernovae and microlensing
- Have looked at “typical use cases” rather than core-survey workflows
 - Unclear if the teams will work in the cloud or other facilities
 - Depends on which teams are selected and what resources they may already have available
 - Expect the survey teams to do extensive survey-level processing beyond the scale of typical users
 - Infrastructure for this processing may or may not be best associated with this general platform
 - May need co-location with extensive simulations or other data sets
 - May make use of other facilities associated with survey teams
 - Even if in AWS, likely beneficial to optimize based on specific needs
- The concept presented here does not preclude very large “consortium” allocations, if needed
- Or separate AWS accounts, access to NASA HPC facilities, ...

Tier Concept

- Entry tier
 - Anyone with a MyST account
 - Examples: Filter the entire HLIS catalog a dozen times, make custom catalogs of few FPAs, registered to an external dataset, Extract & download 1000 cutouts
- Research tier
 - Pro-forma science justification & annual renewal
 - Examples: Custom catalogs on ~ 100 sq deg. Custom processing at the level of 1 minute per FPA for the entire HLIS. Thousands of catalog queries & cutouts.
- Consortium tier
 - Proposals & panel review (lightweight process)
 - Multiple users sharing a single allocation
 - Resource sub-allocations left to the consortium
 - Examples: Signal injection and re-run of pipeline for full HLIS

Summary of Policy Recommendations



- NASA funds the platform through STScI
 - Not via individual grants to PI, because this is simply not worth the administrative overhead of multiple transfers of funds
 - Projects needing exceptional resources could still apply for funds through ROSES process
 - Allow international access in all tiers
- Tier concept to support most users
 - Lightweight process for getting access and increasing allocations
 - Will require user agreements and periodic renewal
 - Regular migration and purging of data to control costs
- Scoped to support ~1000 Roman papers per year
 - Not all computing but stages where proximity of CPU to Roman data is beneficial
 - Not explicitly driven by core survey computational requirements

Big Data Discovery Tools

Courtesy of Josh Peek,
STScI Data Science Mission Office



What is in the archive plan already: MAST adapted for Roman

The SOC will reuse existing MAST functionalities. It will not provide new data mining or exploration functions specifically to assist the community in dealing with the challenges of the large Roman data volume, beyond what is provided by the HLPP ...

What is in MAST for Roman today?

- CASJobs — The first science platform, built by SDSS, based on SQL
- catalogs.mast — A new, lightweight catalog query webpage
- Discovery Portal — A observation-based search tool
- HLSP system — A robust system for ingesting community data

MyDB Local Only

Views

Tables

Functions

Procedures

Sort by... All selected...

| Rows | kB | Name |
|-----------|-----------|--------------------|
| 1,000 | 136 | ajp |
| 3,641,746 | 378,824 | ajpall |
| 1,000 | 136 | ajpXGaiaDR2 |
| 3,991,162 | 506,952 | ajpXGaiaDR2all |
| 1,912,000 | 273,288 | ajpXGaiaDR2allredc |
| 10 | 72 | MyTable |
| 3,641,746 | 378,824 | MyTable_0 |
| 364,673 | 29,832 | MyTable_AiIC |
| 318,097 | 63,688 | MyTable_M83 |
| 318,097 | 1,272,584 | MyTable_M83_all |
| 318,097 | 63,688 | MyTable_M83_W3 |
| 318,097 | 79,560 | MyTable_M83_W3_2 |
| 327,461 | 21,192 | MyTable_V1 |

MyTable_M83_W3

Contains

Notes

Tables

MatchR

float [9

Notes

No notes

New

Add Note

Pan-STARRS Catalog Search

The Panoramic Survey Telescope & Rapid Response System (Pan-STARRS or PS1) is a wide-field imaging facility developed at the University of Hawaii's Institute for Astronomy for a variety of scientific studies from the nearby to the very distant Universe. The PS1 catalog includes measurements in five filters (grizy) covering 30,000 square degrees of the sky north of declination -30 degrees, with typically ~12 epochs for each filter. This interface allows searches for the mean measurements and the deeper stacked measurements from images combining all the epochs. The DR2 release also includes the detection catalog containing all the multi-epoch observations.

