



# How big is big data?



# Working with large astronomy datasets

150,000,000 rows (of 11 double-precision columns) would take ~2 hours over typical home connection, ~10 seconds over internal 10 Gb/s connection, and ~1 s to pull out of RAM from local machine.

We have tables with 15 billion rows (200 hours)

1 TB can be transferred in 20 minutes over internal DL network, pulled from RAM in 1 minute, ~1 week over typical home connection



# Science Platforms

## Why?

- Growing data volume *and* complexity;
- New mode(s) of doing research that are Data-driven and/or Archive-driven

## Goals:

- Maximize scientific output of community: legacy science, versatile tools, collaborative workspaces, joint analysis of large datasets, serendipitous discoveries
- Framework for robust science: data quality, reproducible workflow, data longevity



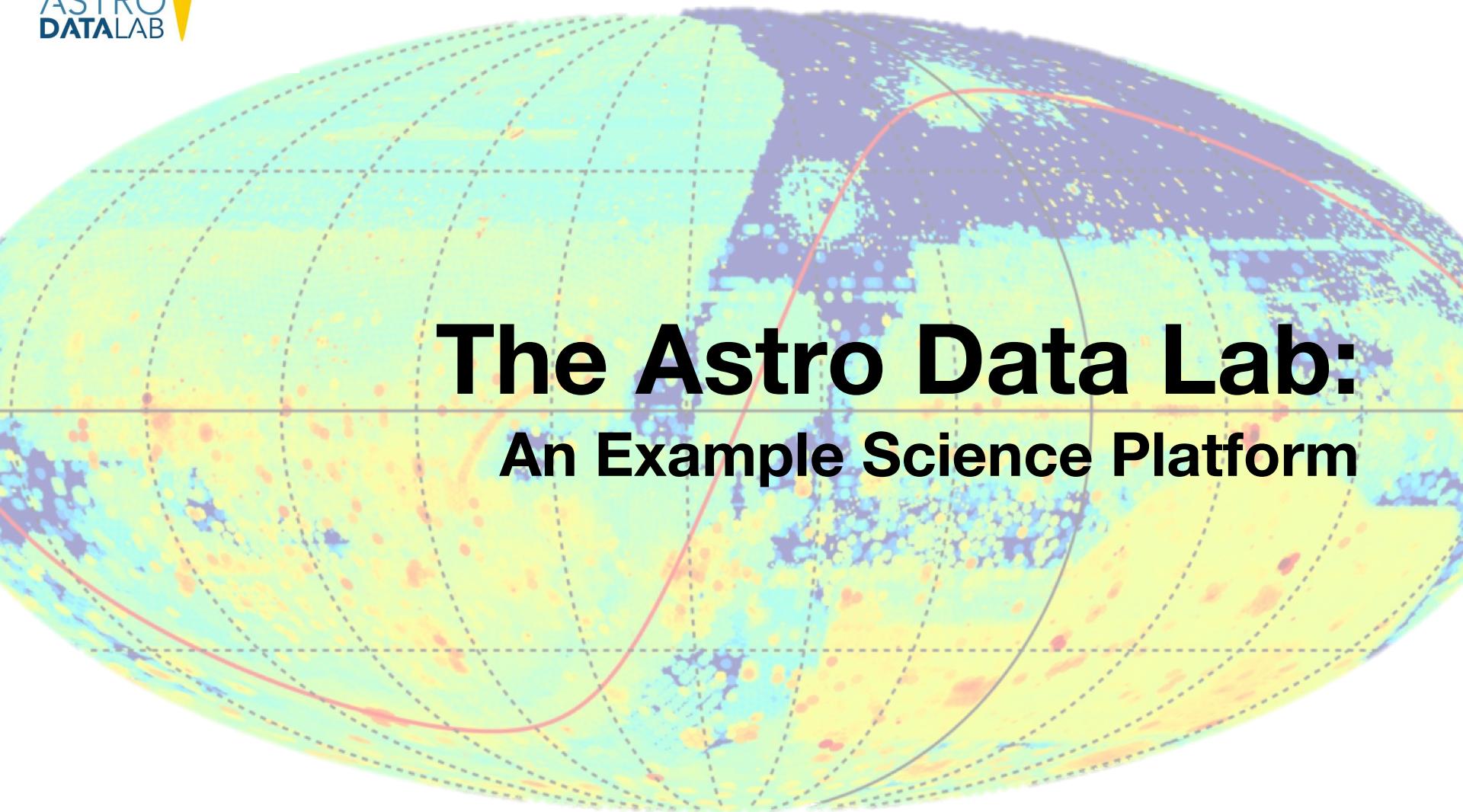
# Science Platforms

## Roles of science platforms:

- Bring the analysis to the data (software tools)
- Lower the barrier of entry: **user-friendly** interfaces, training of workforce at all career stages, tutorials and collaboration with educational institutions or across disciplines
- Coordinate among science platforms: share expertise & lessons learnt, similar interfaces/technologies, data/code transfer

## What?

- Software (& hardware) with a backend and a client side
- Standard protocols (e.g., VO) “under the hood”
- Documentation and tutorials (user manual, helpdesk ++)
- “Everything up to writing the paper”



# The Astro Data Lab: An Example Science Platform



NSF's National Optical-Infrared  
Astronomy Research Laboratory





# Data Lab Team

Cross-section of astronomy, data science & computer science

Current team:

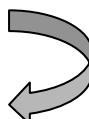
- [Adam Bolton](#), Acting Program Coordinator
- [Mike Fitzpatrick](#), Lead Software Engineer
- [Benjamin Hauger](#), System Administrator
- [Wendy Huang](#), Lead Web Developer, Technical Manager
- [Stephanie Juneau](#), Data Scientist
- [David Nidever](#), Data Scientist
- [Robert Nikutta](#), Project Scientist
- [Adam Scott](#), Database Architect
- [Benjamin Weaver](#), Data Scientist

## Goal:

- Efficient exploration and analysis of **large** astronomy datasets

## Approach:

- Science question
- Data discovery
- Build intuition through interaction with selected catalog and image set(s)
- Develop workflow (automated)
- Scientific discovery!

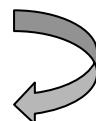


## Goal:

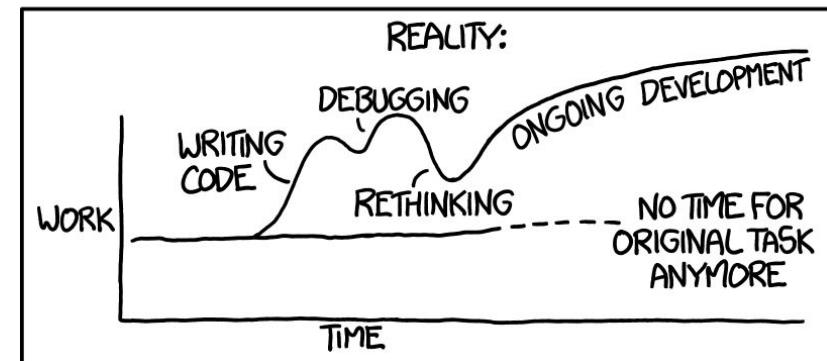
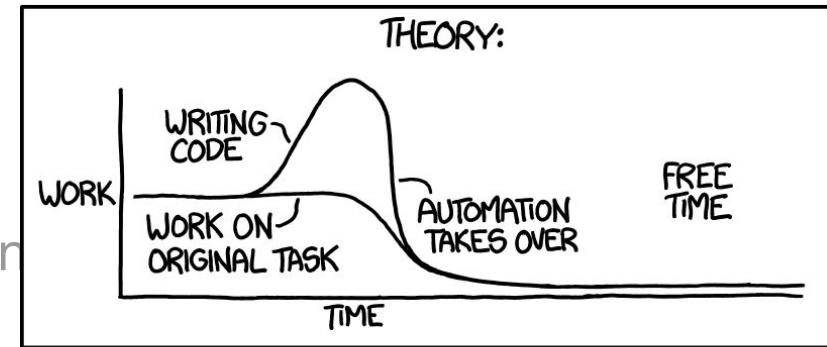
- Efficient exploration and analysis of **large** astronomy datasets

## Approach:

- Science question
- Data discovery
- Build intuition through interaction with image set(s)
- Develop workflow (automated)**
- Scientific discovery!

