# Classification of 4FGL sources with CSCv2 and multi-wavelength surveys
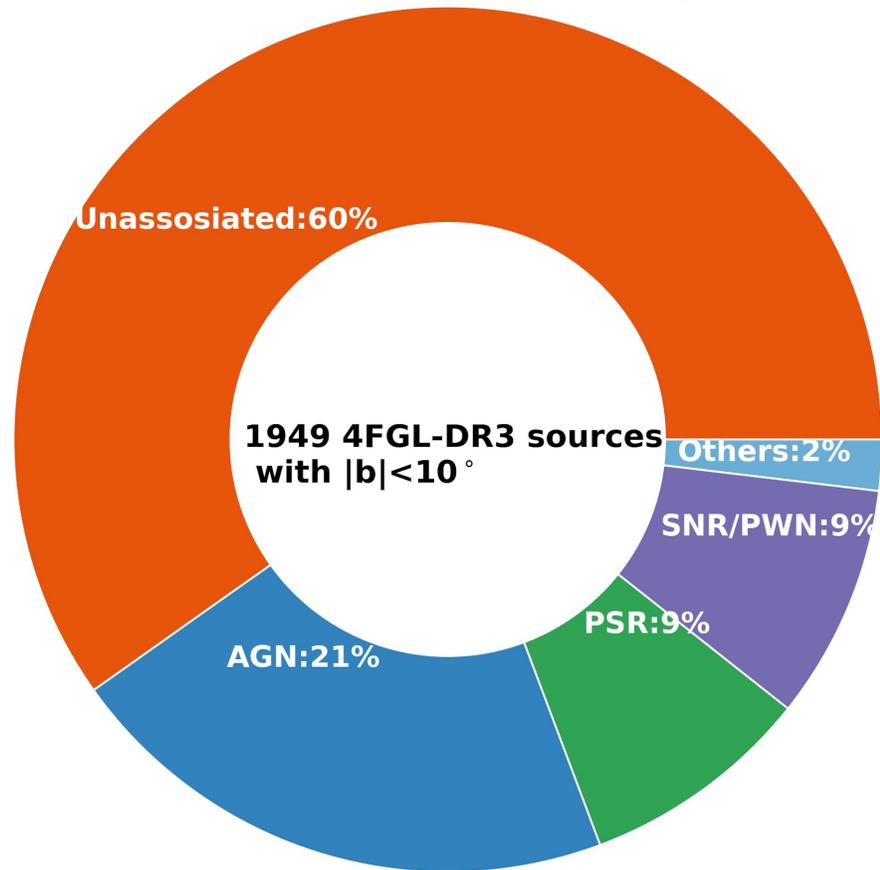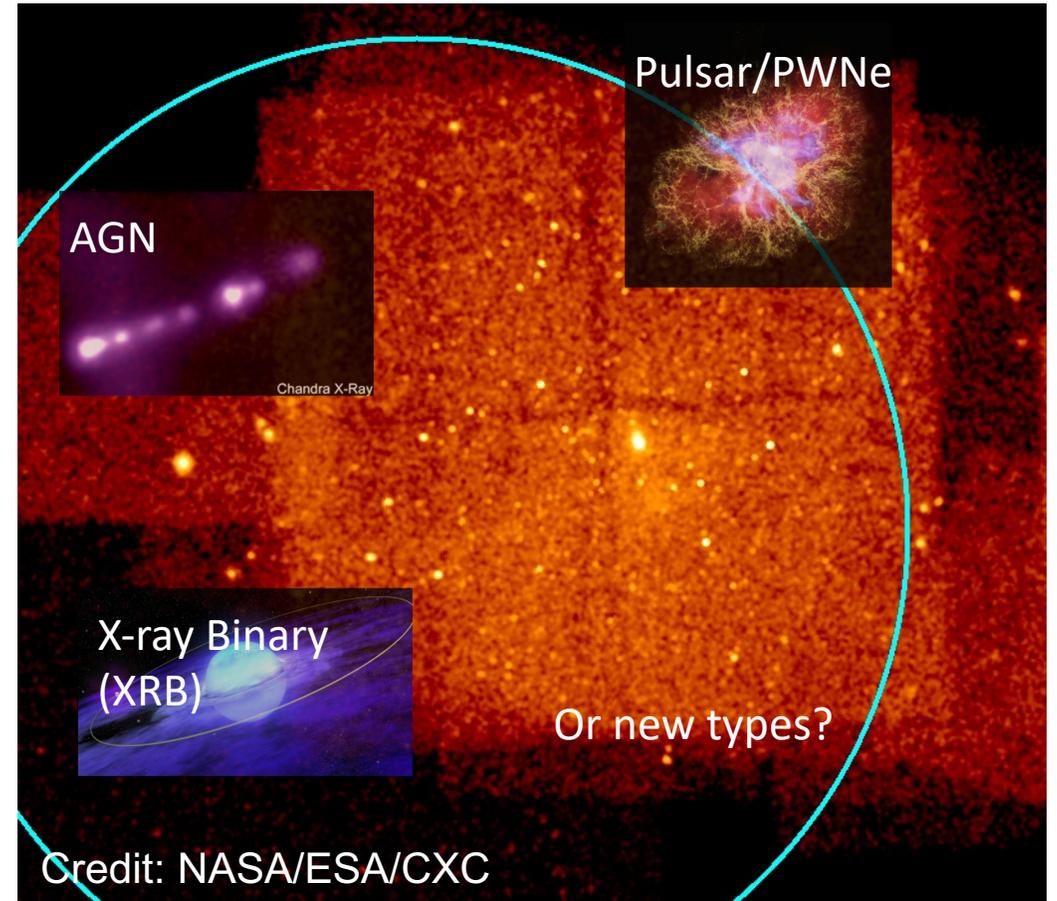
Hui Yang[1], Jeremy Hare[2], Oleg Kargaltsev[1], Steven Chen[1], Igor Volkov[1], Blagoy Rangelov[3]

[1] The George Washington University ,[2] NASA GSFC,[3] Texas State University

Fermi Summer School, Lewes Delaware, 2023

# Unassociated Galactic 4FGL-DR3 sources
## - identifying X-ray sources of $\gamma$-ray emitters



Unassosiated:60%

1949 4FGL-DR3 sources
with |b|<10°

Others:2%

SNR/PWN:9%

PSR:9%

AGN:21%

*Right*: Breakdown of 4FGL-DR3 classifications for Galactic plane sources within |b|<10



Pulsar/PWNe

AGN

Chandra X-Ray

X-ray Binary
(XRB)

Or new types?

Credit: NASA/ESA/CXC

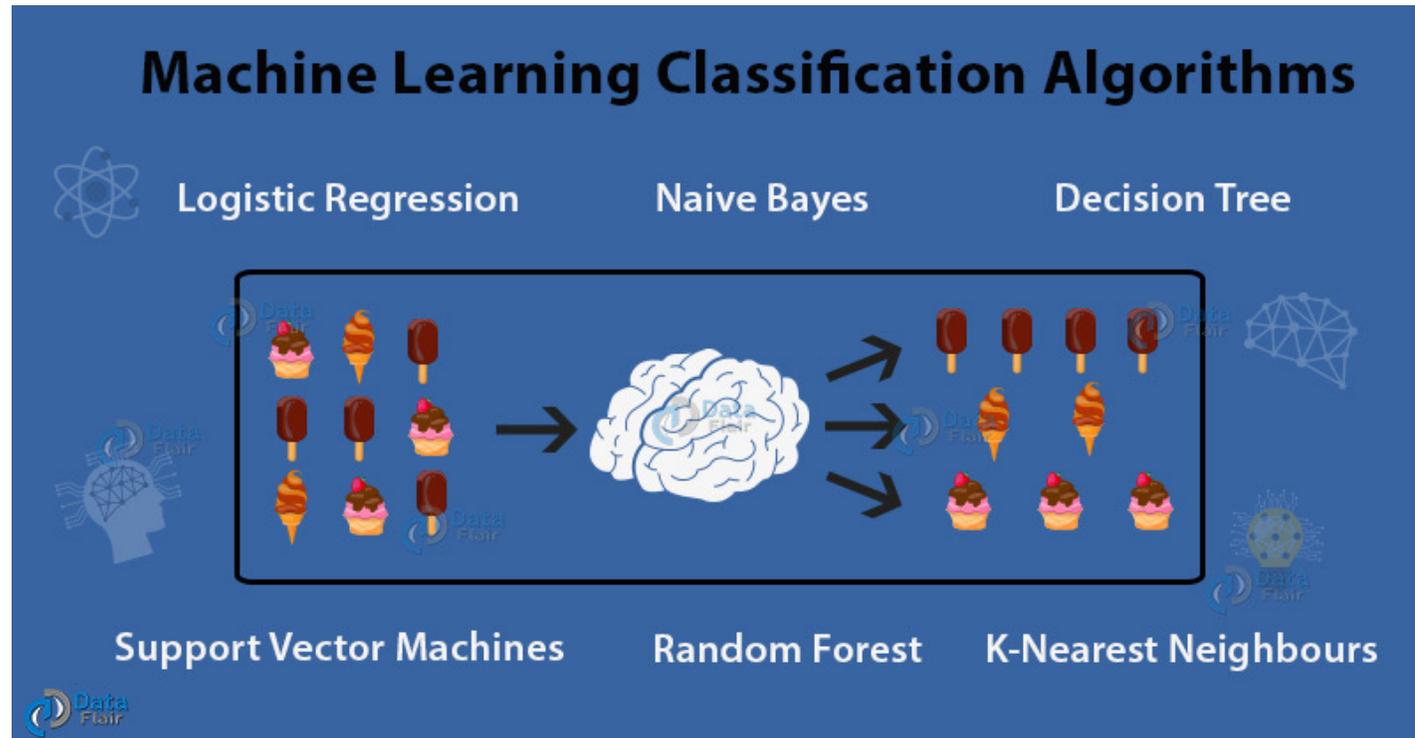Chandra image of $\gamma$-ray sources HESS J1809-1917 overlaid with different types of $\gamma$-ray emitters

