# Active Anomaly Detection Tutorial

**Konstantin Malanchev malanchev@cmu.edu**
**LINCC Frameworks / Carnegie Mellon University**

**AISSAI Anomaly Detection Workshop, March 2024, Clermont-Ferrand**
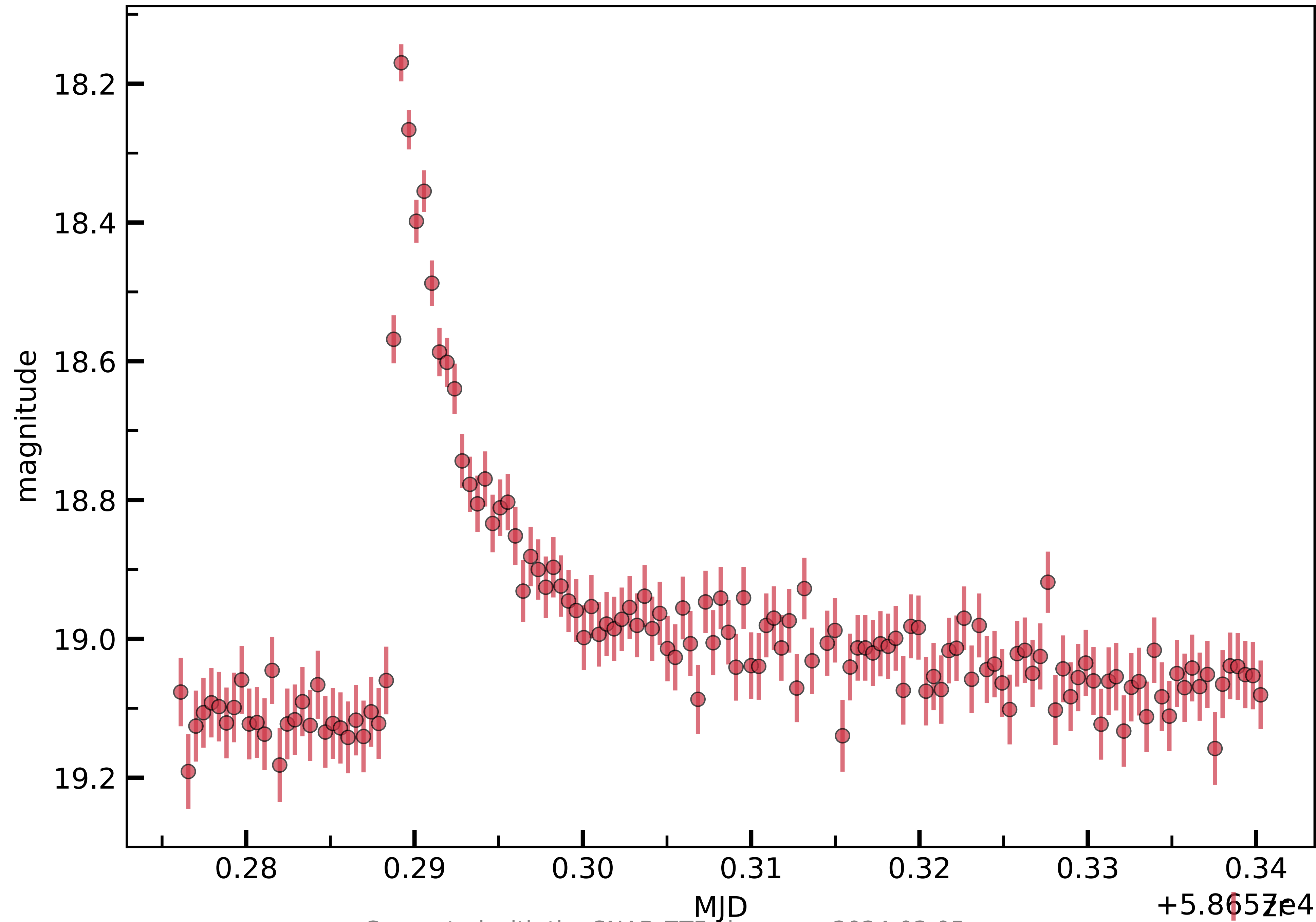
# Before we start
## What we will need

- Google account for Google Colab, OR
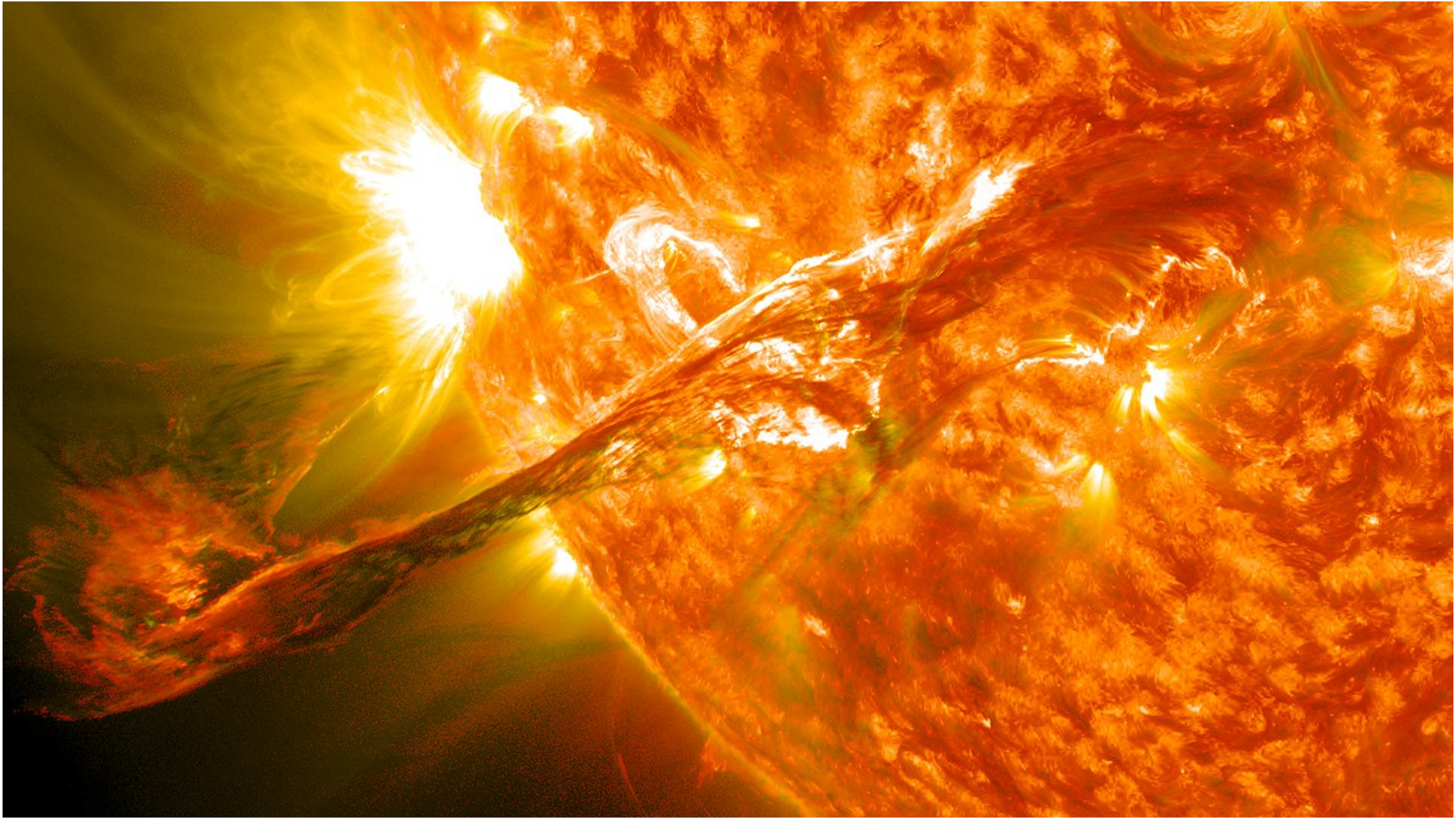
- Python 3.9-3.11 and Jupyter with virtual environment