Target Supply the central coordinates or target name.

Crossmatch a List of Targets Upload a CSV file.

[Choose CSV file](#) [Resolve & validate targets](#) [Target file help](#)

Crossmatch Search Radius Arcseconds Max = 3 Arcseconds

What to Search Select the catalog type and release to search.

Release PS1 DR1 PS1 DR2

Catalog Mean object Stacked object Forced mean object

[View Mean Search Table](#)

Display Columns Select the columns that will be displayed. **67 selected**

Search Conditions Select the rows that will be displayed.

[+ Add condition](#)

is greater than or equal to

is greater than or equal to

[Search Catalog](#)



Select a collection...

MAST Observations by Object Name or RA/Dec

[About Collections...](#)

[Upload Target List](#)

[My Download Basket: 0 files](#)

[Home Page](#) [MAST: m33](#)

and enter target:

[Search](#)

[Show Examples...](#) [Random Search](#) [Advanced Search](#)

[User Manual/Help](#) [Learn Feedback](#) [About This Site](#)

anonymous

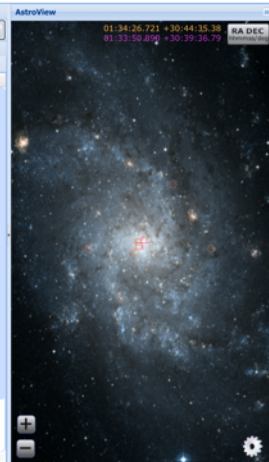
[Login...](#)

[Account Info...](#)

Displaying 19 of 13391 Total Rows

Filters

| Actions | Observation T. | Mission | Provenance Name | Instrument | Project | Filters | Wavelength |
|---------|----------------|---------|-----------------|------------|---------|---------|------------|
| 1 | *** | science | FUSE | FUV | | | UV |
| 2 | *** | science | FUSE | FUV | | | UV |
| 3 | *** | science | FUSE | FUV | | | UV |
| 4 | *** | science | FUSE | FUV | | | UV |
| 5 | *** | science | FUSE | FUV | | | UV |
| 6 | *** | science | FUSE | FUV | | | UV |
| 7 | *** | science | FUSE | FUV | | | UV |
| 8 | *** | science | FUSE | FUV | | | UV |
| 9 | *** | science | FUSE | FUV | | | UV |
| 10 | *** | science | FUSE | FUV | | | UV |
| 11 | *** | science | FUSE | FUV | | | UV |
| 12 | *** | science | FUSE | FUV | | | UV |
| 13 | *** | science | FUSE | FUV | | | UV |
| 14 | *** | science | FUSE | FUV | | | UV |
| 15 | *** | science | FUSE | FUV | | | UV |
| 16 | *** | science | FUSE | FUV | | | UV |





What we need: tools for big data discovery

Roman's greatest impact will be through its archival science.

To amplify this science we need tools that allow collaborations to

- quickly find what they are looking for
- explore the image and catalog space
- integrate *Roman* with other large data sets

Here are 5 main areas of expansion to achieve these goals:

1. Big Data Catalog Discovery

3. Search By Example

2. Big Data Image Discovery

4. Big Data Fast Survey APIs

5. Curated Complementary Data



1. Big Data Catalog Discovery

Users will need a fast source catalog query system that integrates with Jupyter

Cloud-based Catalog Data Infrastructure

- Integration with commercial cloud services, Roman Science Platform
- Asynchronous search
- TAP/ADQL endpoints (consistent with *Gaia*, *Rubin*)
- Gigascale Crossmatching (e.g. GIS-Based Greenplum, AXS)
- Jupyter Viz Stack interface: Notebook, Platform, Webpage



The Jupyter Viz Stack: Notebook, Platform, Webpage



Jdaviz Notebook

This is an auto-generated notebook to access JWST file `jd00736-0039_t001_miri_ch1-long_s3d.fits` in the [Cubeviz](#) Python package.

This notebook attempts to download public JWST data. If the data is not public, it will attempt to authenticate using your MAST API authentication token. If you do not have one, go [here](#) to create a new token and set it to a new environment variable called `MAST_API_TOKEN`.