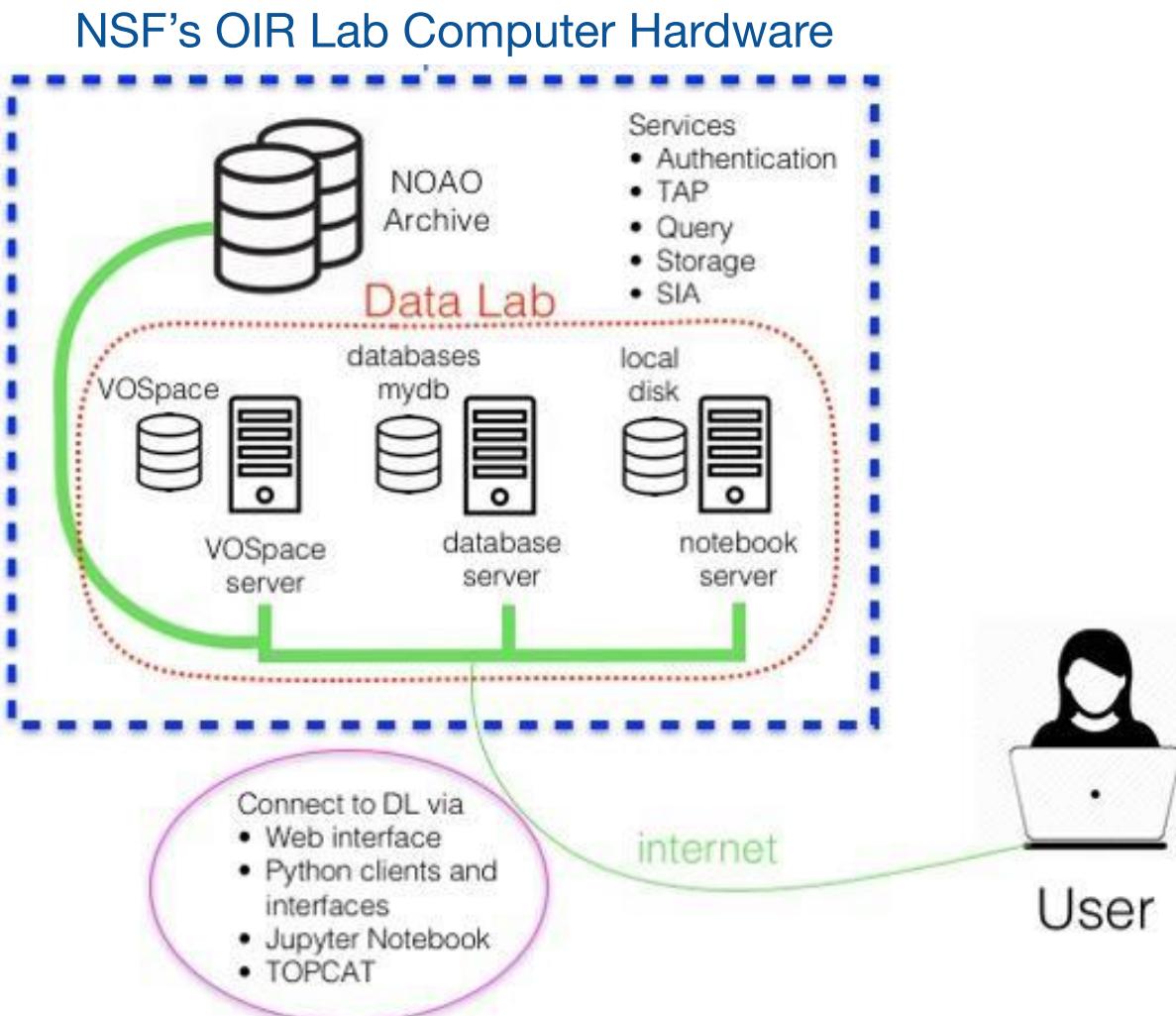


"I SPEND A LOT OF TIME ON THIS TASK.  
I SHOULD WRITE A PROGRAM AUTOMATING IT!"



( [xkcd.com/1319/](http://xkcd.com/1319/) )

# Working close to the data





# Data Lab in a Nutshell

**Large Catalogs** – TB-scale databases: DES, DESI imaging, NSC, AllWISE, SDSS, Gaia, etc.

**Pixel Data** – Get images and cutouts from OIR Lab's Science Data Archive

**X-matching** – Cross-match your catalog with any of ours, very fast

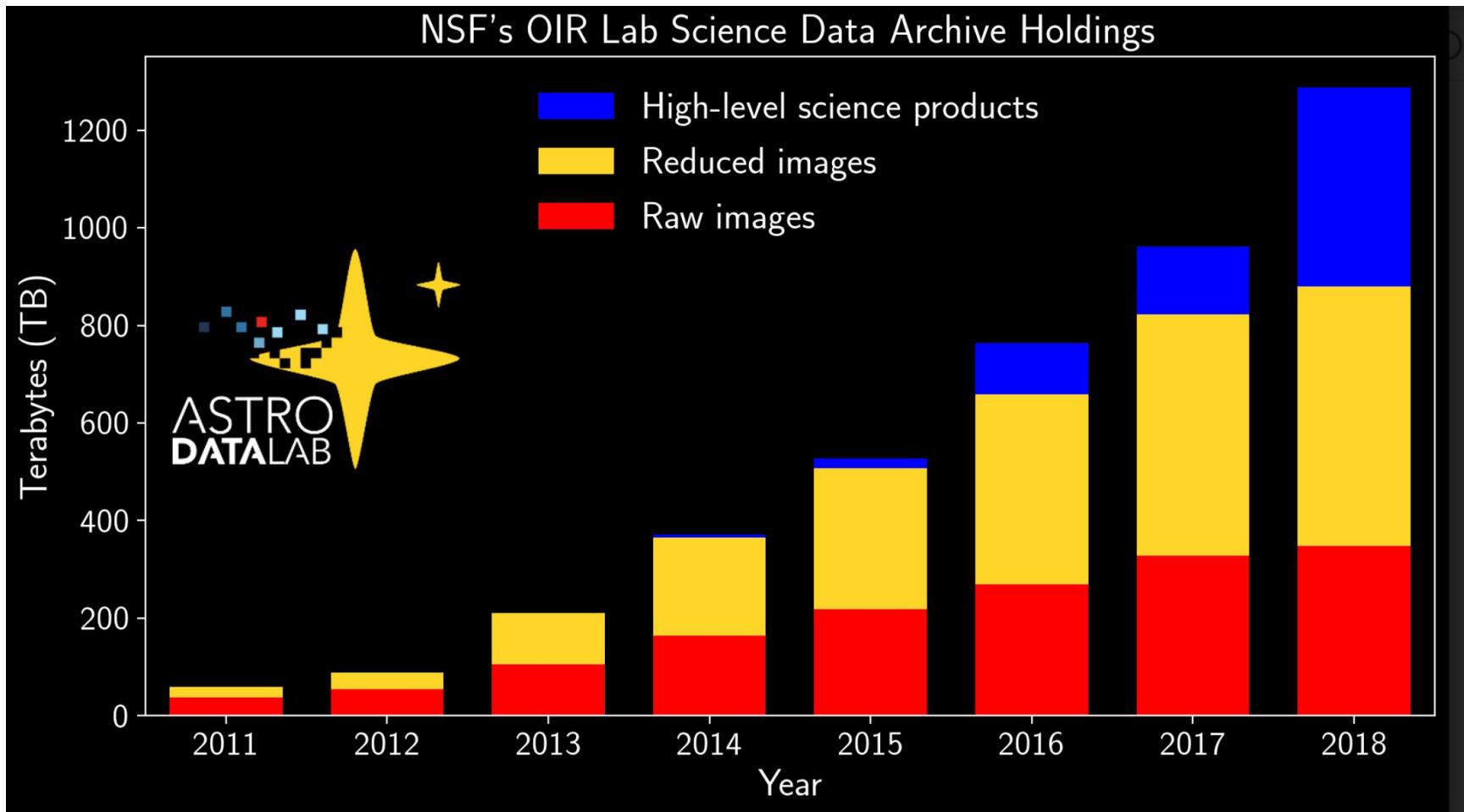
**Visualization** – Explore data interactively

**Analysis** – Run Python code on our Jupyter Notebook server

**Virtual Storage** – Store data & results: 1 TB in VOSpace and 250 GB in myDB per user!



# Data Holdings at OIR Lab Science Archive





# Using the Data Lab

ASTRO Data Lab

Login | Sign up

Service Status:

About Quick Start Tools Survey Data Docs/Help News/Events

Data Lab is hiring a Software Engineer. More information.