Check notebook links on the workshop website, QRs are following
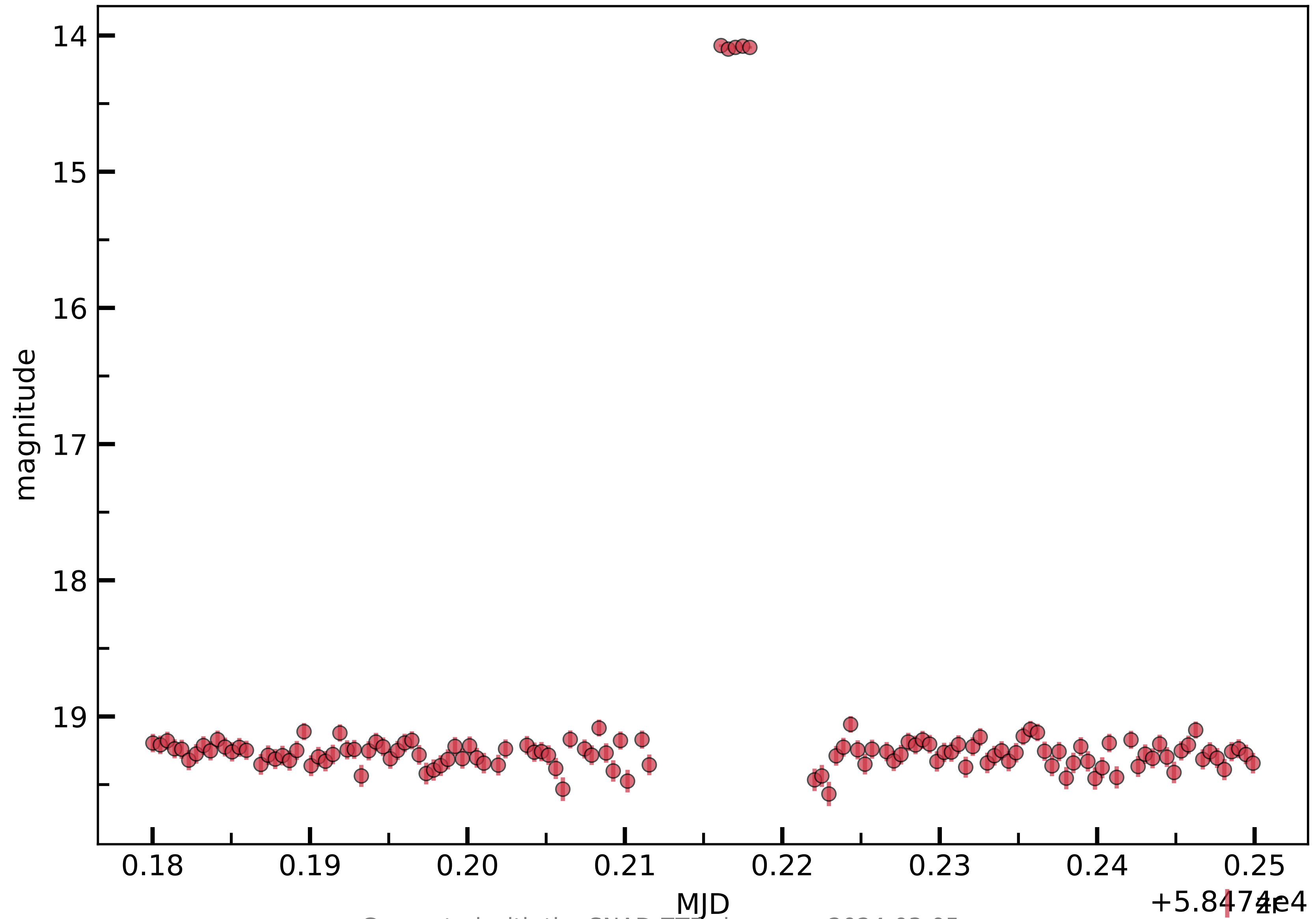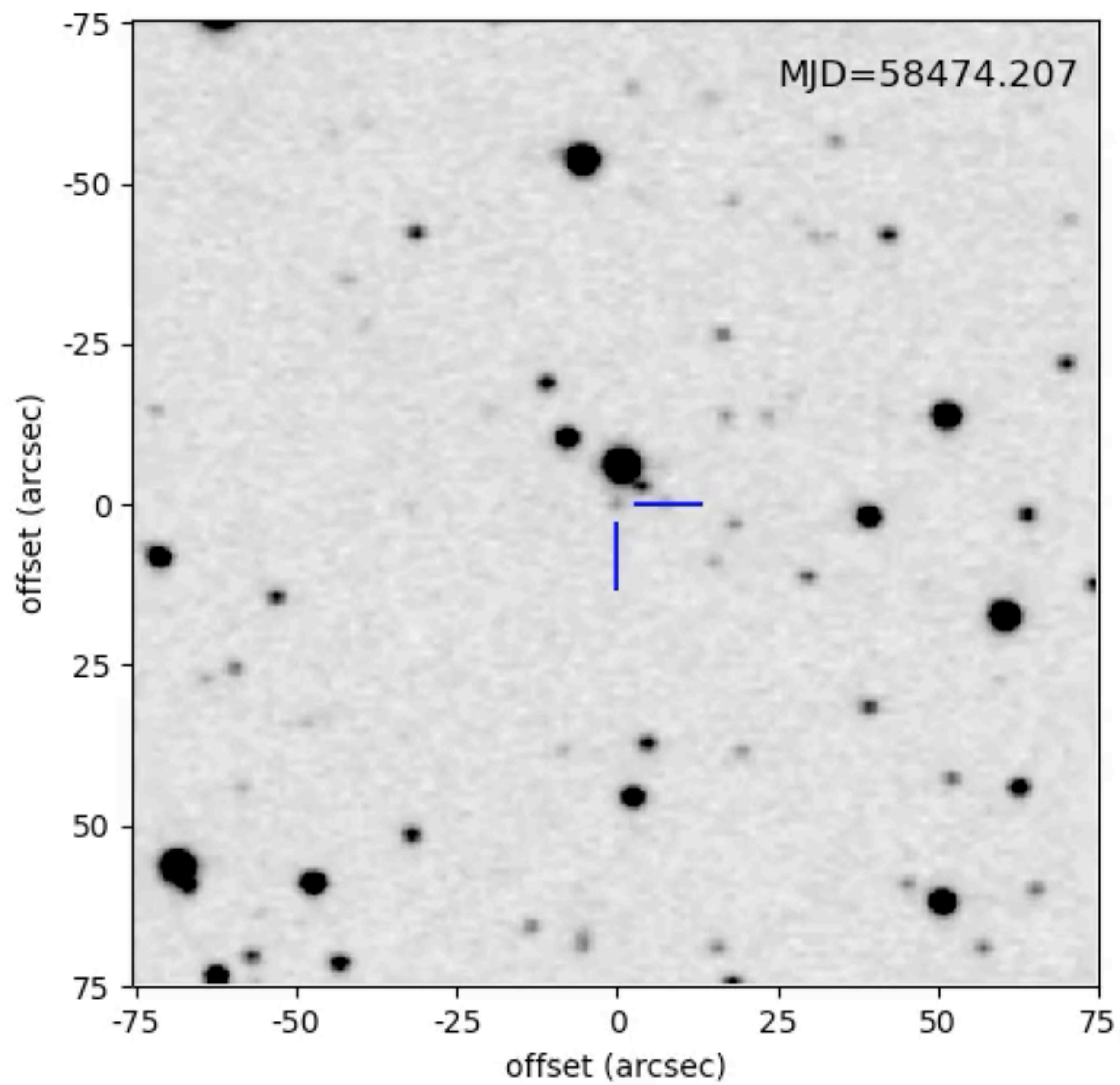
# Why do we want to go active?