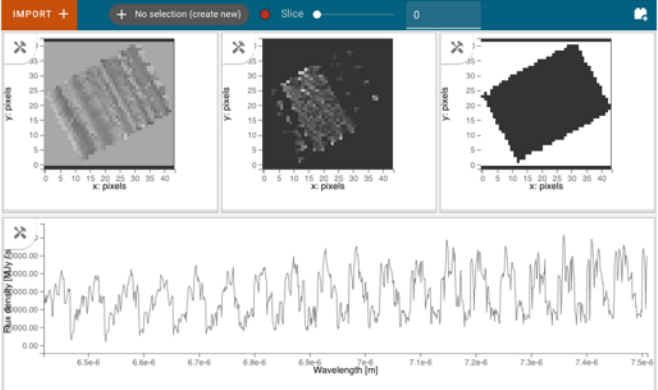
```
[34]: import os
      from astropy.utils.data import download_file
      from jdaviz import Cubeviz

[35]: # access any MAST auth token from the os environment
      auth_token = os.environ.get('MAST_API_TOKEN', '')

      # construct the http path
      path = f'https://mast.stsci.edu/portal_jst/Download/file?uri=MAST/product/jw00736-0039_t001_miri_ch1-long_s3d.fits'

      # download the data file
      try:
          output = download_file(path, cache=True)
      except Exception as e:
          try:
              output = download_file(path, cache=True, http_headers={'Authorization': f'Bearer {auth_token}'})
          except Exception as e:
              if '401' in str(e):
                  print(f'({e}): Please check you have a valid MAST auth token set.')
              else:
                  print(e)

[36]: # load the data file into the helper jdaviz class and display the application
      h = Cubeviz()
      h.load_data(output)
      h.app
```



Jdaviz.MAST

JW00736-0039_T001_MIRI_CH1-LONG_S3D.FITS

FAKE-2-HIP-65426

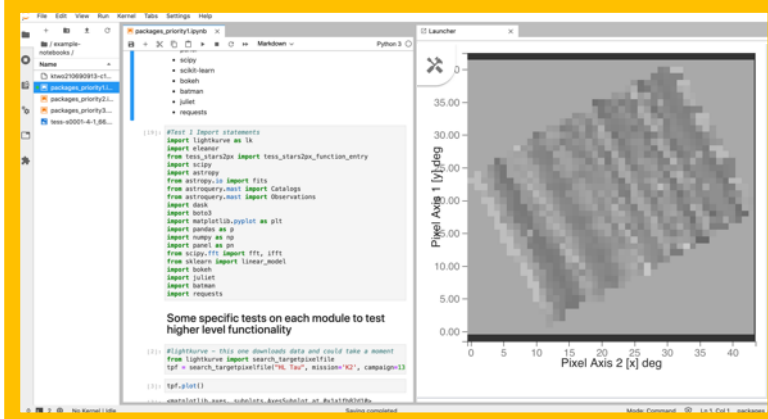
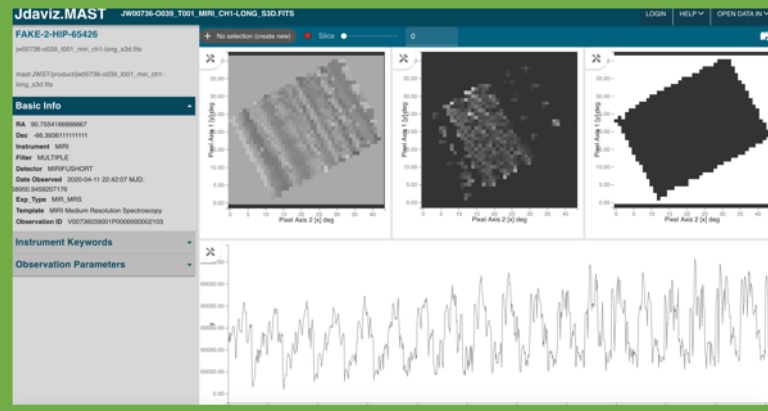
jd00736-0039_t001_miri_ch1-long_s3d.fits