**<https://datalab.noao.edu>**

All-sky NOAO Source Catalog (2.9 billion objects, 34 billion measurements)

About Launch a Jupyter notebook Table Browser Query Interface Image Cutout X-match Service



# Using the Data Lab

## ASTRO Data Lab

Login | Sign up  
Service Status:

About Quick Start Tools Survey Data Docs/Help News/Events

People  
Data Lab Publications  
Acknowledgements  
Disclaimers

datalab.no

**ASTRO Data Lab is hiring a Software Engineer. More information.**

Column Information Query Interface Virtual Storage Job Status

- [allwise](#)
- [dad\\_dr1](#)
- [dad\\_dr2](#)
- [decaps\\_dr1](#)
- [des\\_dr1](#)
- [des\\_sval](#)
- [gaia\\_dr1](#)
- [gaia\\_dr2](#)

[gaia\\_dr2.allwise](#)  
[gaia\\_dr2.cepheid](#)

[gaia\\_dr2.des\\_dr1](#)

[gaia\\_dr2.gaia\\_source](#)

[gaia\\_dr2.iers](#)

[gaia\\_dr2.light\\_curves](#)

**Choose a database in the left panel then select the table you want!**

(The bold columns are indexed columns)

Column Name	Description	Datatype
a_g_percentile_lower	aGVal lower uncertainty	REAL
a_g_percentile_upper	aGVal upper uncertainty	REAL
a_g_val	line-of-sight extinction in the G band, A_G)	REAL
astrometric_chi2_al	AL chi-square value	DOUBLE
astrometric_excess_noise_rms	Browse the schema of Data Lab catalog tables	DOUBLE
astrometric_excess_noise_sig		DOUBLE

About

Launch a Jupyter notebook

Table Browser

Query Interface

Image Cutout

X-match Service

<https://datalab.noao.edu/about.php>



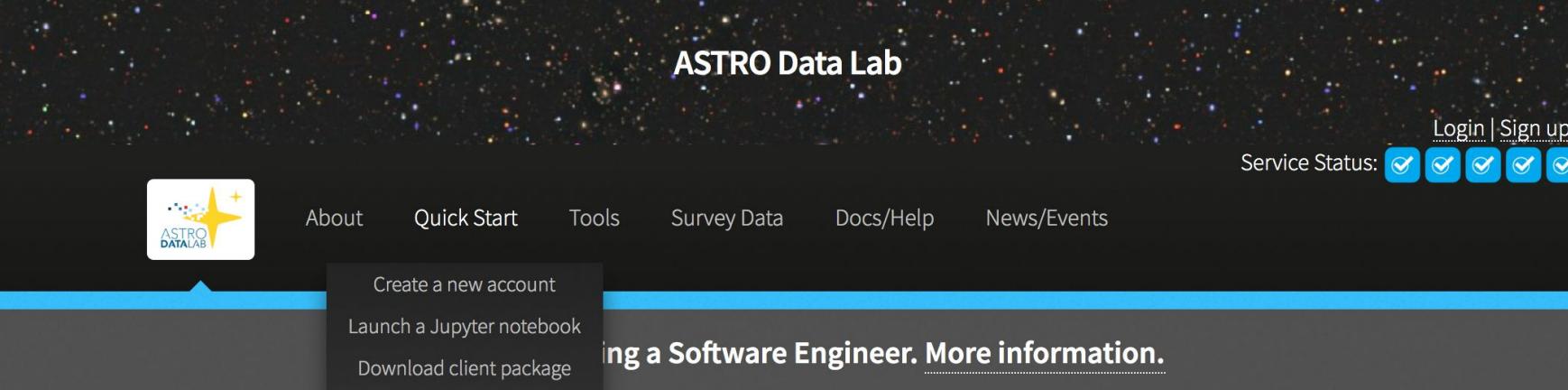
# Using the Data Lab

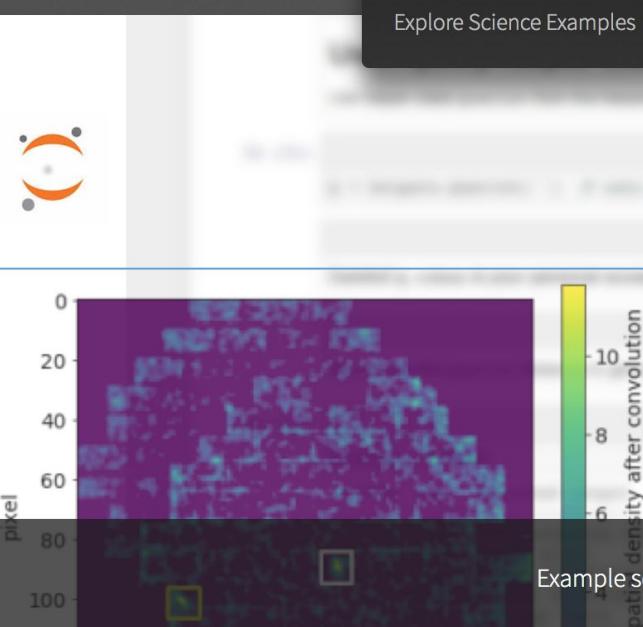
## ASTRO Data Lab

Login | Sign up  
Service Status:

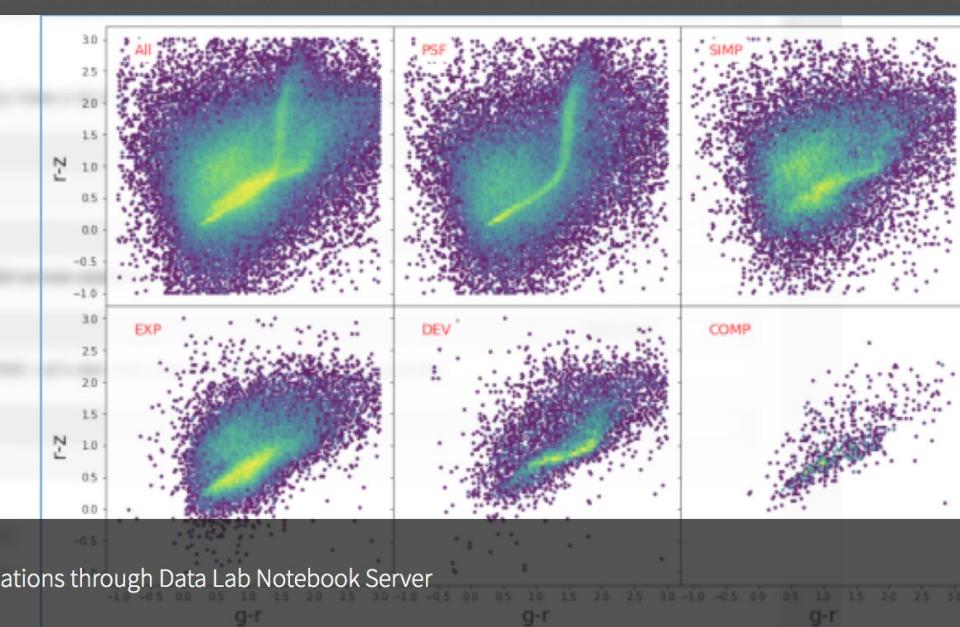
About Quick Start Tools Survey Data Docs/Help News/Events

Create a new account  
Launch a Jupyter notebook  
Download client package  
Explore Science Examples





Example science applications through Data Lab Notebook Server



[About](#) [Launch a Jupyter notebook](#) [Table Browser](#) [Query Interface](#) [Image Cutout](#) [X-match Service](#)

<https://datalab.noao.edu/start.php>



# Using the Data Lab

## ASTRO Data Lab

Login | Sign up

Service Status: ✓ ✓ ✓ ✓ ✓

About Quick Start Tools Survey Data Docs/Help News/Events

Data Lab is hiring engineer. More information.