# Many PSR, AGN, XRBs may look similar

Majority of X-ray sources are faint, and difficult to distinguish from each other.
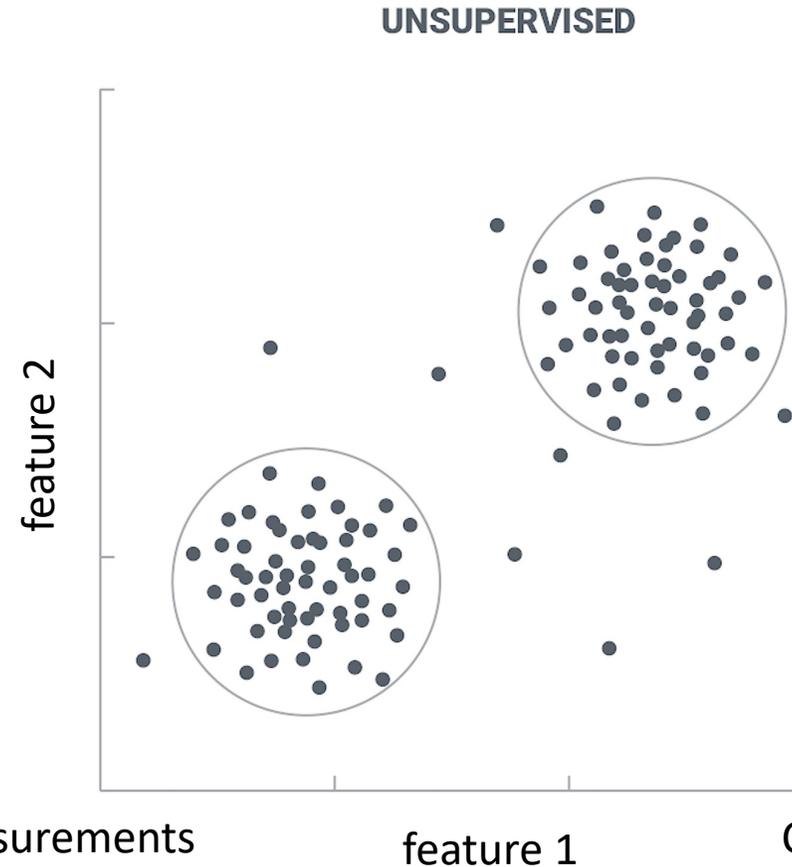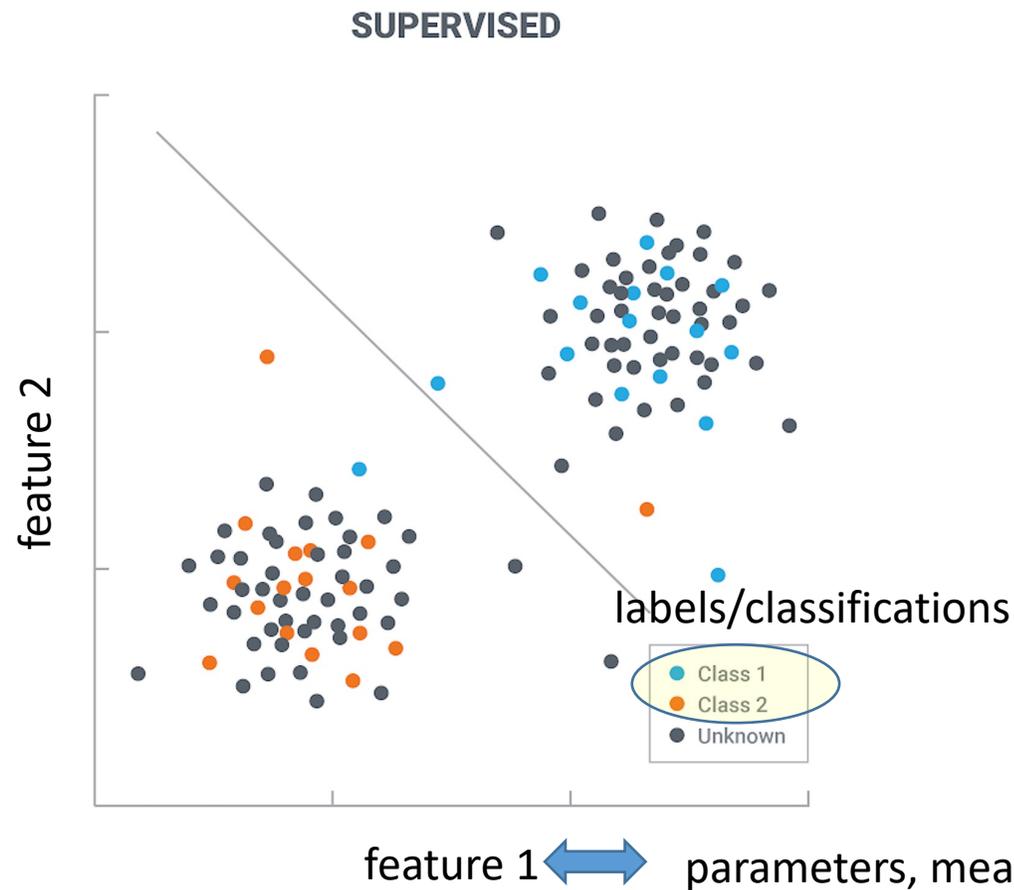


Credit: O. Kargaltsev

# Solution: Machine Learning

ML is needed to enable efficient classification and to reveal dependencies in large datasets **with high dimensionality.**

# Machine Learning Basics

- Supervised Learning vs. Unsupervised Learning (Clustering)



SUPERVISED

UNSUPERVISED

feature 2

feature 1 ⟷ parameters, measurements

labels/classifications

- Class 1
- Class 2
- Unknown

Credit: LAWTOMATED

# Classifying Unidentified X-Ray Sources in the Chandra Source Catalog Using a Multiwavelength Machine-learning Approach

Hui Yang[1] , Jeremy Hare[2,3] , Oleg Kargaltsev[1] , Igor Volkov[1], Steven Chen[1] , and Blagoy Rangelov[4]

[1] Department of Physics, The George Washington University, 725 21st Street NW, Washington, DC 20052, USA; huiyang@gwmail.gwu.edu

[2] NASA Goddard Space Flight Center, Greenbelt, MD 20771, USA

[3] NASA Postdoctoral Program Fellow

[4] Department of Physics, Texas State University, 601 University Drive, San Marcos, TX 78666, USA

## Abstract

The rapid increase in serendipitous X-ray source detections requires the development of novel approaches to efficiently explore the nature of X-ray sources. If even a fraction of these sources could be reliably classified, it would enable population studies for various astrophysical source types on a much larger scale than currently possible. Classification of large numbers of sources from multiple classes characterized by multiple properties (features) must be done automatically and supervised machine learning (ML) seems to provide the only feasible approach. We perform classification of Chandra Source Catalog version 2.0 (CSCv2) sources to explore the potential of the ML approach and identify various biases, limitations, and bottlenecks that present themselves in these kinds of studies. We establish the framework and present a flexible and expandable Python pipeline, which can be used and improved by others. We also release the training data set of 2941 X-ray sources with confidently established classes. In addition to providing probabilistic classifications of 66,369 CSCv2 sources (21% of the entire CSCv2 catalog), we perform several narrower-focused case studies (high-mass X-ray binary candidates and X-ray sources within the extent of the H.E.S.S. TeV sources) to demonstrate some possible applications of our ML approach. We also discuss future possible modifications of the presented pipeline, which are expected to lead to substantial improvements in classification confidences.