Basic info

- RA: 161.76541699999997
- Dec: -66.32061111111111
- Instrument: MIRI
- Filter: MULTIPLE
- Detector: MIRI2-CMORT
- Date Observed: 2020-04-11 22:42:57 MJD
- ISSN: 9458207178
- Exp. Type: MRS_MRS
- Template: MIRI Medium Resolution Spectroscopy
- Observation ID: JW00736003901P0000000002103

Instrument Keywords

Observation Parameters



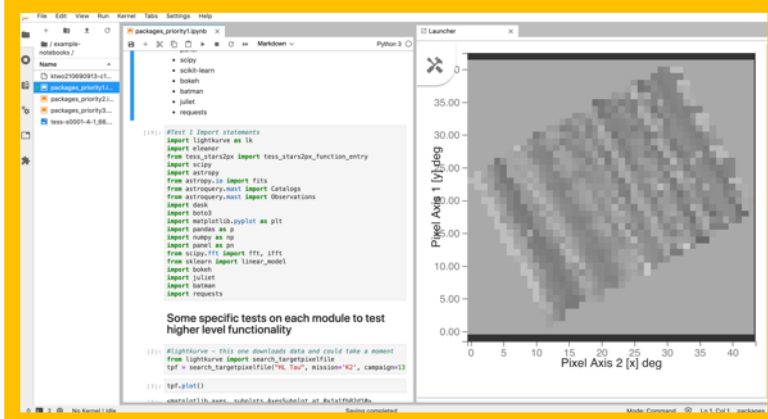
```
[1]: # Import statements
import lightcurve as lc
import astropy as ap
from test_stars2px import test_stars2px_function_entry
import numpy as np
import requests

from astropy.io import fits
from astropy.model import Catalogs
from astropy.model import Observations
import dash
import boto3
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import requests as req
from scipy.fft import fft, ifft
from sklearn import LinearModel
import boto3
import jupyterlab
import boto3
import requests

Some specific tests on each module to test higher level functionality

[2]: #lightcurve - this one downloads data and could take a moment
from lightcurve import search_target_in_file
lcf = search_target_in_file("file", "classical", "K2", "campaign13")

[3]: lcf.plot()
```