All-sky NOAO Source Catalog (2.9 billion objects, 34 billion measurements)

[About](#) [Launch a Jupyter notebook](#) [Table Browser](#) [Query Interface](#) [Image Cutout](#) [X-match Service](#)

# Using the Data Lab: Survey Data

## NSF's OIR Lab facilities

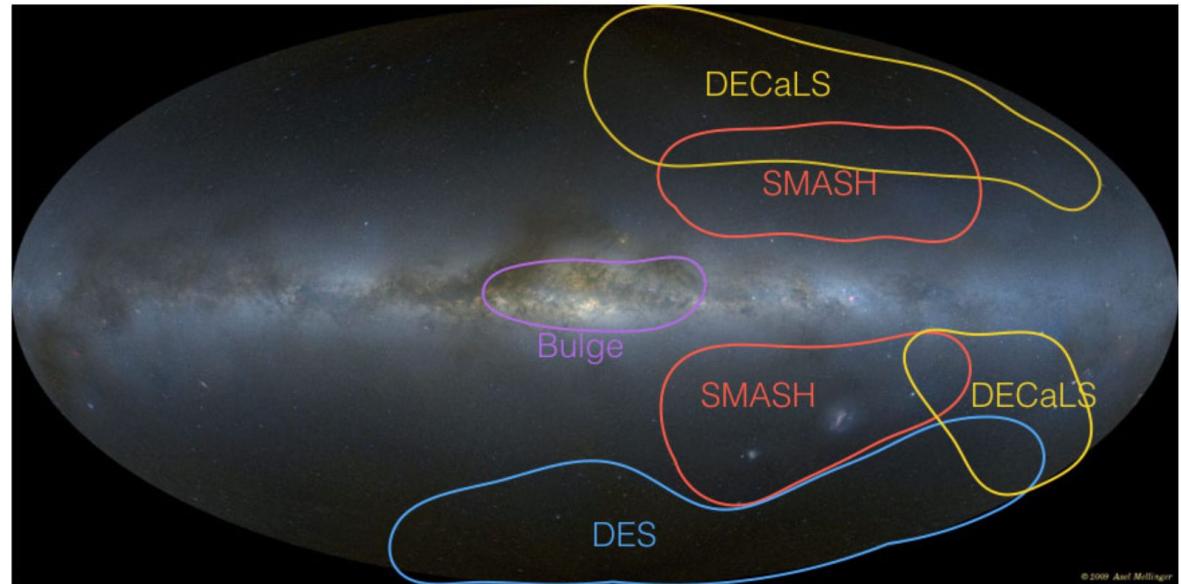
- [DECaPS](#)
- [DES](#)
- [Legacy Survey](#)
- [NOAO Source Catalog](#)
- [SMASH](#)

## External Datasets

- [DECam Asteroid Database](#)
- [DES SVA1](#)
- [GAIA](#)
- [PHAT](#)
- [SDSS DR13](#)
- [SPLUS](#)
- [UNWISE](#)
- [USNO](#)

## Survey Coverage

The map below shows the areas covered by surveys with catalog data currently available through the ASTRO Data Lab (SMASH, DECaLS) and those expected within approximately a year (DES, Bulge). Hover over an outline to see the survey name or click on an outline or on the sidebar links to go to the page for that survey. For pure image data from these and other observing programs, visit the [NOAO Science Archive](#).





# Using the Data Lab

**ASTRO Data Lab**

Login | Sign up

Service Status:

About Quick Start Tools Survey Data Docs/Help News/Events

Helpdesk User manual API documentation

Data Lab is hiring a Software information.

Example science applications through Data Lab Notebook Server

pixel

0 20 40 60 80 100

spatial density after convolution

0 10 8 6 4 2 0

g-r

R-Z

All PSF SIMP EXP DEV COMP

Launch a Jupyter notebook

About Table Browser Query Interface Image Cutout X-match Service

<https://datalab.noao.edu/docs.php>

The screenshot displays the homepage of the ASTRO Data Lab. At the top, there's a navigation bar with links for About, Quick Start, Tools, Survey Data, Docs/Help, and News/Events. A dropdown menu from the Docs/Help link shows options for Helpdesk, User manual, and API documentation. A banner in the center says "Data Lab is hiring a Software" and "information.". Below this, there are two main sections: one showing a heatmap of spatial density after convolution with axes R-Z and g-r, and another showing six scatter plots comparing R-Z vs g-r for different datasets: All, PSF, SIMP, EXP, DEV, and COMP. At the bottom, there are links for Launch a Jupyter notebook, Table Browser, Query Interface, Image Cutout, and X-match Service, along with a link to the documentation page.



# Datasets available at the Data Lab

**NSF's OIR Lab facilities** – MzLS (Mayall 4m, Kitt Peak), DECaLS, SMASH, DES (Blanco 4m, Chile), NSC (all)

**Steward Obs.** – BASS (Bok 2.3m, Kitt Peak)

**Public surveys** – SDSS DR13-16, GAIA DR1-2, AllWISE, unWISE, USNO and ++

Hundreds of TB more coming  
Total holdings at PB scale (including pixel data)



# Datasets available at the Data Lab

## OIR Lab/Steward Facilities Featured Surveys:

- [DESI imaging Legacy Survey \(LS\)](#): ~1.6 billion objects in DR8
- [SMASH](#): ~370 million objects in DR2
- [DES](#): ~400 million objects in DR1
- [DECaPS](#): ~2 billion objects
- [NOAO All-Sky Source Catalog \(NSC\)](#): ~2.9 billion objects
- [NSC \(+ single-epoch\)](#): ~3 billion objects (~30 billion measurements)

## Additional Surveys:

- selected tables from [SDSS/BOSS](#) DR13, 14 & 16, [All-WISE](#), [unWISE](#), [GAIA](#) DR1 & DR2, DES SVA1, the [Allen NEO catalog \(DAD\)](#), and USNO-A2/B, etc.