# The MUltiWavelength CLASSification Pipeline (MUWCLASS): Training Dataset (TD)

**Supervised** Machine Learning (ML) approach requires training dataset (TD)

| Source Type | Number of CSCv2 sources |
|---|---|
| Active galactic nucleus (AGN) | 1390 |
| Cataclysmic variable (CV) | 44 |
| High-mass star (HM-STAR) | 118 |
| High-mass X-ray binary (HMXB) | 26 |
| Low-mass star (LM-STAR) | 207 |
| Low-mass X-ray binary (LMXB) * | 65 |
| Pulsar and isolated neutron star (NS) | 87 |
| Young stellar object (YSO) | 1004 |
| Total    **8 source classes** | 2941 |

*LMXB also includes non-accreting X-ray binaries (e.g., red-back and black widow systems)

# MUWCLASS: (29) Features/Properties
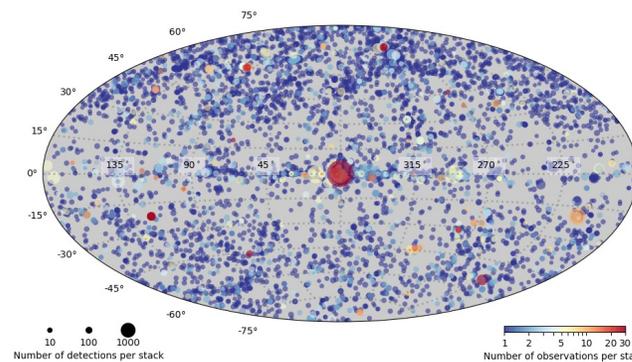
Various Colors

Mid-Infrared from WISE

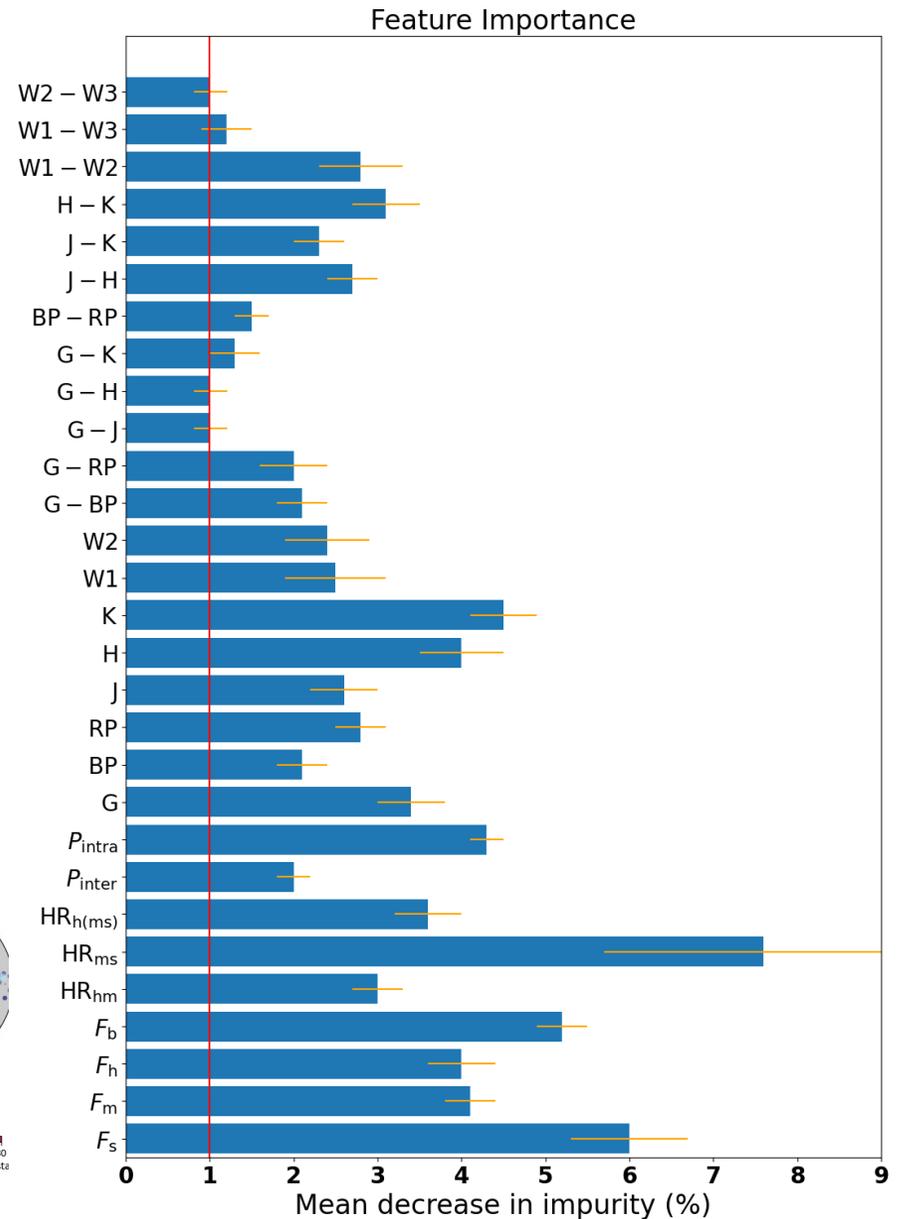Near-Infrared from 2MASS

Optical from Gaia

X-ray variability