637212400010948

Generated with the SNAD ZTF viewer on 2024-03-05

807203300039547

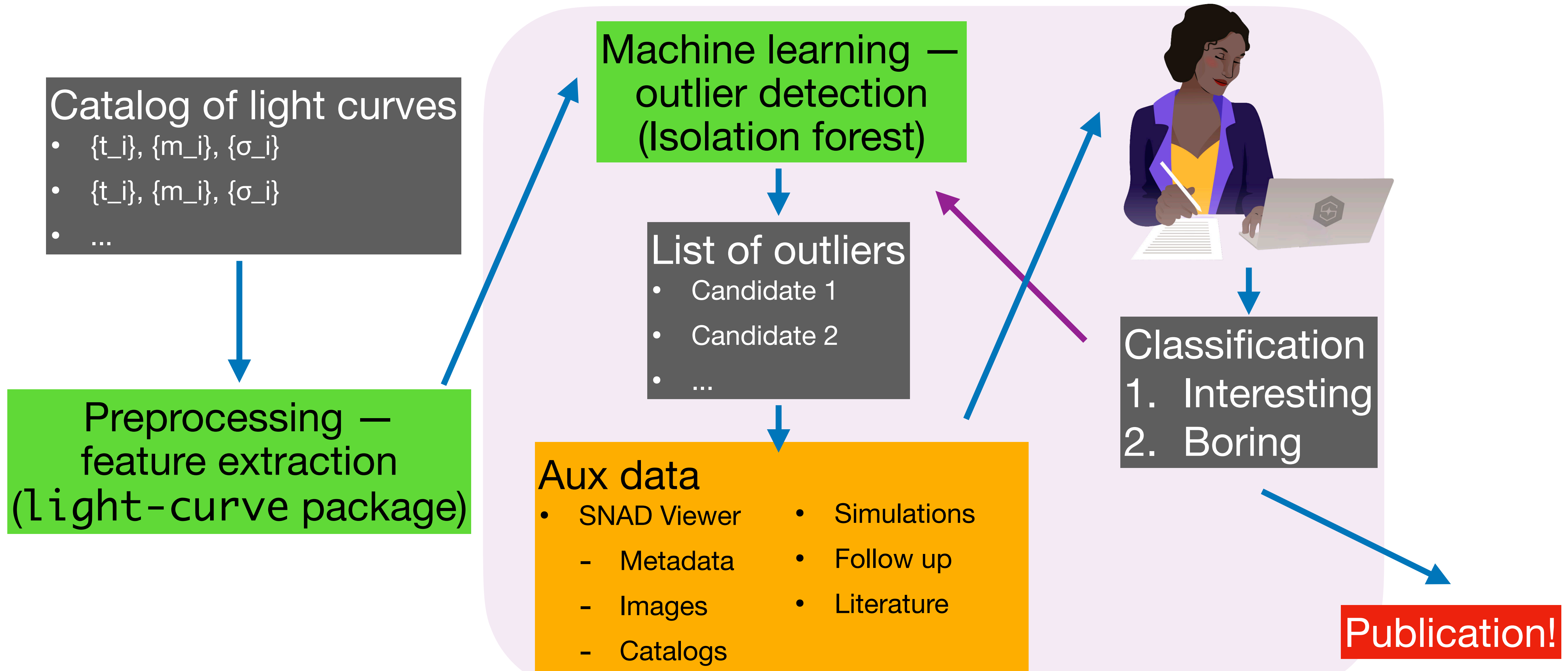Generated with the SNAD ZTF viewer on 2024-03-05