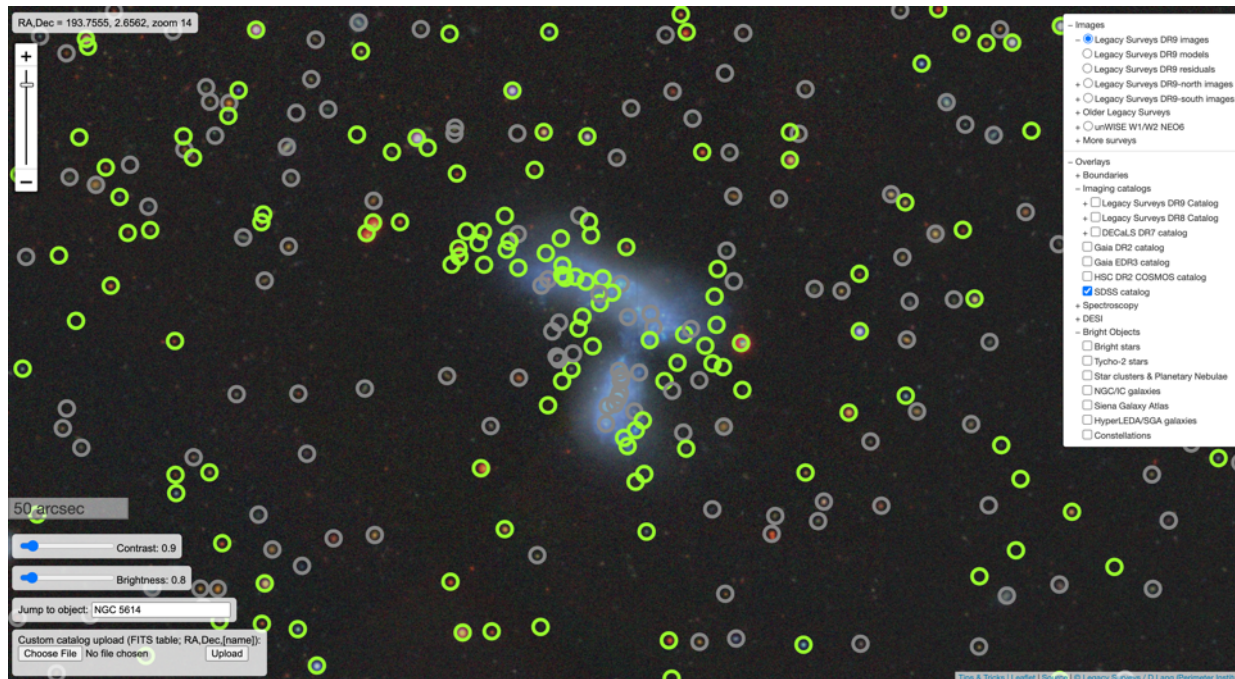




2. Big Data Image Discovery

Users will need to be able to quickly explore the surveys and image cutouts

- Survey Viewer



*example:
legacy survey*



2. Big Data Image Discovery

Users will need to be able to quickly explore the surveys and image cutouts

- Survey Viewer
- **Cutout Viewer**

table border="1">| obj list | page 1 |
| --- | --- |
| 274-51913-230 J103915.59-003918 | 275-51910-275 J104412.23+000907.1 | 275-51910-525 J104657.36+005334.7 | 276-51909-19 J105621.6-005320.4 | 278-51900-39 J111352.79+000014.4 |
| 278-51900-112 J111222.08-001518 | 278-51900-225 J110821.84-001257.5 | 278-51900-430 J110827.36+001456.3 | 279-51984-456 J111549.43+005136 | 279-51984-520 J111753.28-000025.2 |
| 281-51614-230 J112426.16-002537.2 | 282-51658-167 J113535.51-003505.9 | 285-51930-309 J115537.91-004615.5 | 286-51999-359 J120105.03+000650.3 | 288-52000-173 J121920.87-001431.1 |
| 349-51699-562 J170208.88-541221.6 | 353-51703-328 J170256.87+603346.8 | 353-51703-365 J170437.67+603506 | 355-51788-167 J171556.15+571416.7 | 355-51788-563 J172029.03+584749.1 |
| 358-51818-349 J172343.2+570025.1 | 387-51791-72 J000258.56+000831.1 | 389-51795-481 J001529.76+003823.9 | 390-51900-196 J002043.91-002823.9 | 390-51900-464 J002143.68+001745.5 |