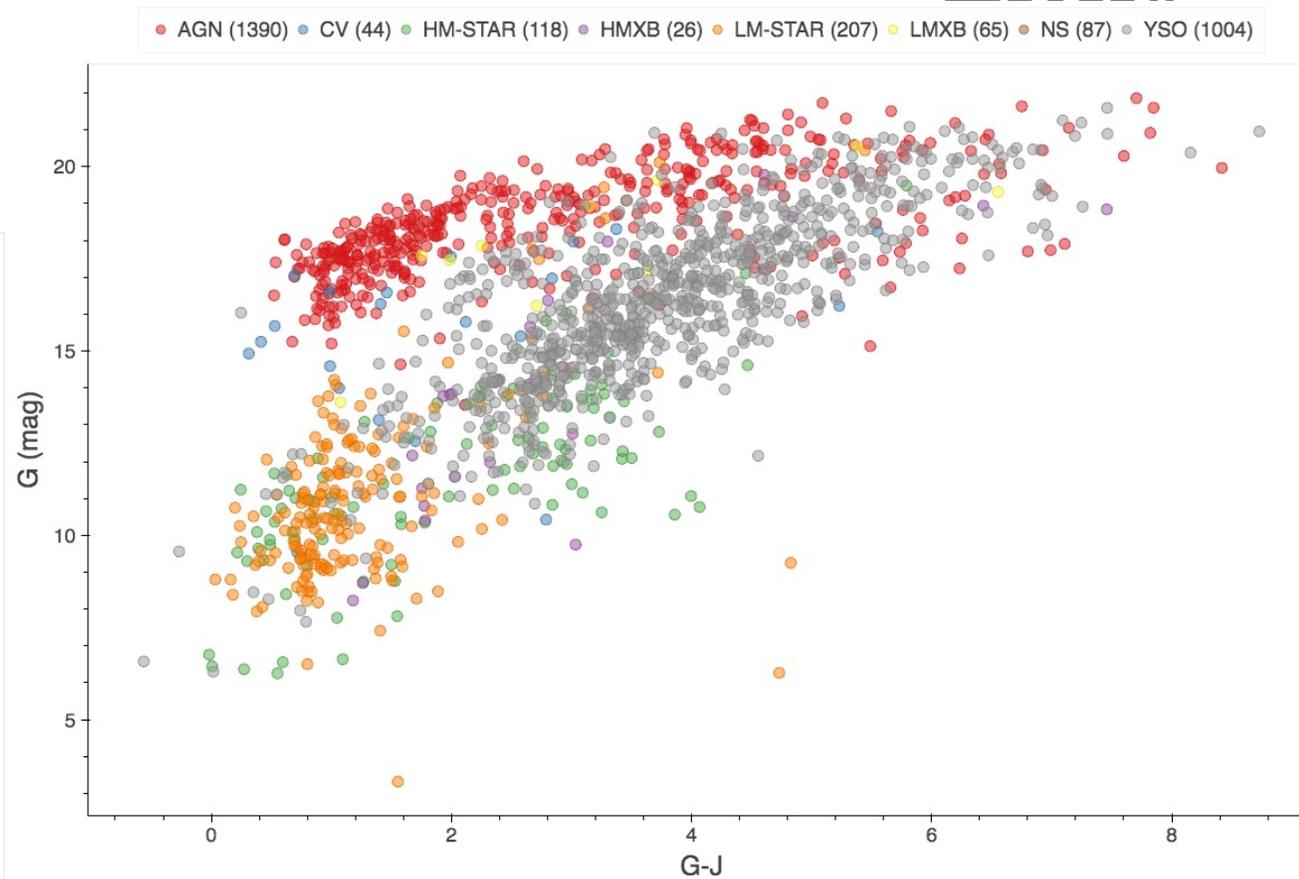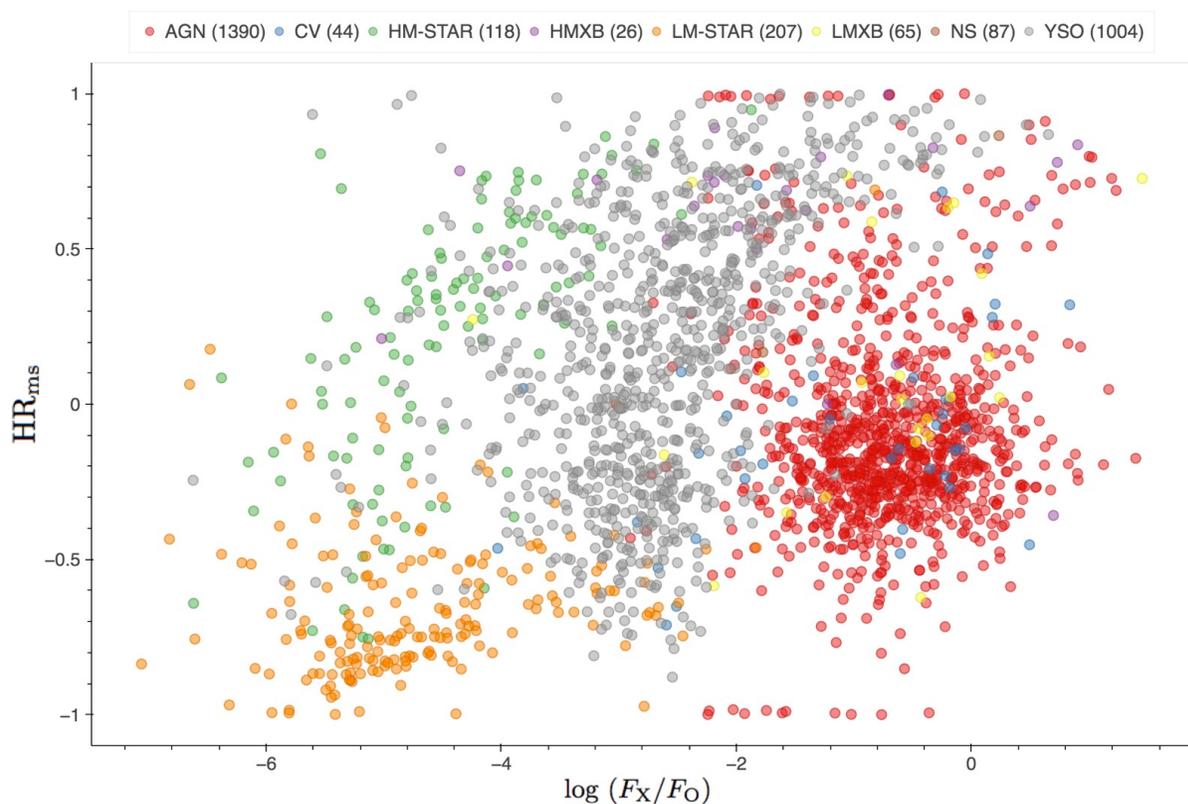Hardness ratios

Chandra X-ray fluxes

CSC v2.0 Detection Map from CXO

# 2D representations of our TD

Typically traditional classification consists
of using multi-wavelength parameter plots to separate
source classes



Hard to comprehend more than 2 or 3 dimensions for human

Explore the TD yourself using the visualization GUI at
**https://home.gwu.edu/~kargaltsev/XCLASS/** (Yang et al. 2021)

# MUWCLASS pipeline flow chart



**Unclassified sources**

**Training Data**

MC sampling ?

Yes → Feature Uncertainty?

Sample features from their PDFs

No

Apply extinction/absorption on AGNs in TD

Standardize the data

Oversample TD with SMOTE

Replace missing data with -100

Train the Random Forest classifier

Apply classifier on unclassified sources to obtain probabilities belonging to any of 8 classes