MJD=58474.207

# Pipeline

# Anomaly Detection for Light Curves
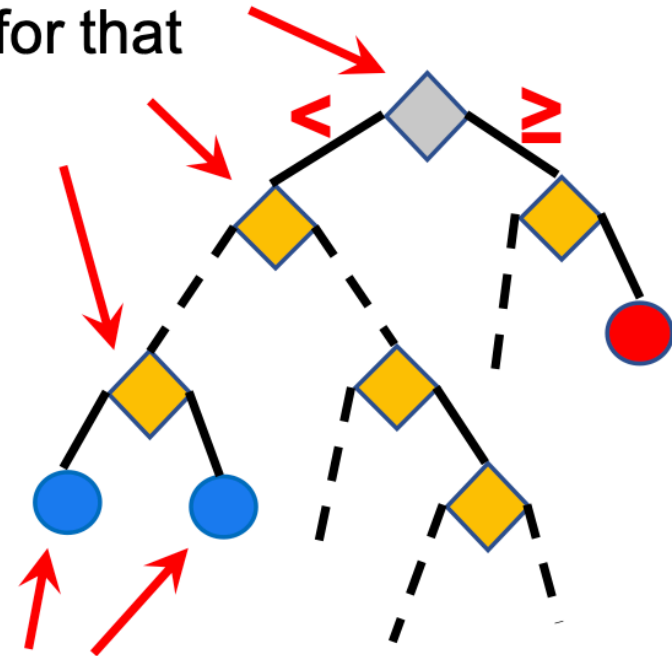## SNAD papers: Pruzhinskaya+19,22; Ishida+19; KM+21; Aleo+22.