*example:
SDSS Imaging*

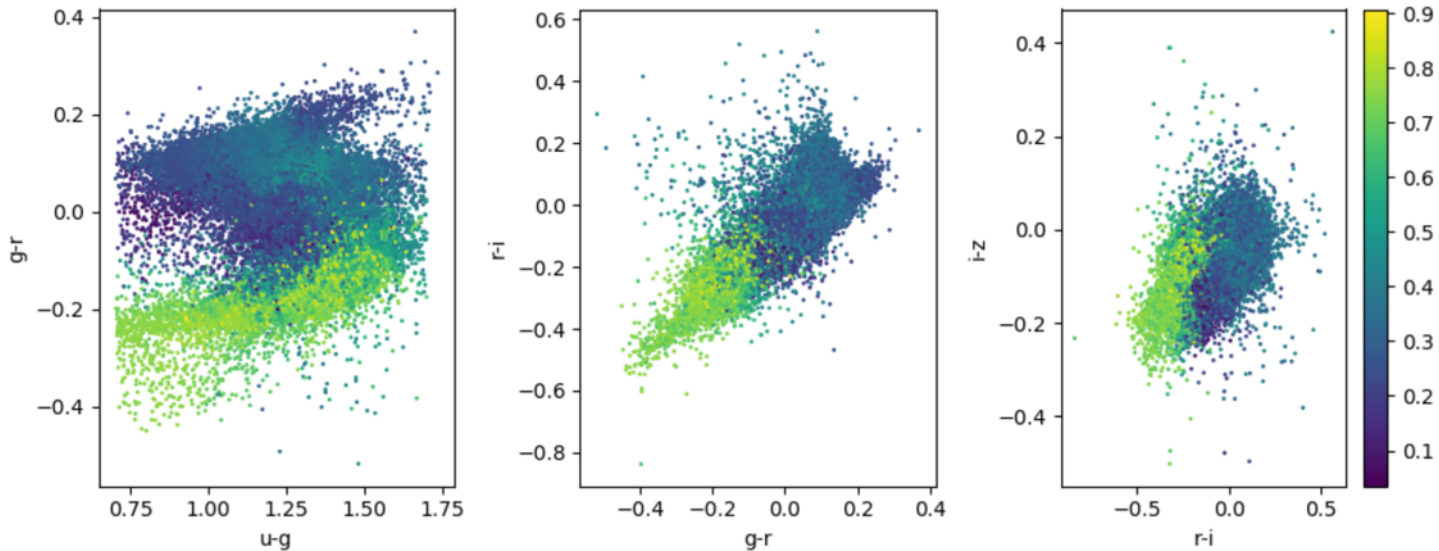




3. Machine Learning Search By Example

By being able to put in the coordinates, metadata, or images of scientifically relevant users can retrieve the full wealth of similar structures within the Roman archive

- Search by sources with catalog metadata (Classical Machine Learning)



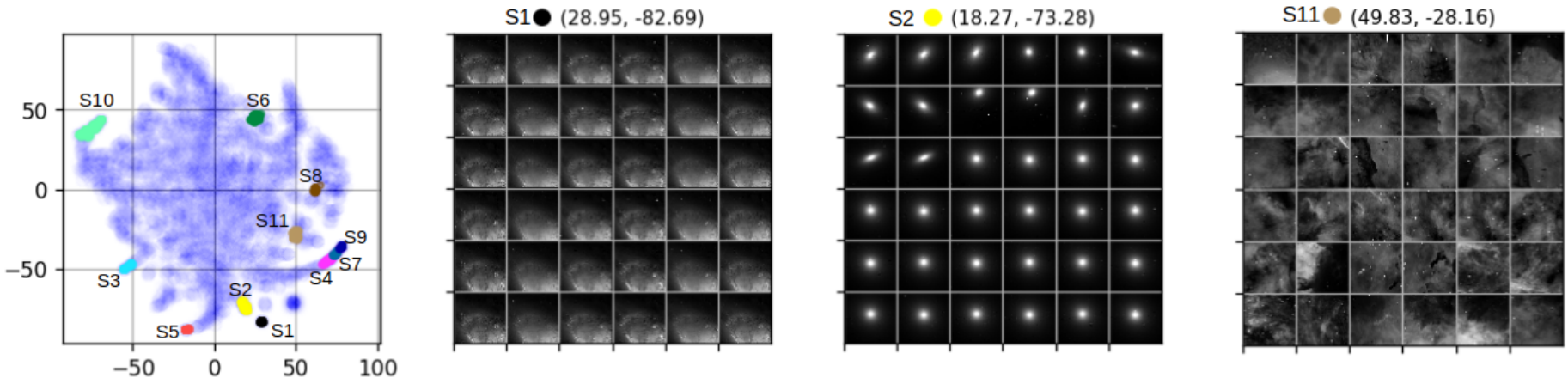
example: finding BHBs with KNN



3. Machine Learning Search By Example

By being able to put in the coordinates, metadata, or images of scientifically relevant users can retrieve the full wealth of similar structures within the Roman archive

- Search by sources with catalog metadata (Classical Machine Learning)
- **Search by Image (Deep Learning)**