scikit-learn

*Machine Learning in Python*

Getting Started | Release Highlights for 0.23 | GitHub

- Simple and efficient tools for predictive data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
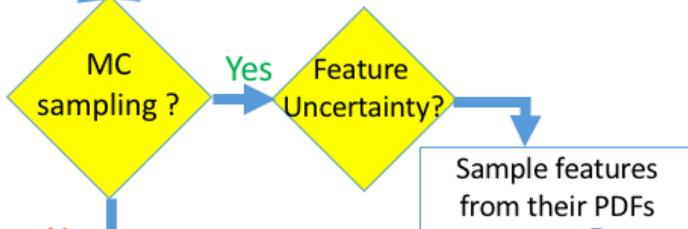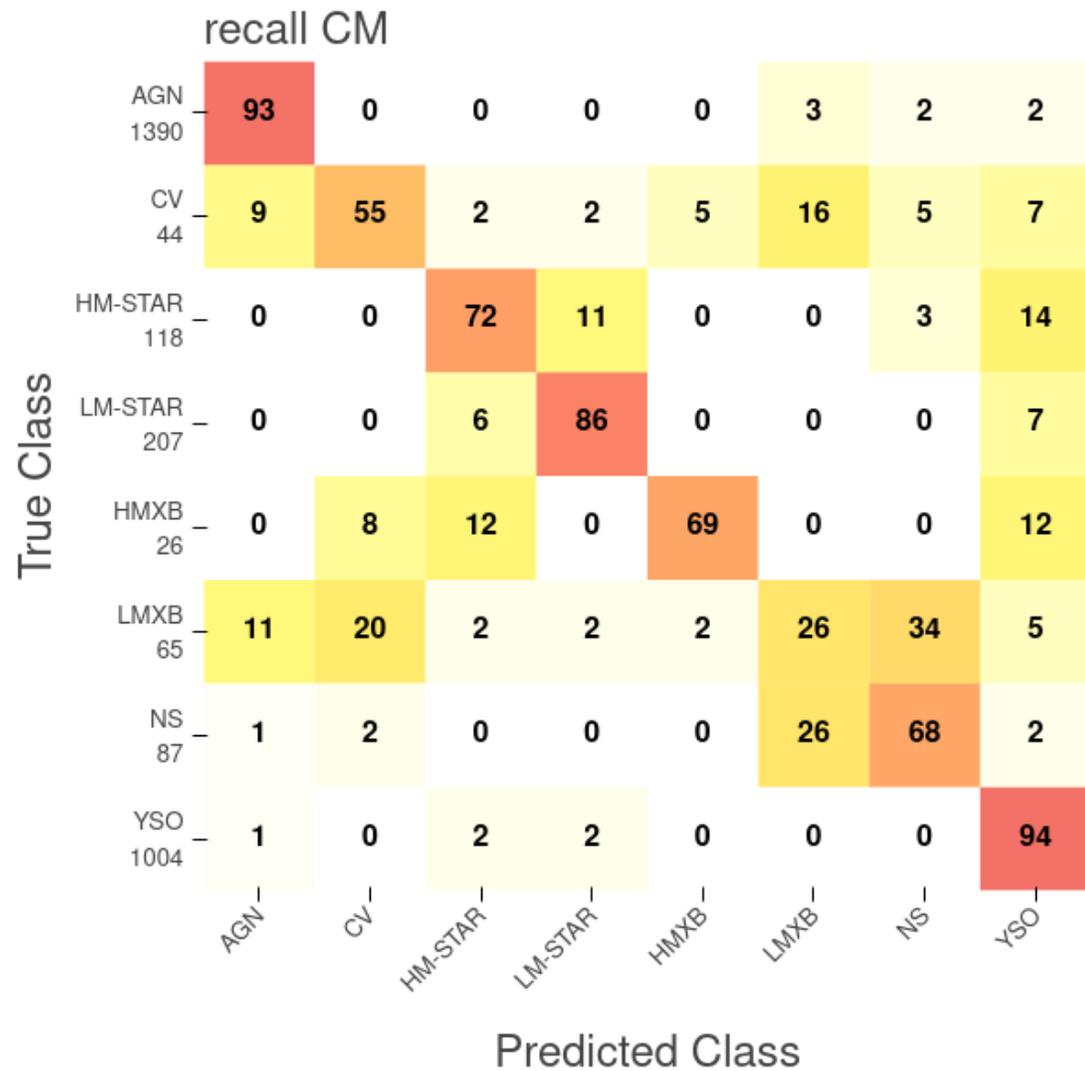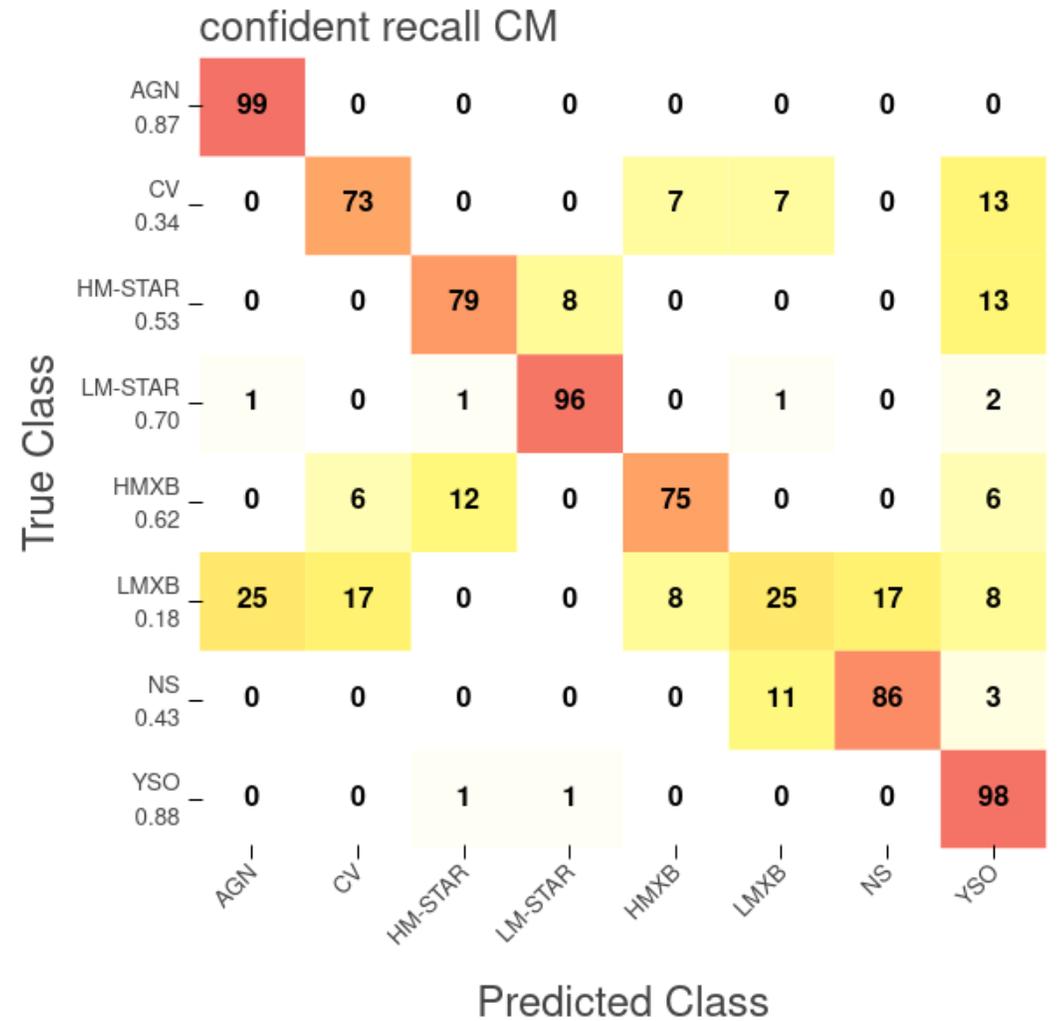- Open source, commercially usable - BSD license

https://scikit-learn.org/stable/



X-ray flux PDF

X-ray flux distribution/$10^{-14}$ erg cm$^{-2}$s$^{-1}$

Feature measurement uncertainties are taken into account by sampling from feature PDFs via Monte-Carlo.

AGN 1390

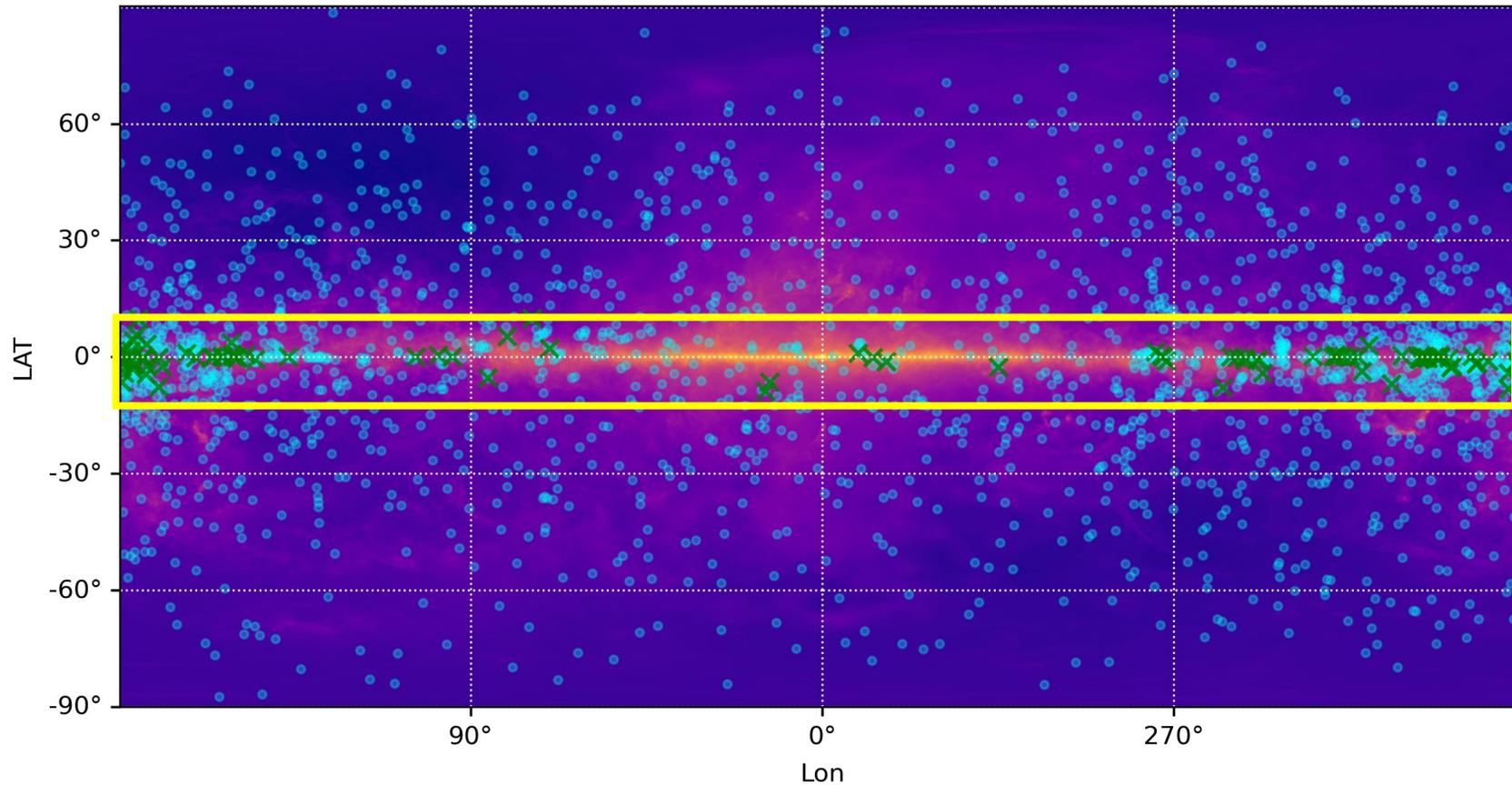Confusion Matrix (true vs. predicted class)