# Isolation Forest and Pineforest
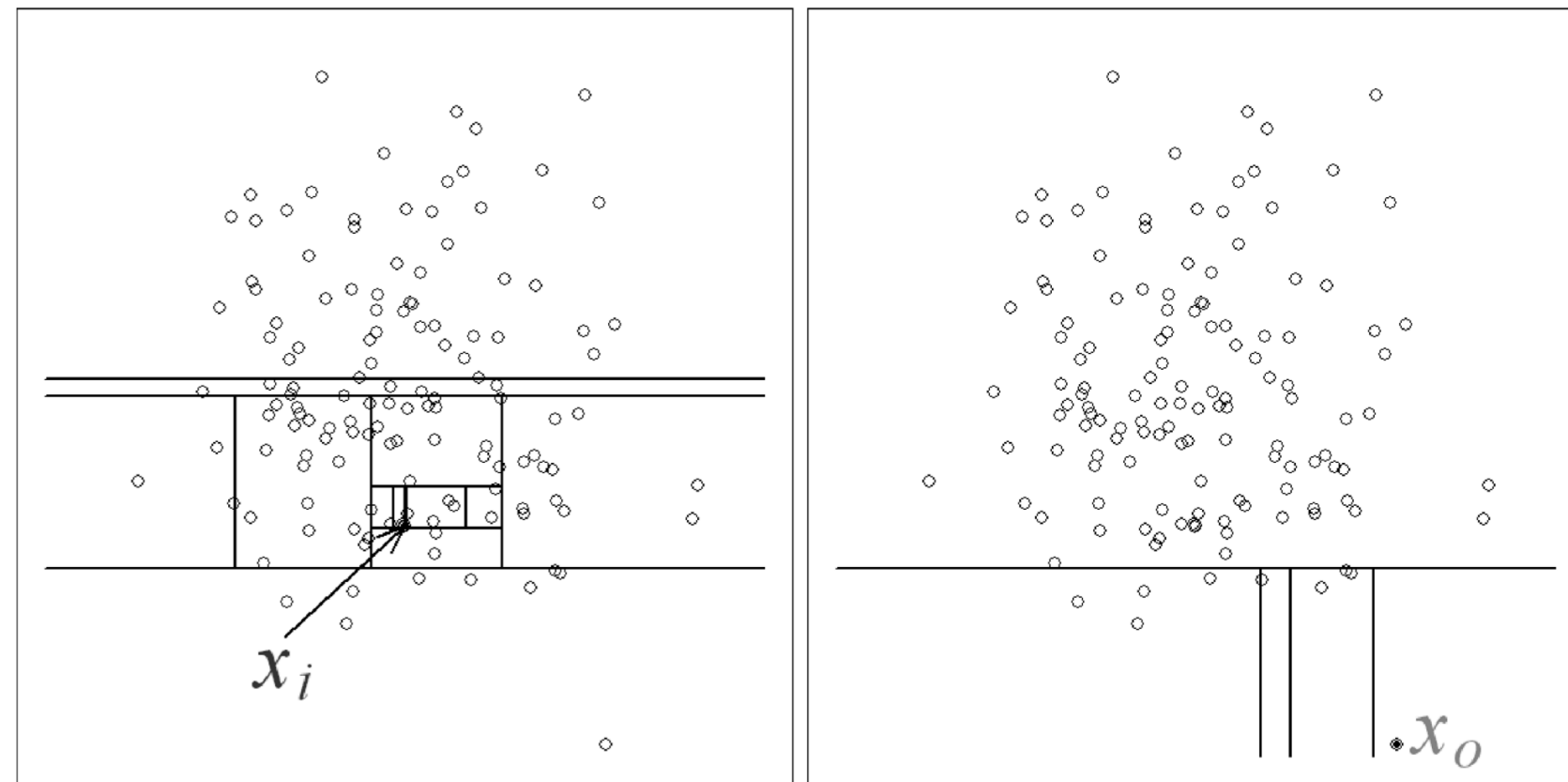
# Isolation Forest

## iTree

Select a random feature at each node, and a random split point for that feature

Shallower leaf nodes have higher