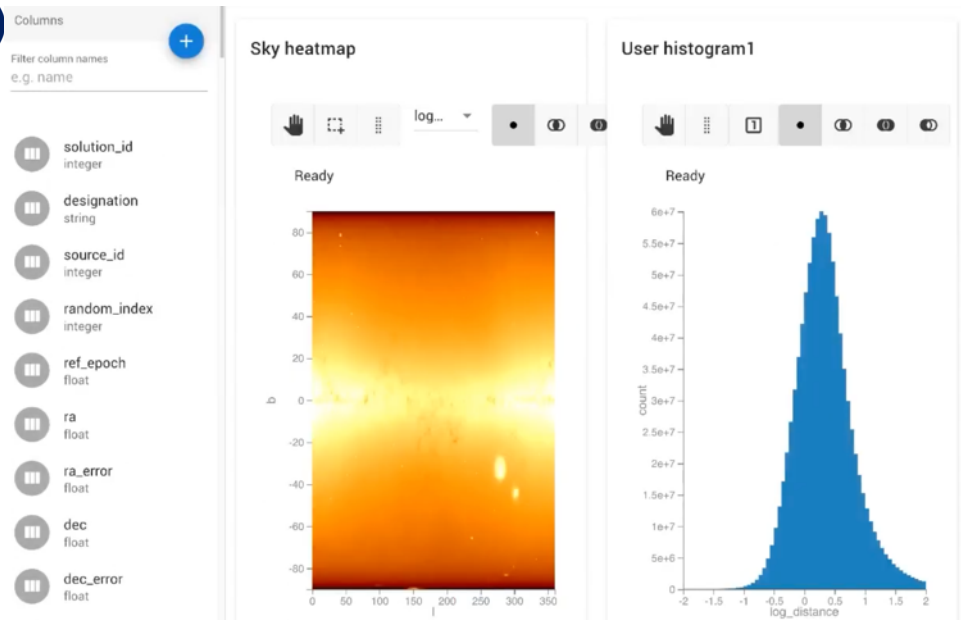
example: Finding similar ACS images with transfer learning



4. Big Data Fast Survey APIs

By providing fast, optimized services for common computationally intensive tasks, Roman can allow for advanced processing workflows both within and beyond the Roman Science Platform

- Histogram API (e.g. AXS, vaex)



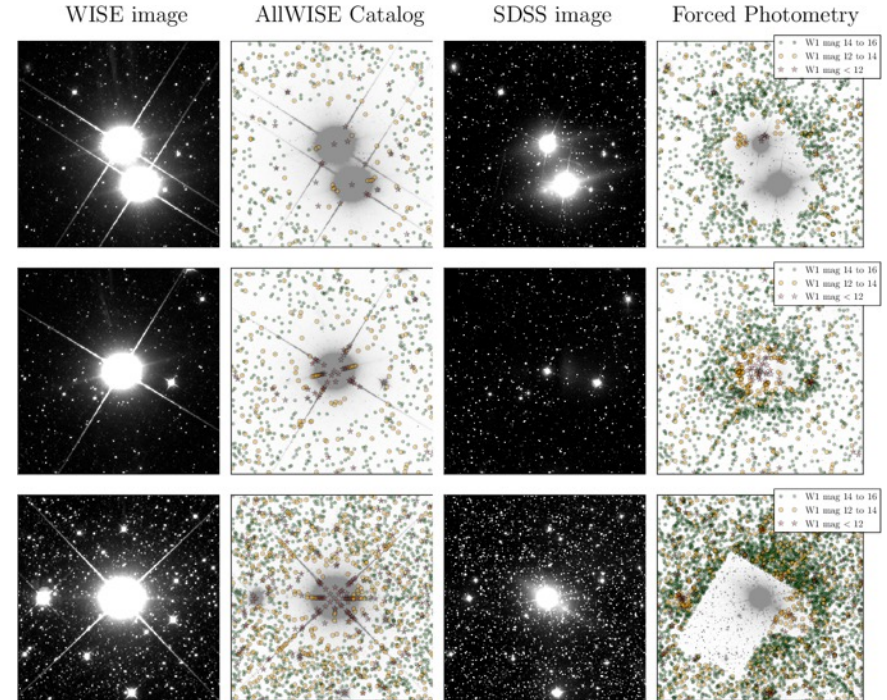
*example:
vaex+vuetify on Gaia*



4. Big Data Fast Survey APIs

By providing fast, optimized services for common computationally intensive tasks, Roman can allow for advanced processing workflows both within and beyond the Roman Science Platform

- Histogram API (e.g. AXS, vaex)
- **Forced Photometry API**





5. Curated Complementary Data

While Roman will be combined with almost every other astronomical data set, hosting copies of specific, well-chosen data sets for fast comparisons will dramatically amplify the scientific impact of Roman

- LSST/Rubin co-adds on *Roman* survey area
- Euclid Vis & IR co-adds on *Roman* survey area
- eROSITA on *Roman* survey area
- Subaru HSC on *Roman* survey area
- HST on *Roman* survey area