Average Accuracy: 88.6%
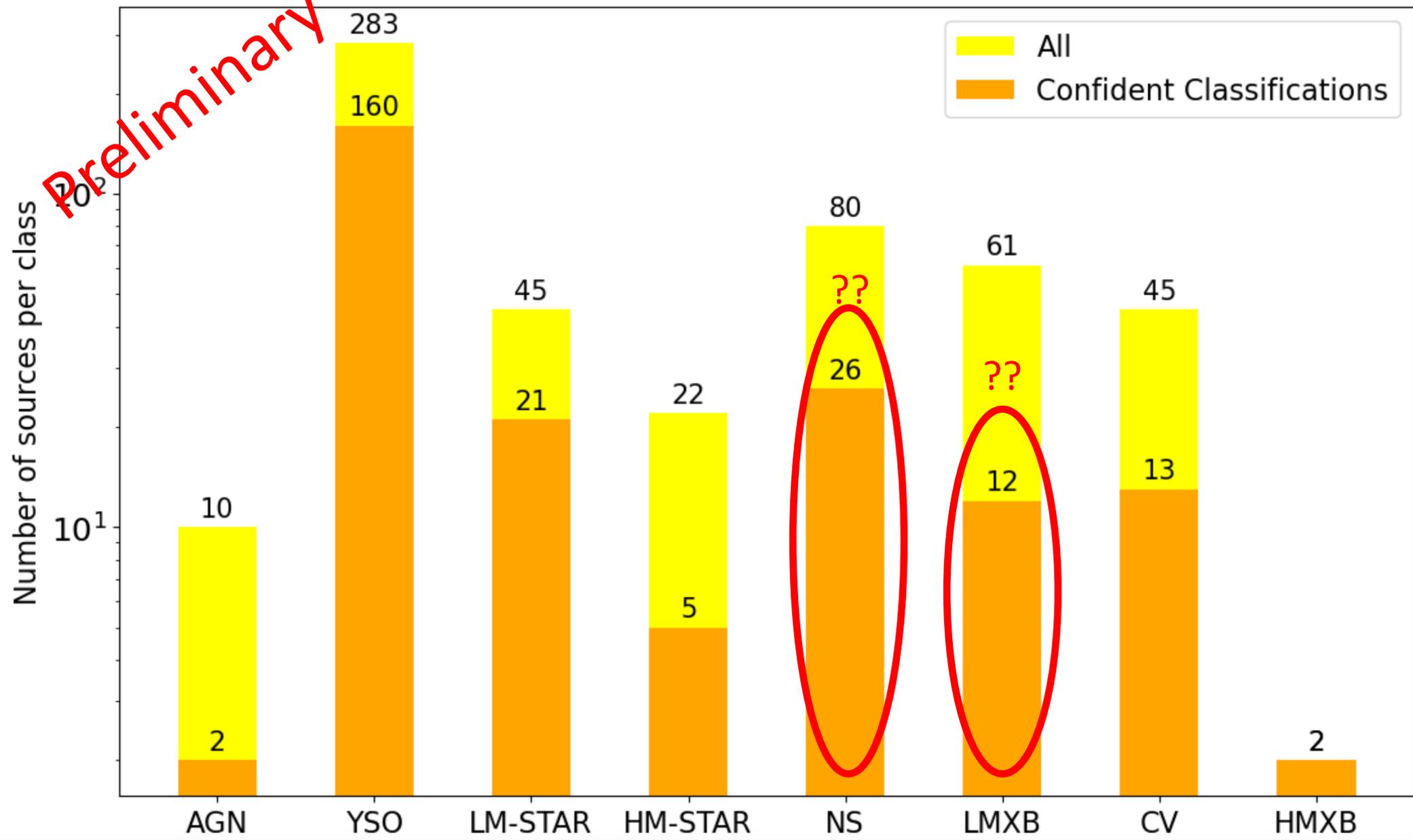
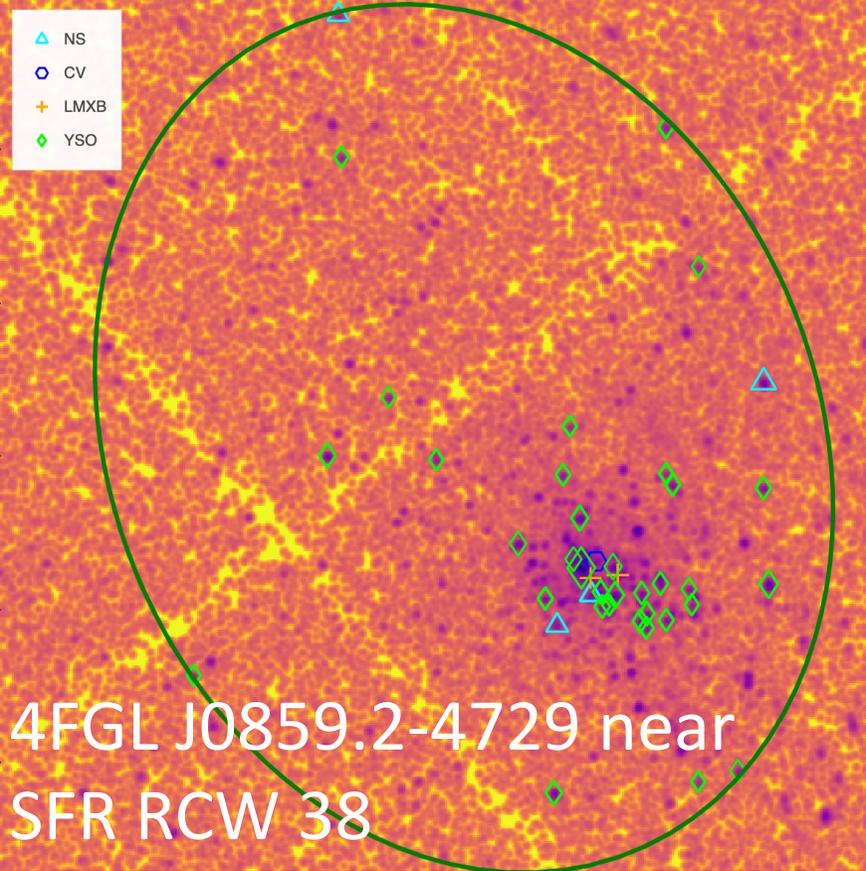Confusion Matrix of **Confident** Classifications

Average Accuracy: 97%

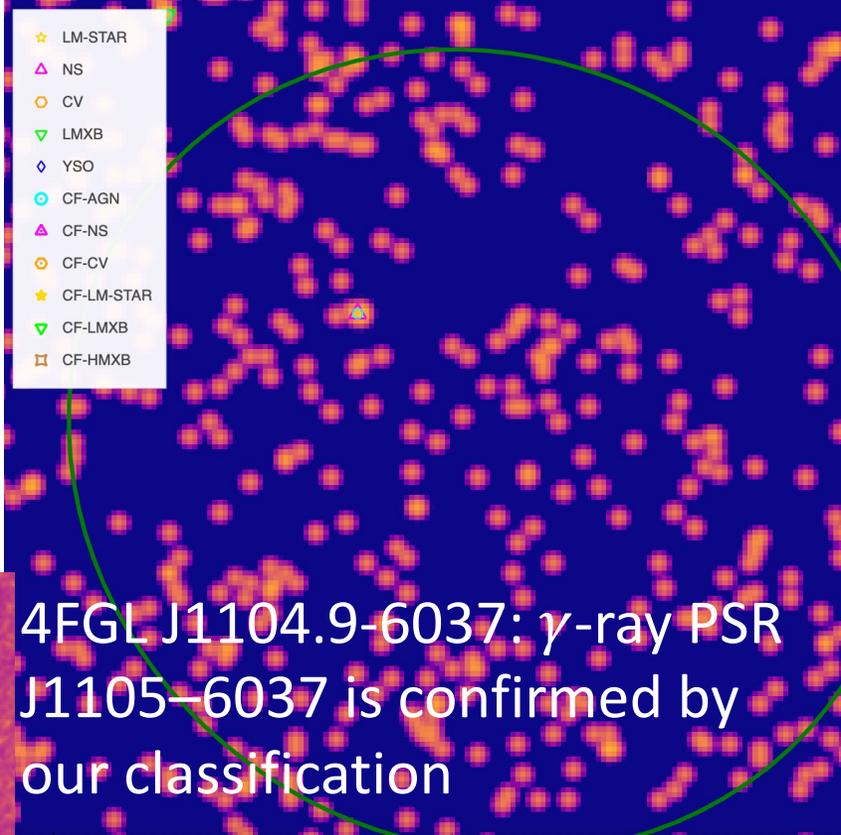# Classification of CSCv2 Sources within unassociated 4FGL-DR3 fields



- 37 unassociated 4FGL-DR3 sources within |b|<10° with ≥5 CSCv2 sources within their error ellipses;

- 548 significant (X-ray S/N>5) CSCv2 sources within these 37 FGL sources.