anomaly scores, whereas, deeper leaf nodes have lower anomaly scores.
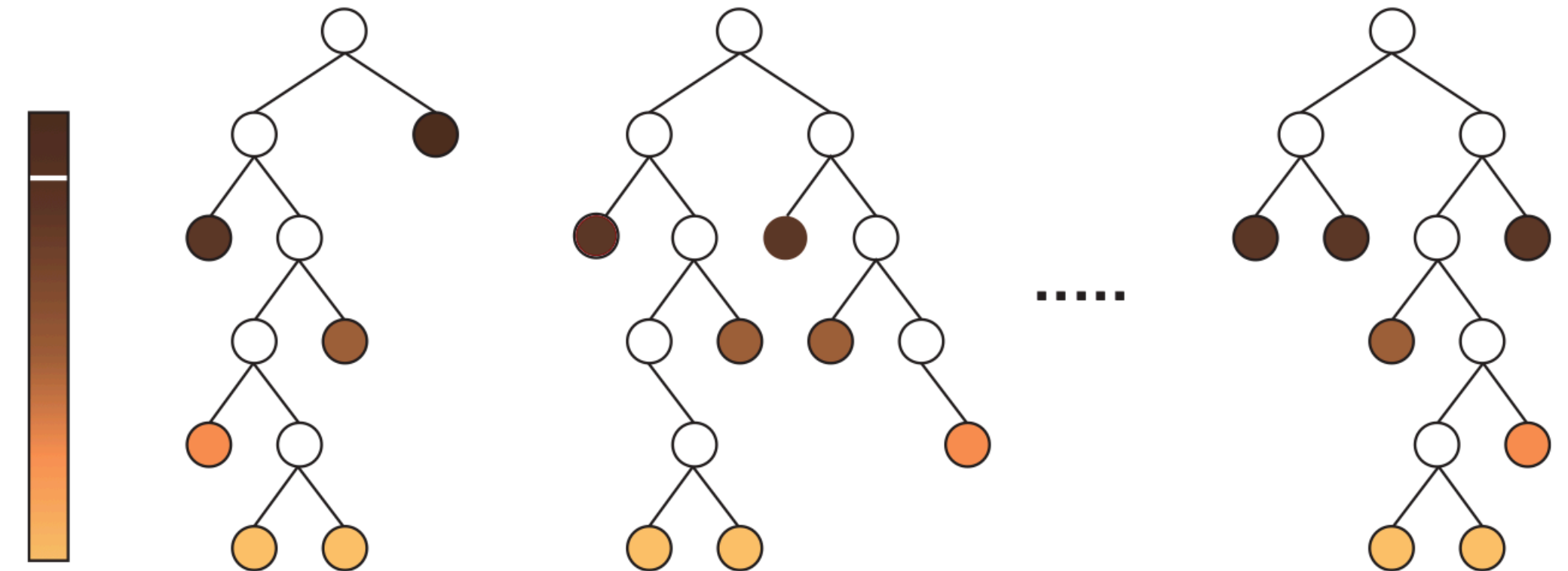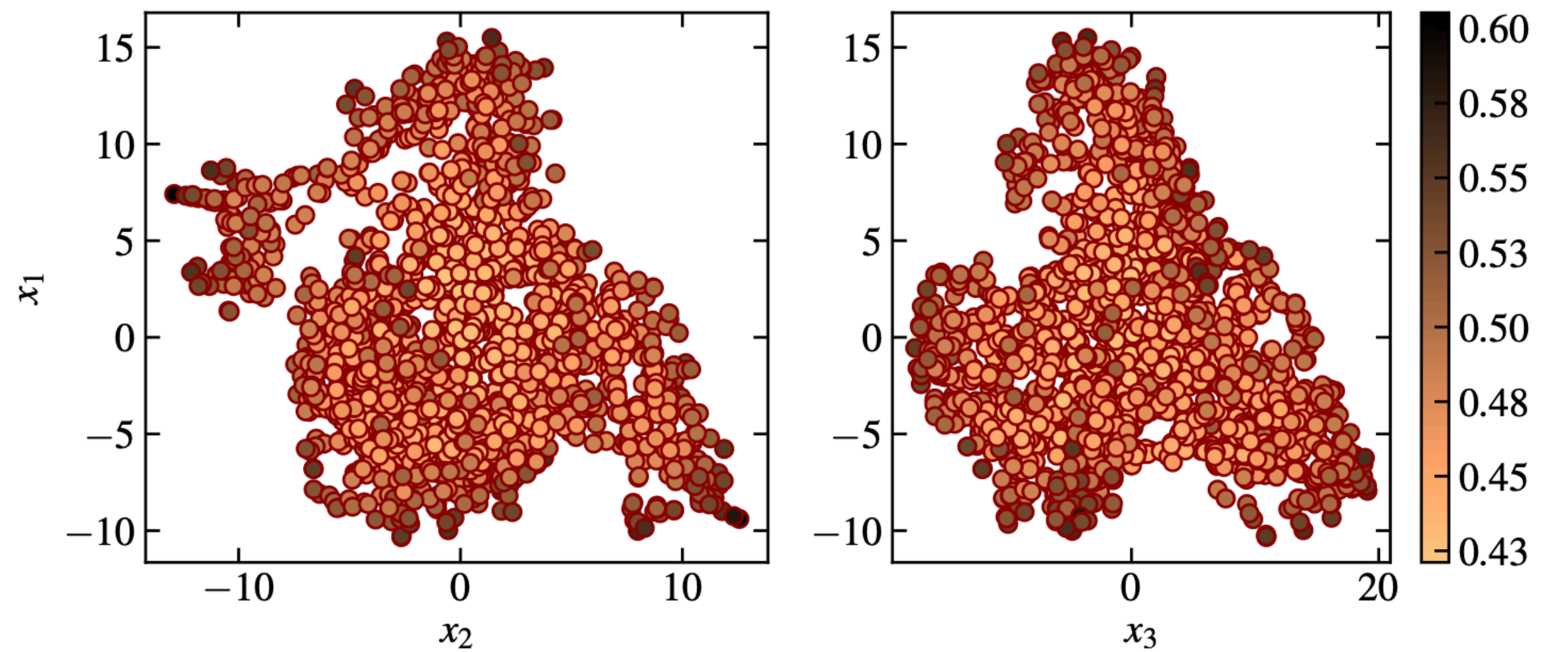
Leaf instance
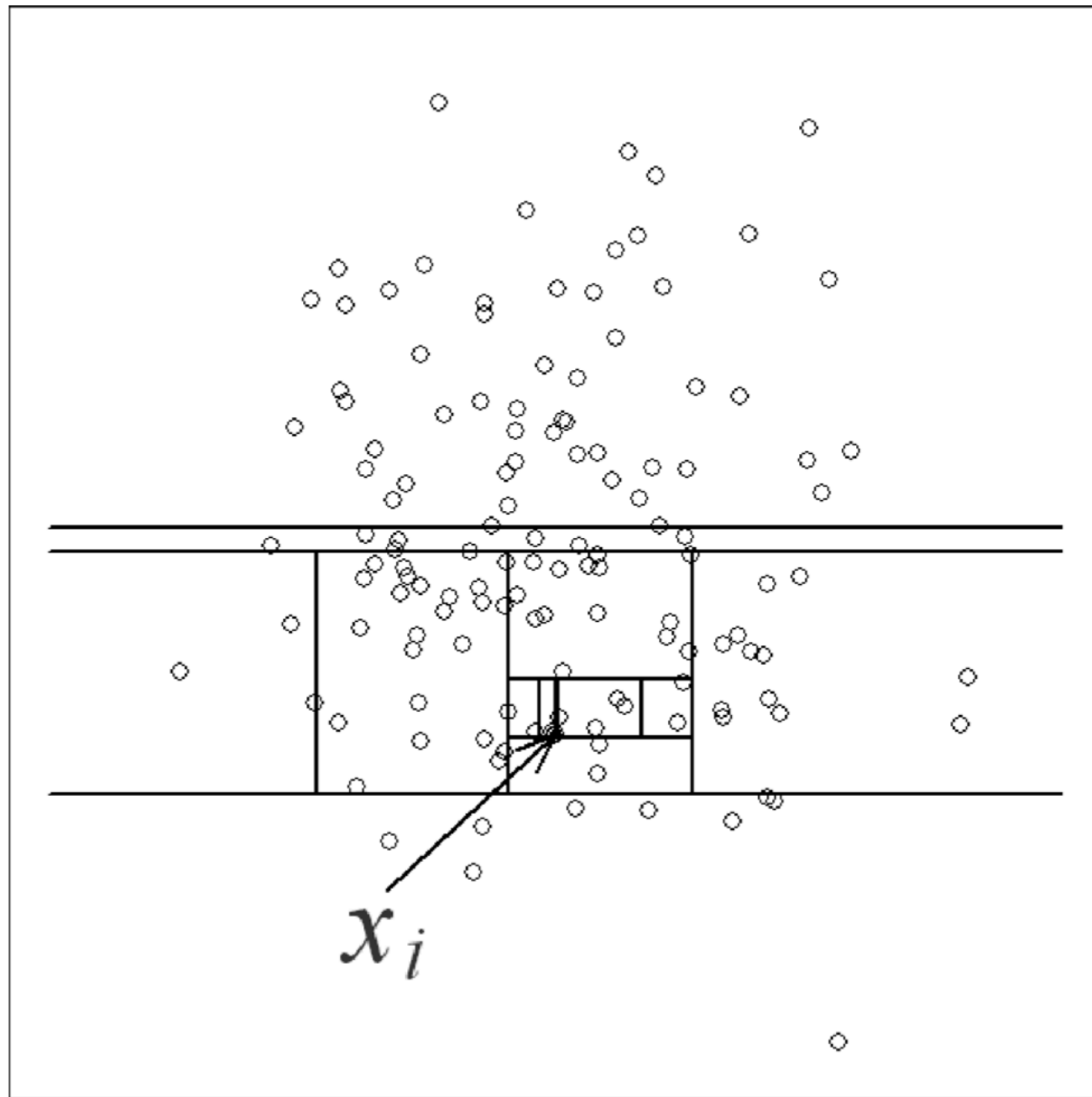
arXiv:1708.0944