# Mayall 4m & Bok 2.3m telescopes on Kitt Peak, AZ



NSF's National Optical-Infrared  
Astronomy Research Laboratory





# Blanco 4m telescope in Cerro Tololo, Chile

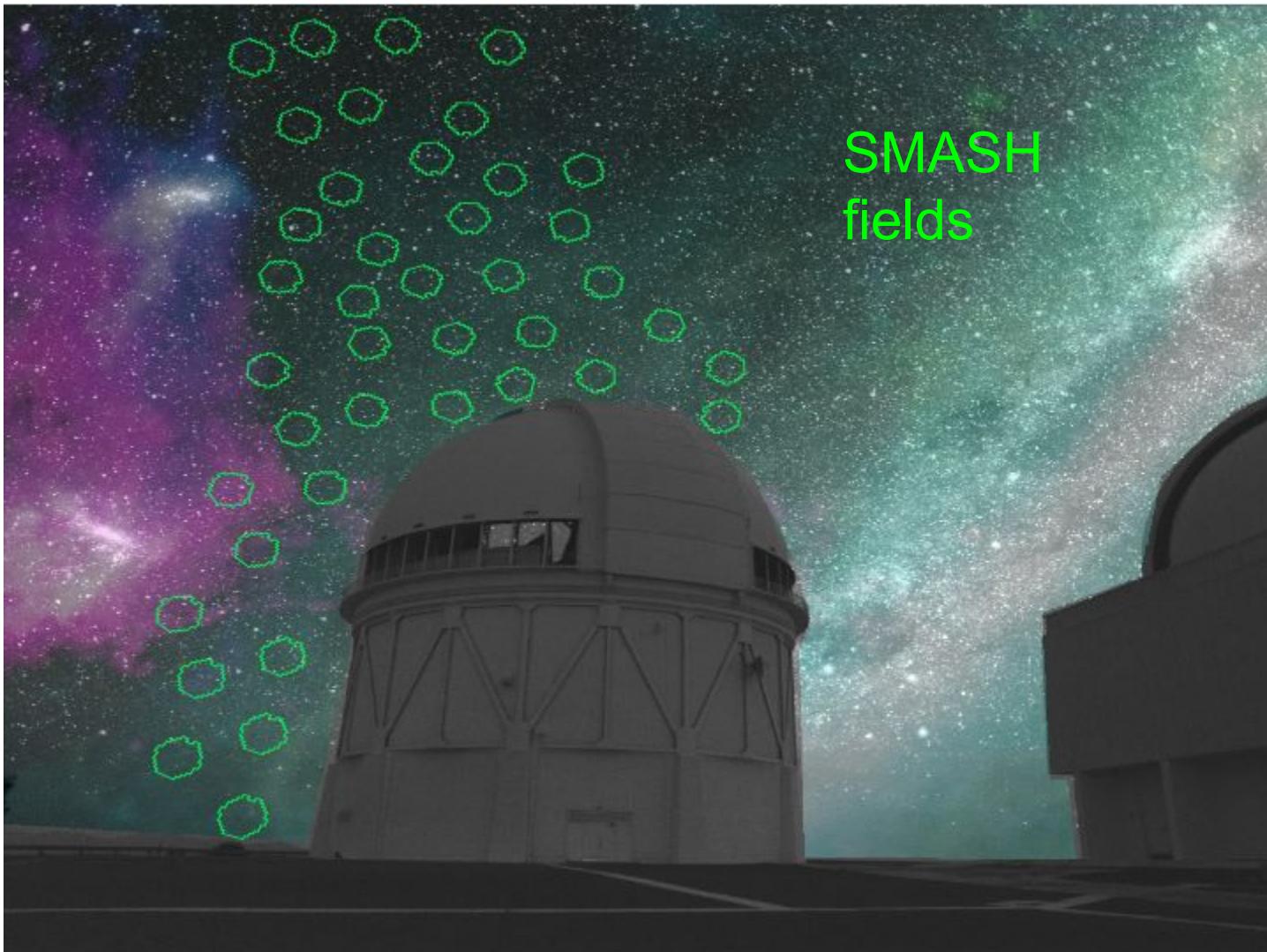


NSF's National Optical-Infrared  
Astronomy Research Laboratory





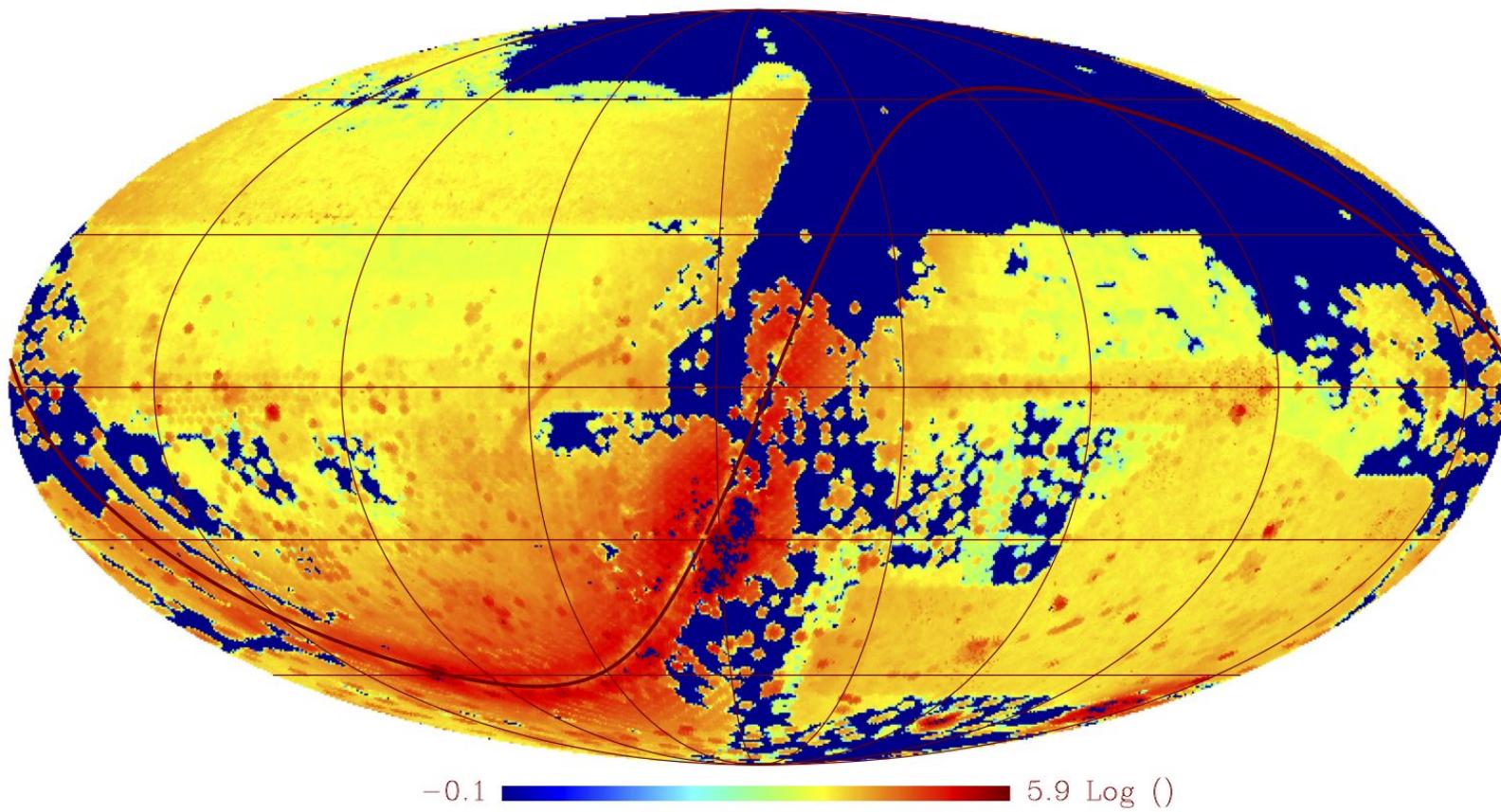
# Blanco 4m telescope in Cerro Tololo, Chile



NSF's National Optical-Infrared  
Astronomy Research Laboratory



# All-Sky NOAO Source Catalog: NSC



- ~2.9 billion objects, ~30 billion measurements; basic aperture photometry
- DR1 available (Nidever et al 2018, AJ); DR2 in progress



# Simulated dataset at the Data Lab

## TRILEGAL simulated LSST catalog:

- produced by Leo Girardi and Giada Pastorelli
- 19 billion stars for use in investigating LSST science cases

[datalab.noao.edu/tap](http://datalab.noao.edu/tap)

Column Information   Query Interface   Virtual Storage   Job Status

- [allwise](#)
- [dad\\_dr1](#)
- [dad\\_dr2](#)
- [decaps\\_dr1](#)
- [des\\_dr1](#)
- [des\\_sva1](#)
- [gaia\\_dr1](#)
- [gaia\\_dr2](#)
- [hipparcos](#)
- [hipparcos2](#)
- [ls\\_dr3](#)
- [ls\\_dr4](#)
- [ls\\_dr5](#)
- [ls\\_dr6](#)
- [ls\\_dr7](#)
- [ls\\_dr8](#)
- [lsst\\_sim](#)
- [lsst\\_sim.simdr1](#)
- [nsc\\_dr1](#)
- [phat\\_v2](#)
- [sdss\\_dr13](#)
- [sdss\\_dr14](#)
- [skymapper\\_dr1](#)
- [smash\\_dr1](#)
- [smash\\_dr2](#)
- [splus\\_dr1](#)
- [splus\\_edr](#)
- [stripe82](#)
- [tap\\_schema](#)

Choose a database in the left panel then select the table you want!