Legend (top-left panel):
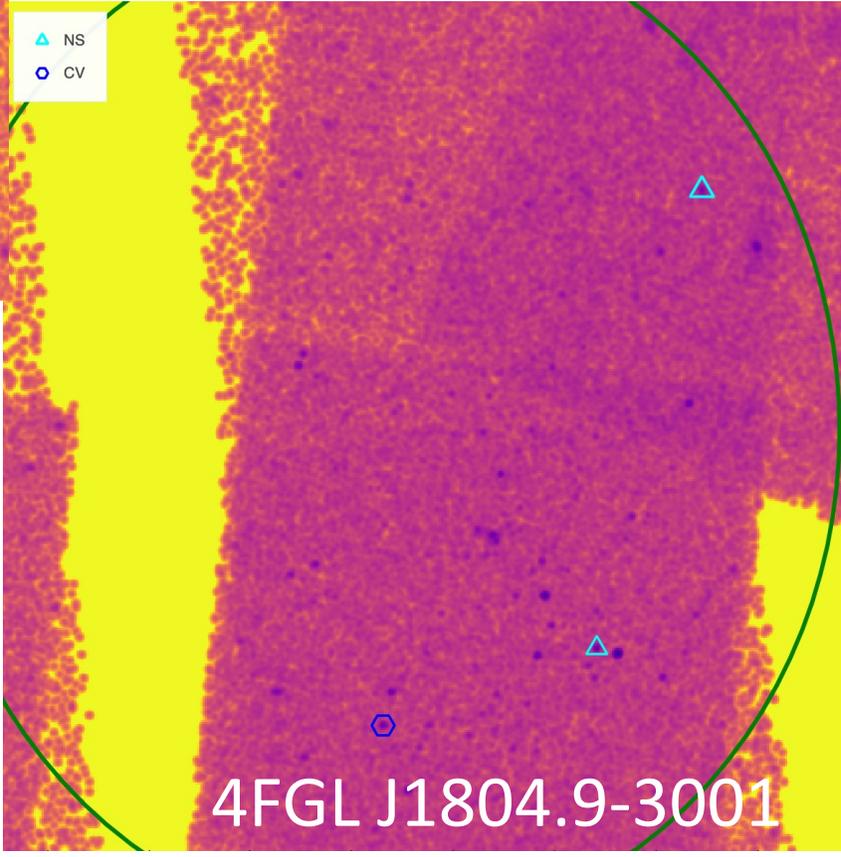- NS
- CV
- LMXB
- YSO

Legend (top-right panel):
- LM-STAR
- NS
- CV
- LMXB
- YSO
- CF-AGN
- CF-NS
- CF-CV
- CF-LM-STAR
- CF-LMXB
- CF-HMXB

Legend (center panel):
- NS
- CV

4FGL J0859.2-4729 near SFR RCW 38

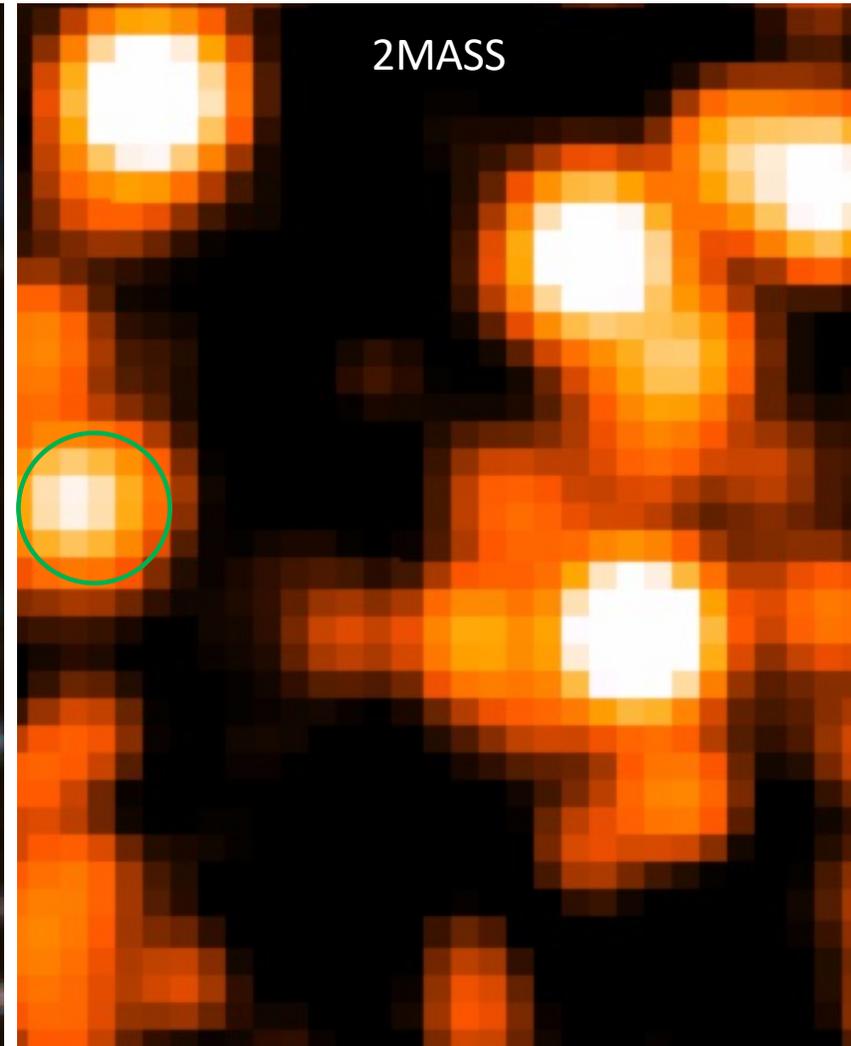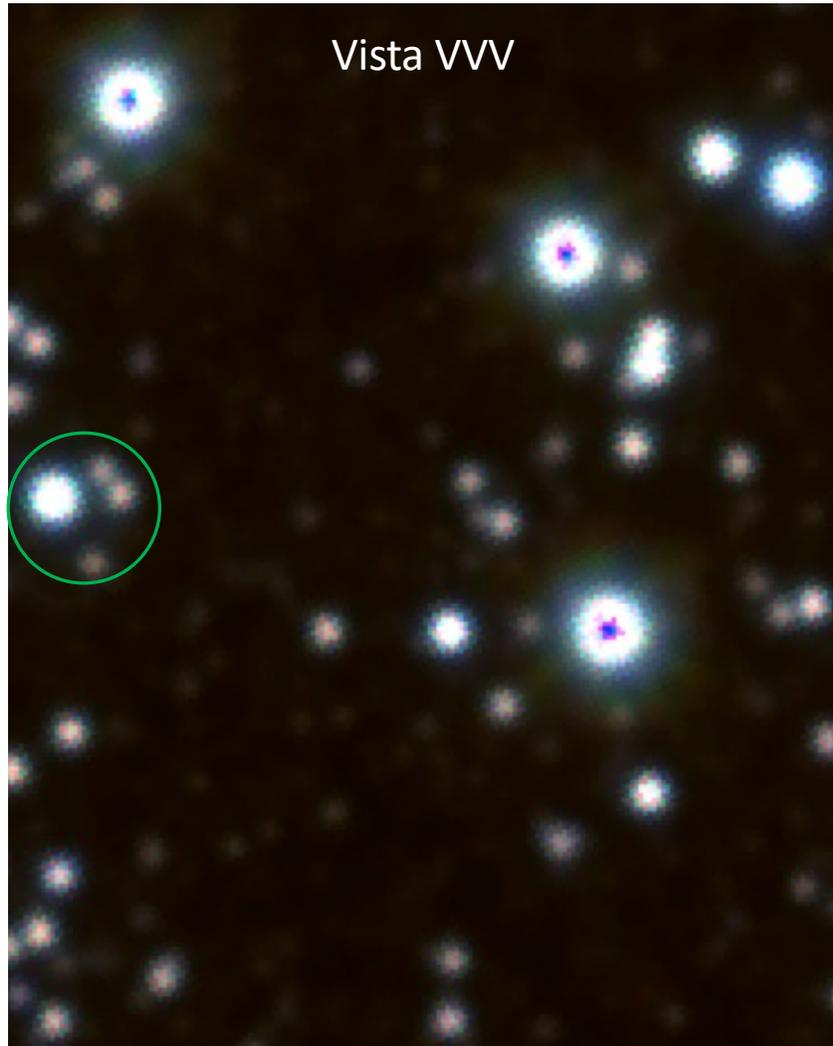4FGL J1104.9-6037: $\gamma$-ray PSR J1105−6037 is confirmed by our classification