(a) Isolating $x_i$

(b) Isolating $x_o$



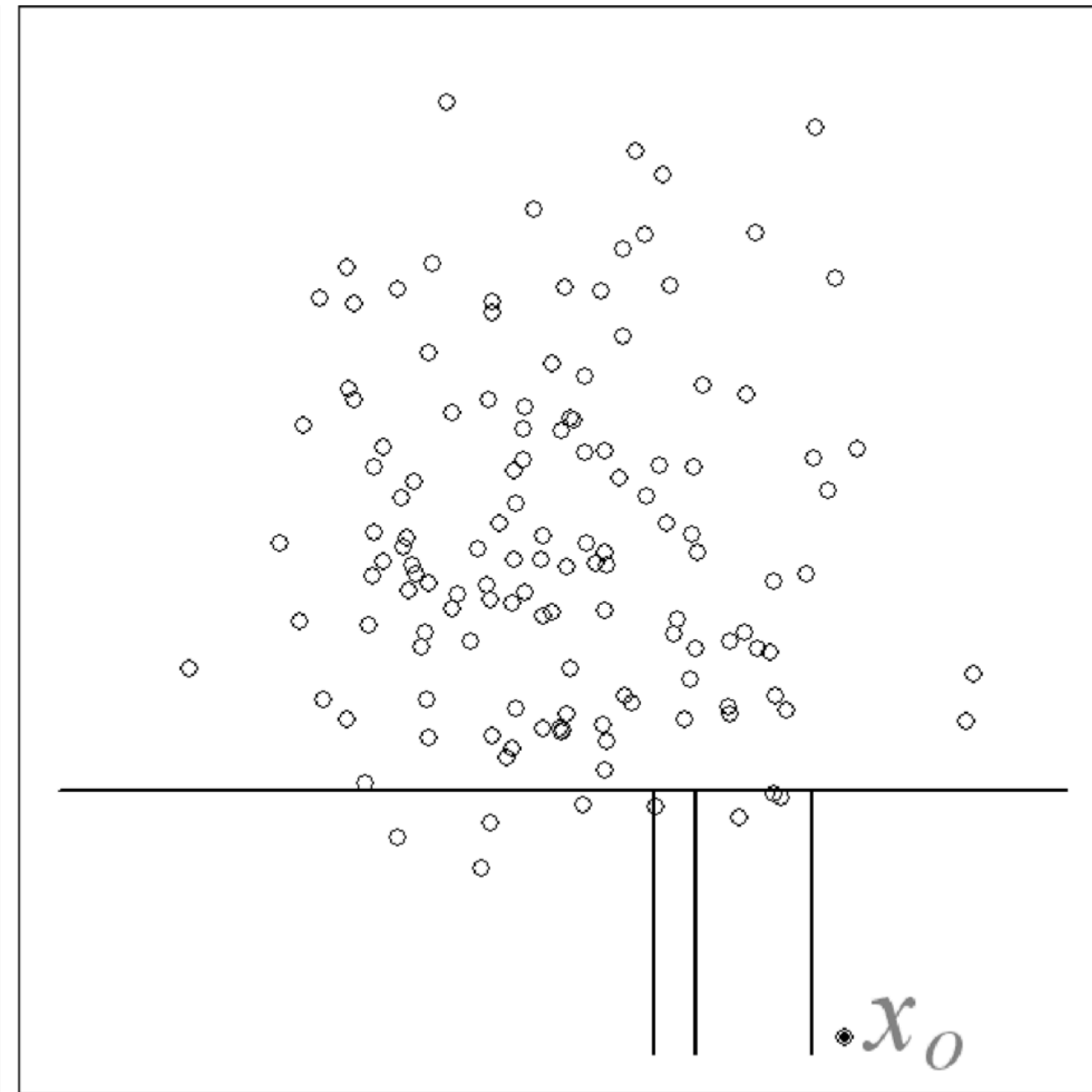## Darker is more anomalous



Liu+ 2008, Liu+ 2012

arXiv:1905.11516

# Isolation Tree



(a) Isolating $x_i$

(b) Isolating $x_o$

$$c(\psi) = \begin{cases} 2H(\psi - 1) - 2(\psi - 1)/\psi & \text{for } \psi > 2, \\ 1 & \text{for } \psi = 2, \\ 0 & \text{otherwise,} \end{cases}$$
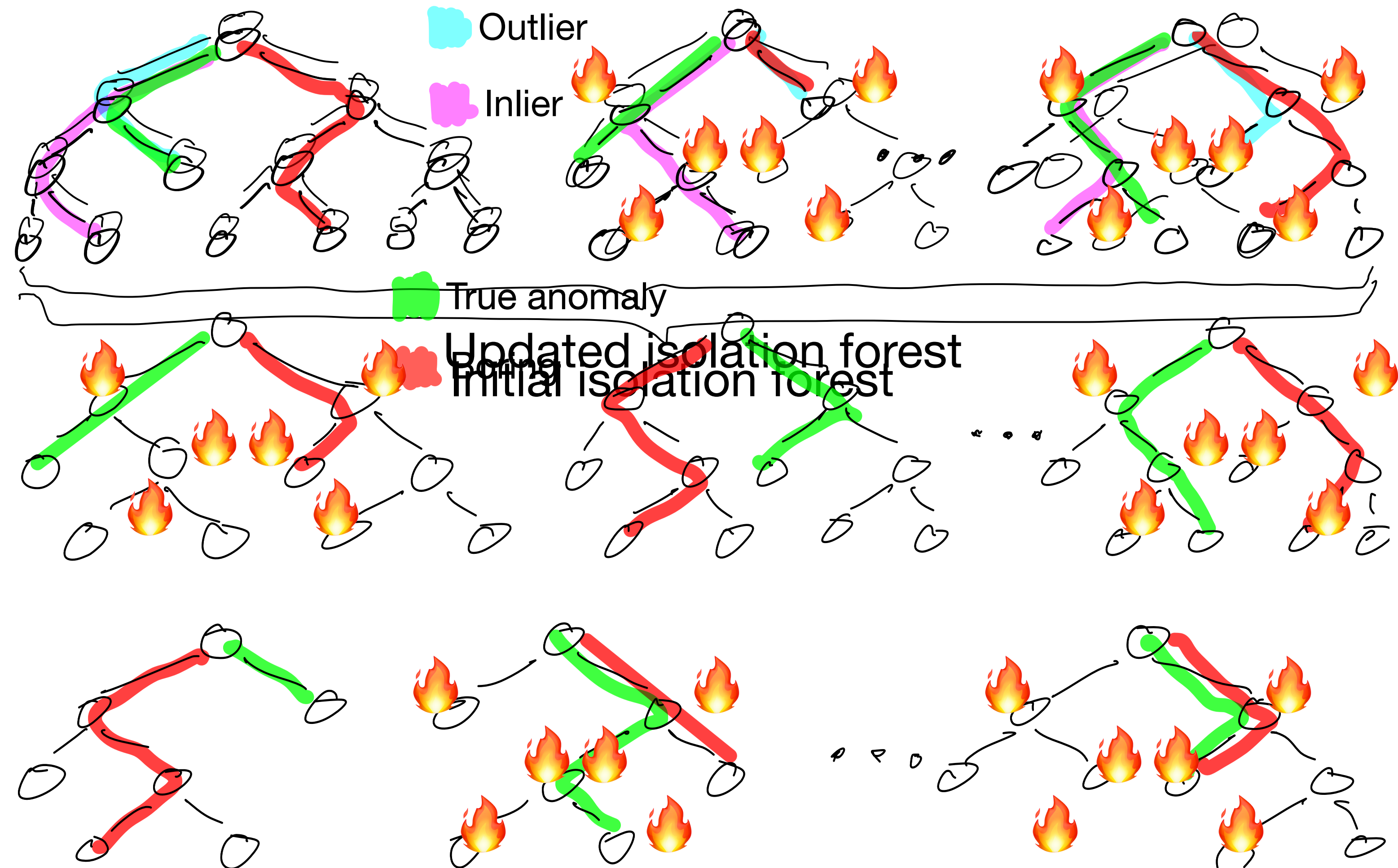
$$s(x, \psi) = 2^{-\frac{E(h(x))}{c(\psi)}},$$