(The bold columns are indexed columns)

Column Name	Description	Datatype
av	V-band Extinction	REAL
c_o	Surface C/O ratio by number, n_C/n_O	REAL
cexcess	Carbon excess at surface, log(n_c-n_O)-log(n_H)+12, -1 for O-rich stars	REAL
comp	Flag for binarity, single star=0, primary=1, secondary=2, entire binary=3	SMALLINT
g_bpmag	Gaia G-BP color (Vegamag)	REAL
g_rpmag	Gaia G-RP color (Vegamag)	REAL
gaia_gmag	Gaia G magnitude (Vegamag)	REAL
galb	Galactic latitude	DOUBLE
gall	Galactic longitude	DOUBLE
gc	Galactic component, thin disk=1, thick disk=2, halo=3, bulge=4, MCs=5	SMALLINT
gmag	LSST g-band magnitude (ABmag)	REAL
htm9	HTM index (order-9 -> ~10 arcmin size)	INTEGER
imag	LSST i-band magnitude (ABmag)	REAL



NSF's National Optical-Infrared  
Astronomy Research Laboratory





# DESI

## Dark Energy Spectroscopic Instrument



- 14,000 square degrees
- >35 million spectra of galaxies and quasars!
- 10 million spectra of stars
- Commissioning started in 2019 (survey 2020-2025)

Object Class	Number of Spectra	Redshift Range
bright galaxies, $r < 19.5$	10 million	$0 < z < 0.4$
luminous red galaxies (LRGs)	4.2 million	$0.4 < z < 1.0$
emission line galaxies (ELGs)	18 million	$0.6 < z < 1.6$
quasars (QSOs)	2.4 million	$0.5 < z < 3.5$
Milky Way stars	10 million	---



→ [desi.lbl.gov](http://desi.lbl.gov)



NSF's National Optical-Infrared  
Astronomy Research Laboratory





# DESI at the Data Lab

- Host DESI imaging Legacy Surveys (DECaLS, BASS/MzLS)
  - Databases (latest: ls\_dr8)
  - Images in Science Archive (raw + processed)
- Host a copy of DESI targets
  - Database for final, public set of targets
- Host DESI redshifts
  - Database for public releases of redshift catalogs
  - Tools for spectra visualization/analysis
- Create example Notebooks & workflows
- Users can work with all data products

*now!*

*future*



NSF's National Optical-Infrared  
Astronomy Research Laboratory





# Summary of Current Functions

Function	Method
Sky exploration	Image discovery tool Catalog overlay tool Catalog visualization tool (prototype)
Authentication	Web interface datalab command Python authClient, DL interface
Catalog query	<b>Web interface</b> datalab command line (CLI) <b>Python queryClient</b> , DL interface TOPCAT
Image query	Simple Image Access (SIA) service Web Cutout Tool
Query result storage	myDB Virtual storage space
File transfer	datalab command and Virtual storage space
Analysis	Jupyter notebook server

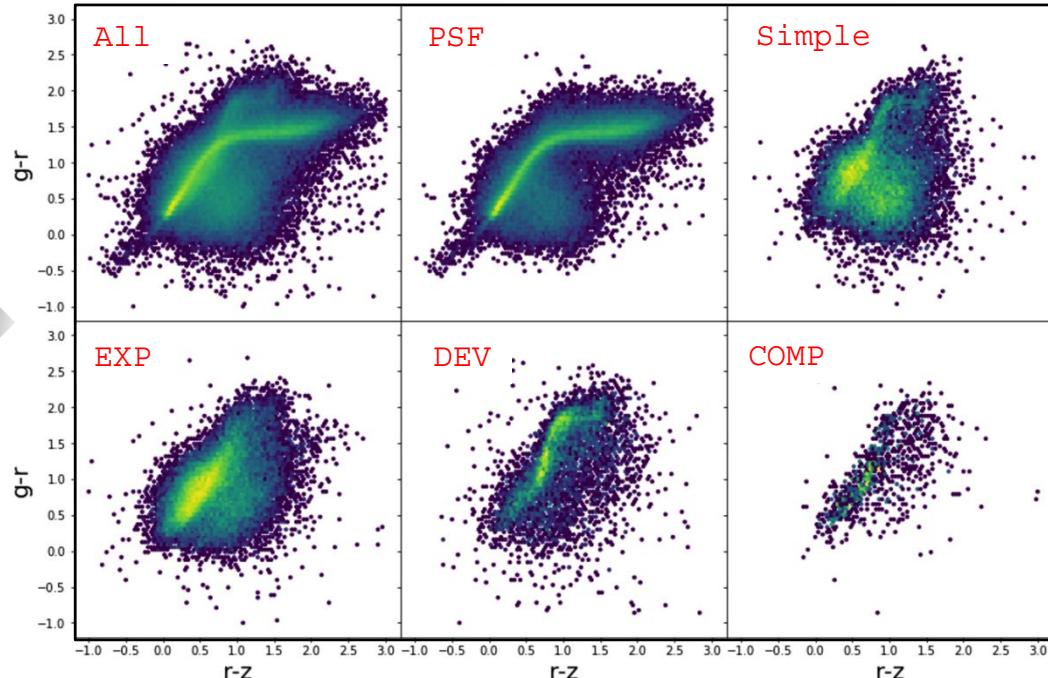
## Query to database:

*magnitudes and object shape (type)*

```
query = """
    SELECT dered_mag_g as gmag, dered_mag_r as rmag,
           dered_mag_z as zmag,
           dered_mag_w1 as w1mag, dered_mag_w2 as w2mag, type,
           snr_g, snr_r, snr_z, ra, dec
      FROM ls_dr3.tractor_primary
     WHERE (snr_g>3) and (snr_r>3) and (snr_z>3)
       LIMIT 200000"""

# dered_mag_g,r,z = AB mag in DECam g,r,z bands corrected
#                   for Galactic reddening
# dered_mag_w1,w2 = AB magnitudes in WISE bands W1 & W2
#                   corrected for Galactic reddening
# type             = object type (PSF, SIMP, EXP, DEV, COMP)
# snr_g,r,z        = signal-to-noise ratios (S/N) in g,r,z bands
# ra,dec          = celestial coordinates
#
# WHERE: requirement that S/N>3 in each DECaLS band
# LIMIT: returns 200,000 rows that satisfy the query
```

## Analysis: color-color plot per type



## Example Workflow

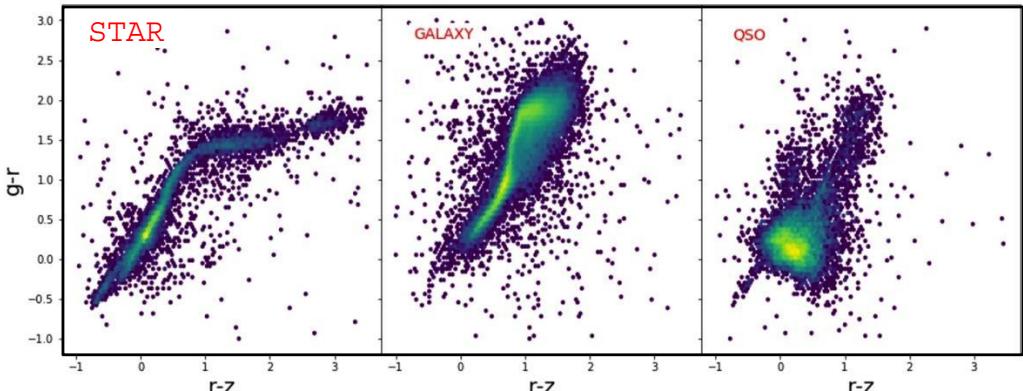
### Machine-Learning:

*Confusion matrix (spectroscopic training set)*

	GALAXY	0.982	0.008	0.001
QSO		0.087	0.878	0.035
STAR		0.018	0.012	0.97

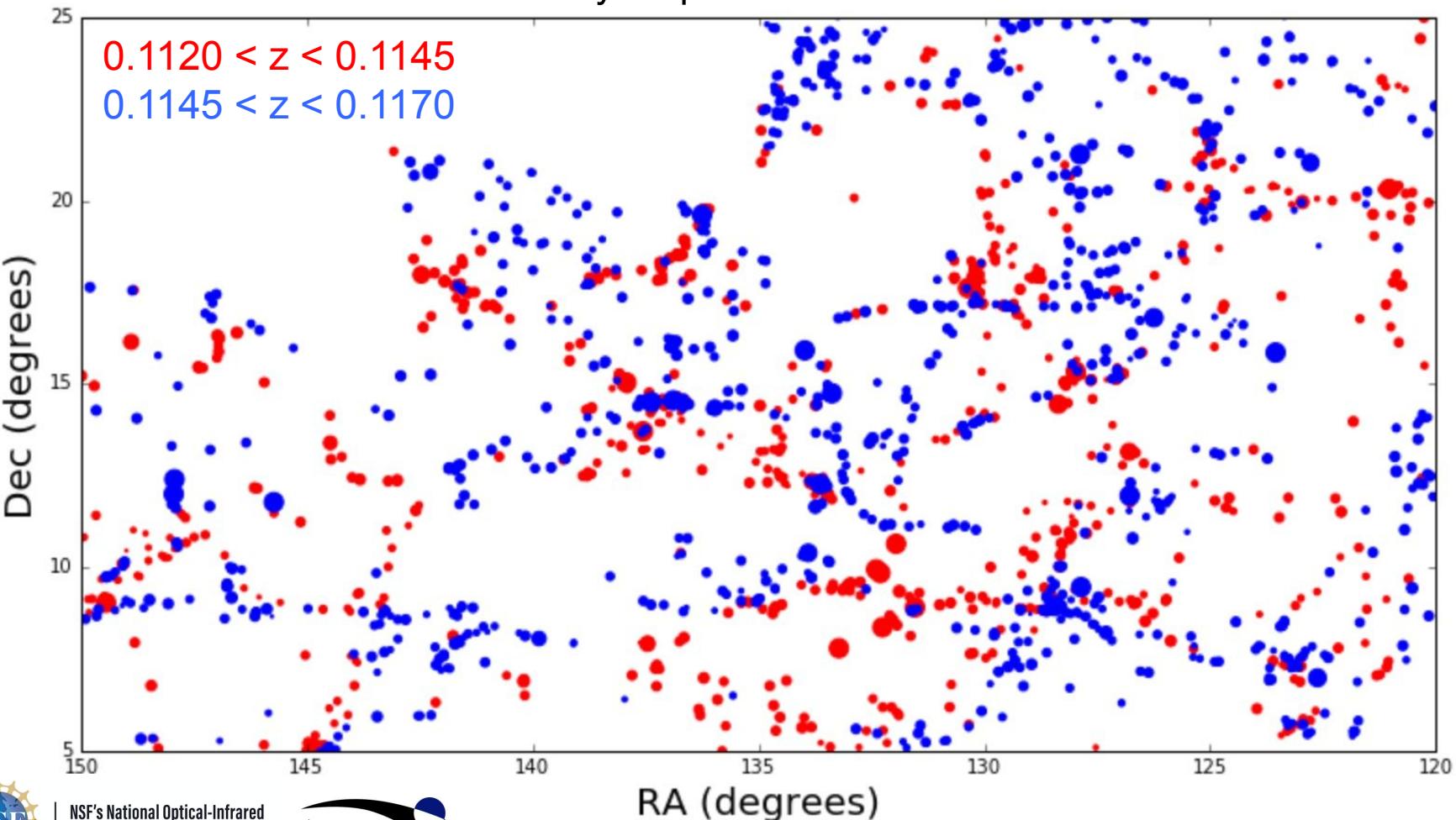
### Joint query:

*cross-match with SDSS spectroscopic class*



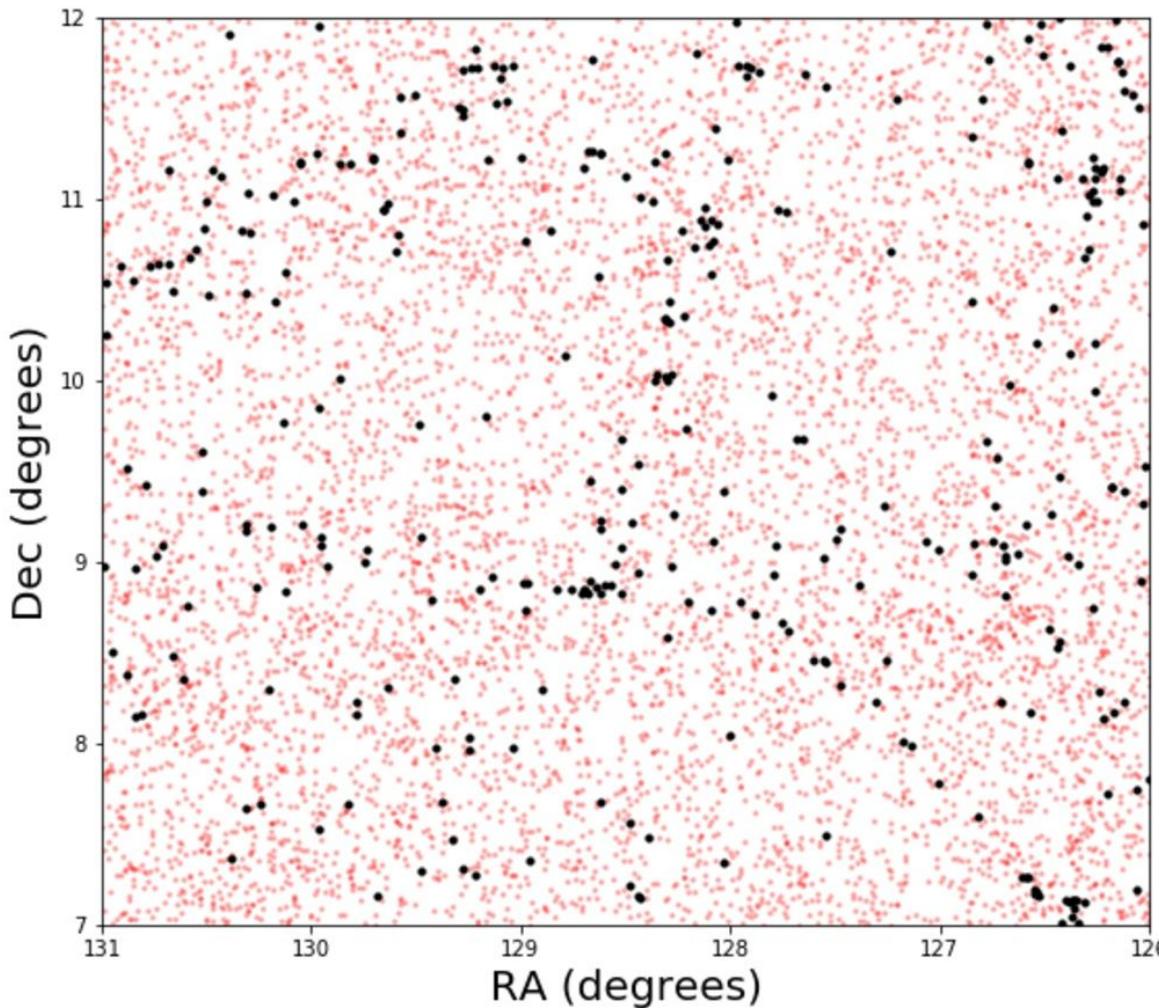
# Example Application: Large Scale Structures

- Query SDSS (specObj DR13)
- Galaxies in 2 thin redshift slices (750 km/s)
- Point-size coded to velocity dispersion