4FGL J1804.9-3001

# Issues and Biases

- NS virtually all too faint to be detected by multi-wavelength surveys used in training dataset

- Many LMXBs also have counterparts too faint to be detected by multi-wavelength surveys used in training dataset, hence they become confused with the NS class

- Sources too faint to be detected by these MW surveys (e.g., M-dwarfs, absorbed AGN) will be preferentially classified as NS/LMXBs

# Future Improvements: deeper surveys

- Update to more sensitive surveys (e.g., Pan-STARRs, DECaps, Vista VVV)

# Future Improvements: New Multi-wavelength Features

### ASKAP



https://www.atnf.csiro.au/projects/askap/index.html

### Gaia



https://upload.wikimedia.org/wikipedia/en/0/01/Gaia_spacecraft.jpg
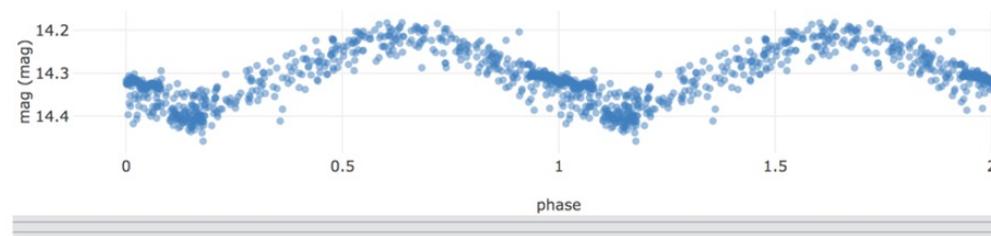
### ZTF



https://www.ipac.caltech.edu/project/ztf

- Inclusion of new radio surveys:

- Australian SKA Pathfinder Telescope (ASKAP)

- VLA All-sky Survey (VLASS)

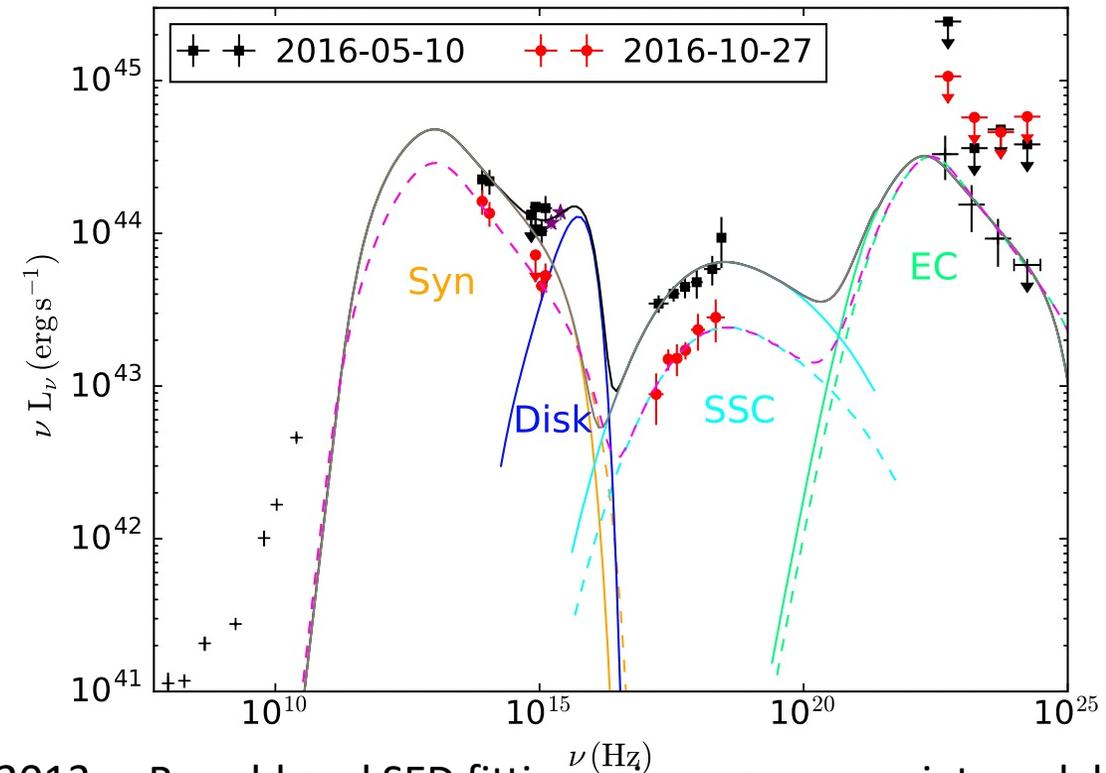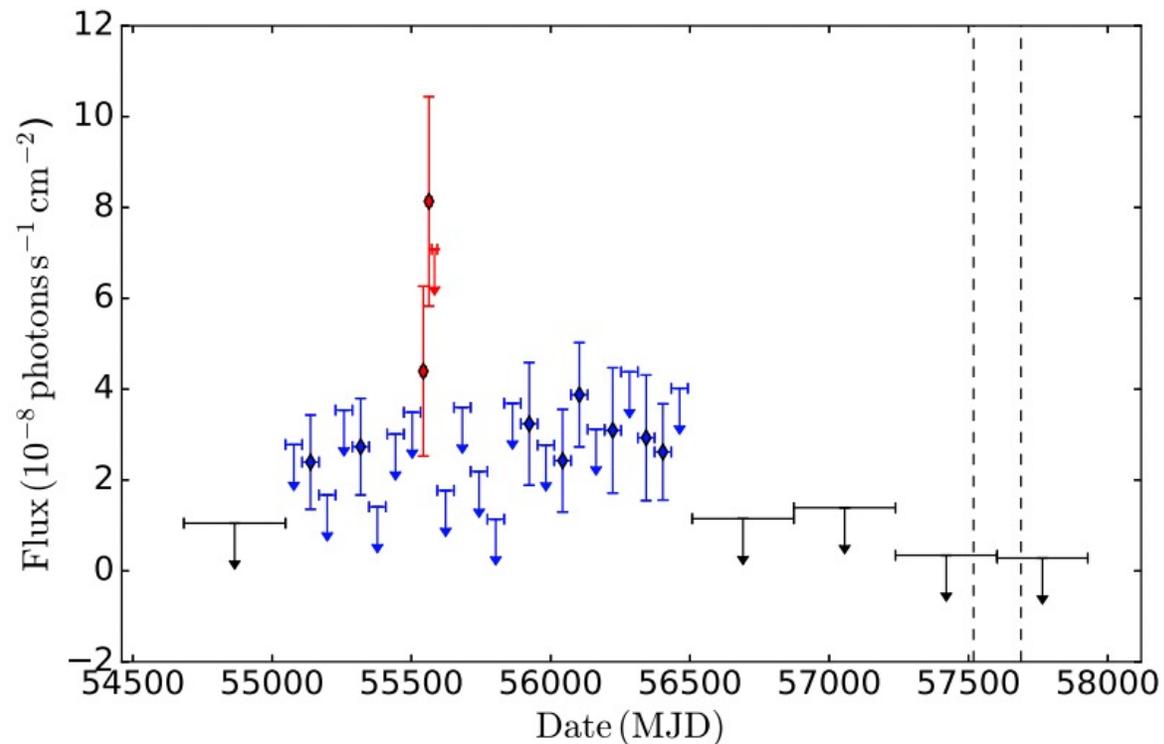- MeerKAT source catalog

- Distances and proper motions from Gaia DR3



- Large field of view optical time domain surveys:

- ZTF

- TESS

- VCRO

# SDSS J211852.96−073227.5: a new γ-ray flaring narrow-line Seyfert 1 galaxy

Hui Yang,[1,2]★ Weimin Yuan,[1,2]★ Su Yao,[3,4] Ye Li,[3,4] Jin Zhang,[1] Hongyan Zhou,[5,6] S. Komossa,[7] He-Yang Liu[1,2] and Chichuan Jin[1]

Fermi-LAT light curve shows monthly γ-ray flaring activities during 2009-2013.   Broad-band SED fitting using a one-zone jet model.