# Example Application: Large Scale Structures

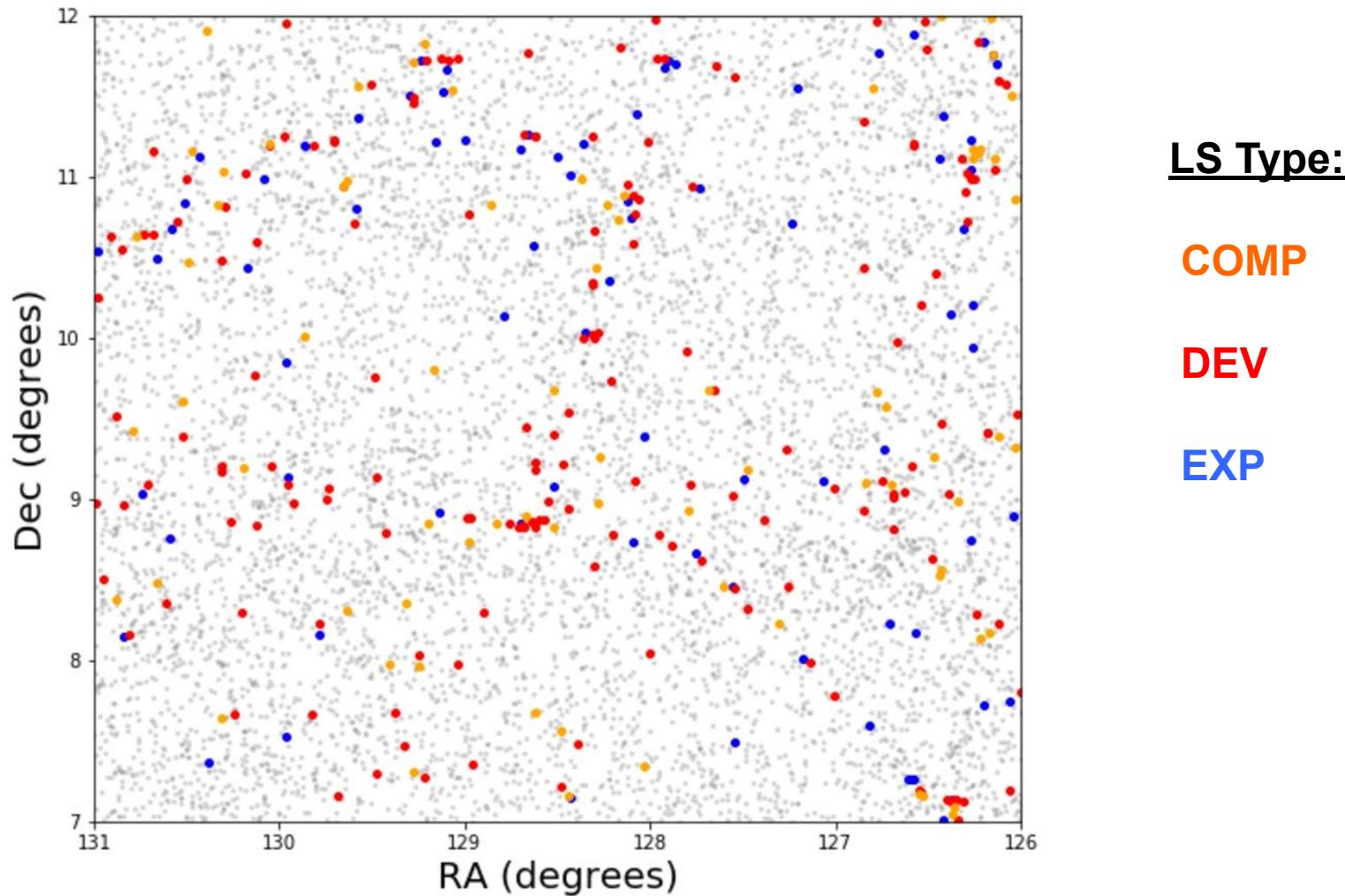
- Joint Query SDSS (specObj DR13) – ls\_dr3



All spec-z  
 $0.105 < z < 0.125$

# Example Application: Large Scale Structures

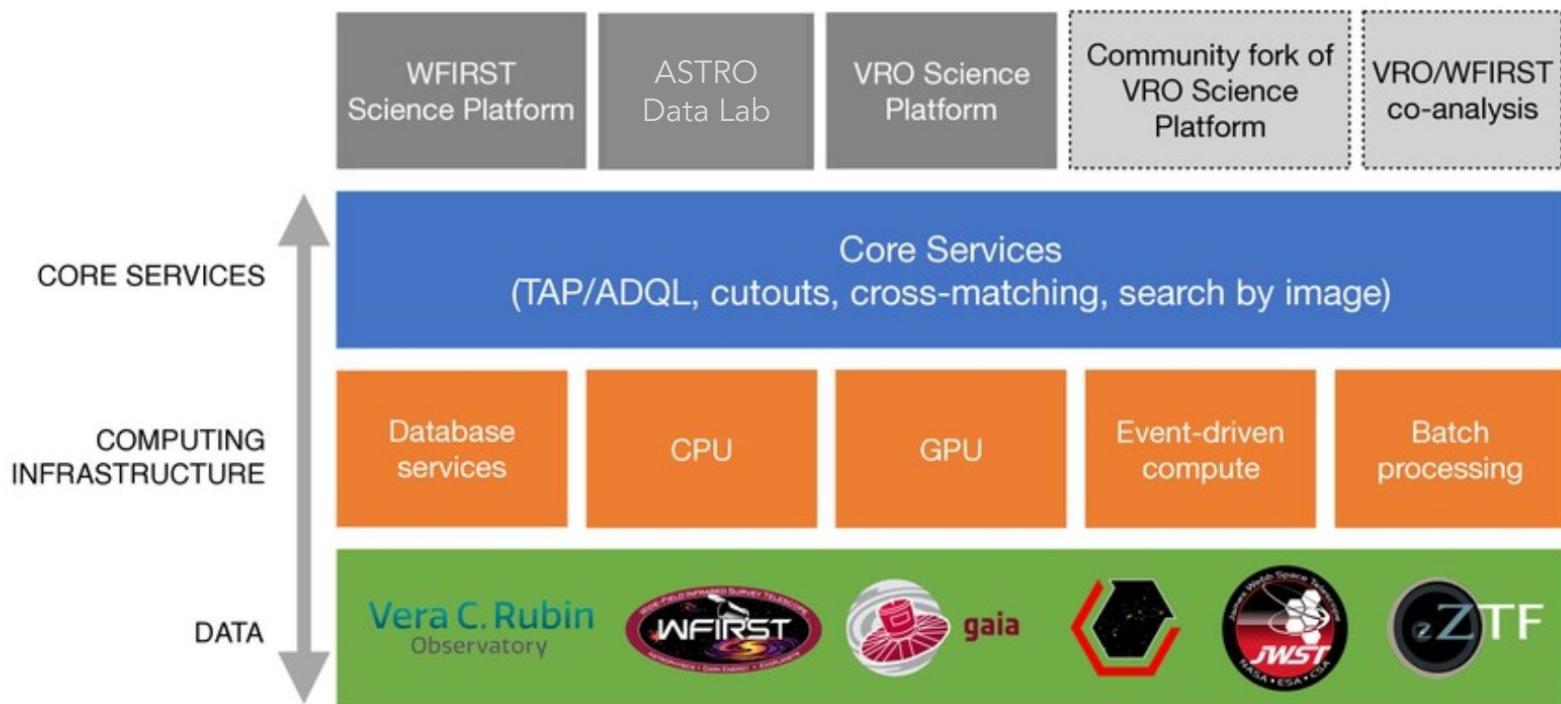
- Joint Query SDSS (specObj DR13) – ls\_dr3





# Future Astronomy Archives / Science Platforms

## Astronomy data commons



*...co-locate **data** with **cloud computing** infrastructure and commonly used **software services**, **tools & apps** for managing, analyzing and sharing data to create an **interoperable resource**...*



NSF's National Optical-Infrared  
Astronomy Research Laboratory



Slide by Arfon Smith (STSci)



# Try it out and get in touch!

- Web: [datalab.noao.edu](http://datalab.noao.edu)
- Email: [datalab@noao.edu](mailto:datalab@noao.edu)
- GitHub: <https://github.com/noaodatalab>
- Twitter: @DataLabAstro



NSF's National Optical-Infrared  
Astronomy Research Laboratory