**Liu et al 2008, Liu et al 2012**

# Active Anomaly Detection with Pineforest
## Based on Isolation Forest (Liu+08) and inspired by AAD (Das+17)

1. Build an isolation forest

2. Select the best outlier from the unlabeled data

3. Ask the expert to classify

4. Build more trees

5. Rank the trees with labeled data

6. Select the best trees and prune the rest

7. Go to 2.



Outlier

Inlier

True anomaly

Updated isolation forest

Initial isolation forest

# Tools we are going to use

# Coniferest Package

**Docs: https://coniferest.snad.space**
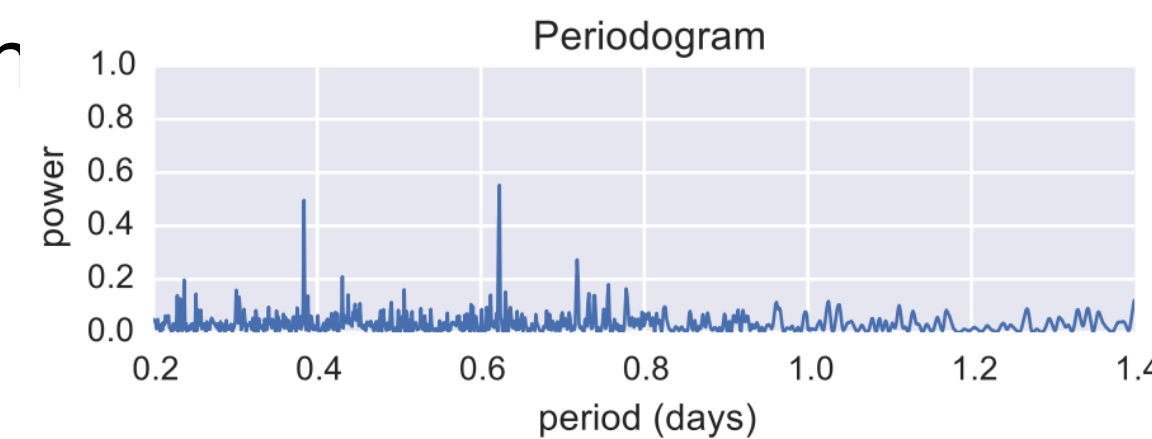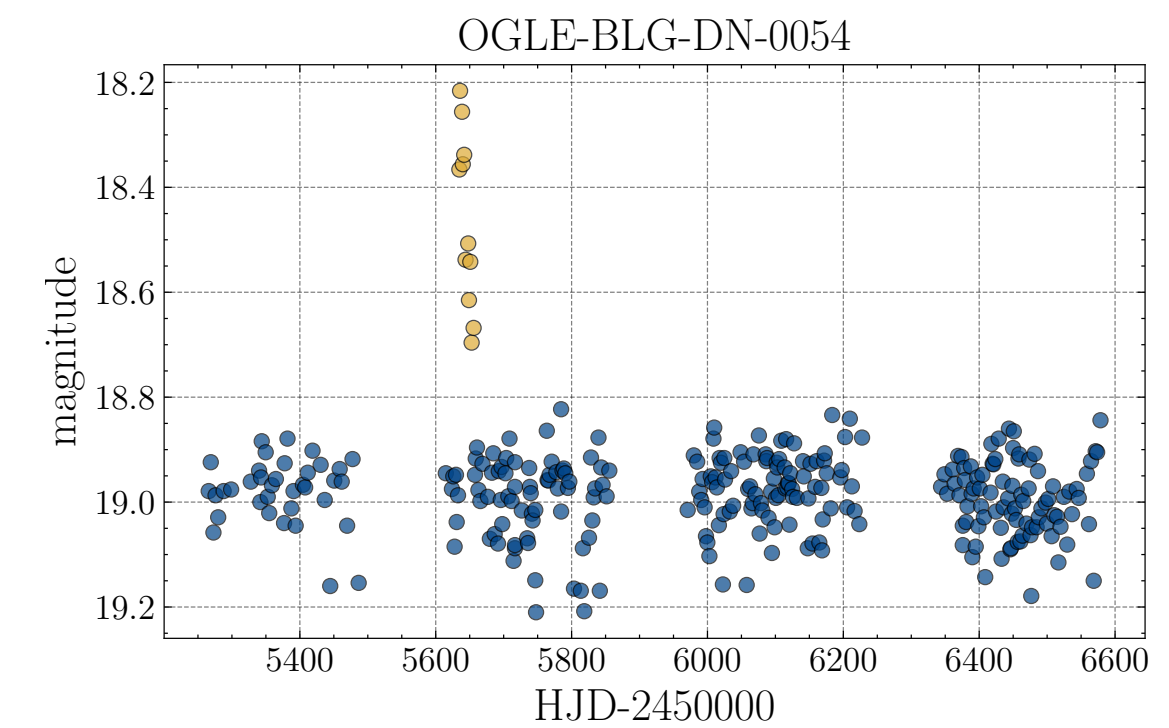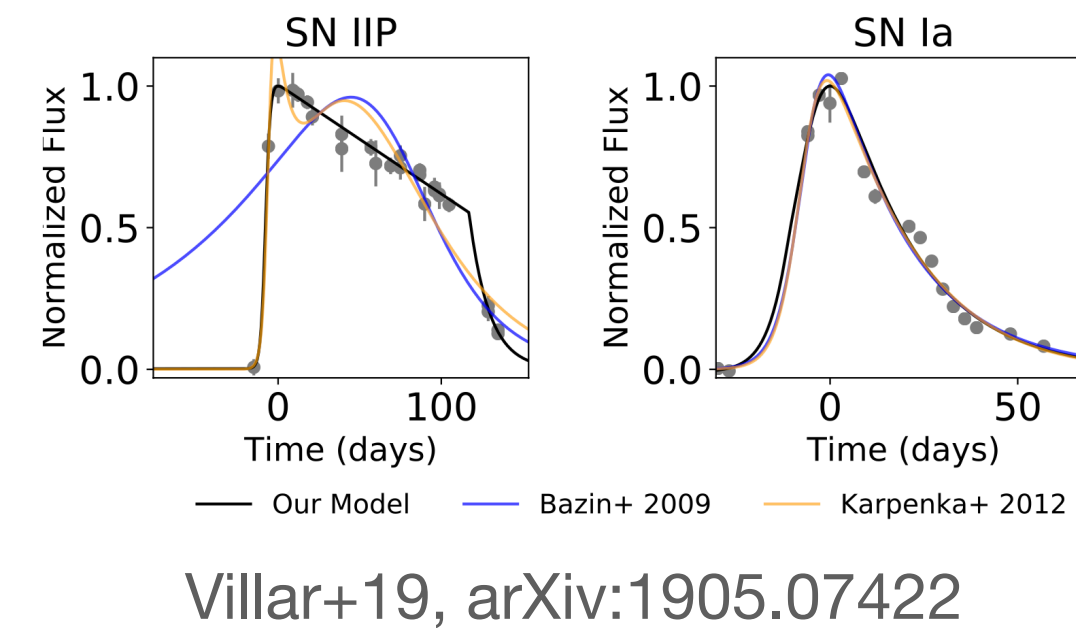**Code: https://github.com/snad-space/coniferest**

- Performant re-implementation of scikit-learn's IsolationForest

- Two "active" algorithms atop of it: AAD (Das+2017) and Pineforest

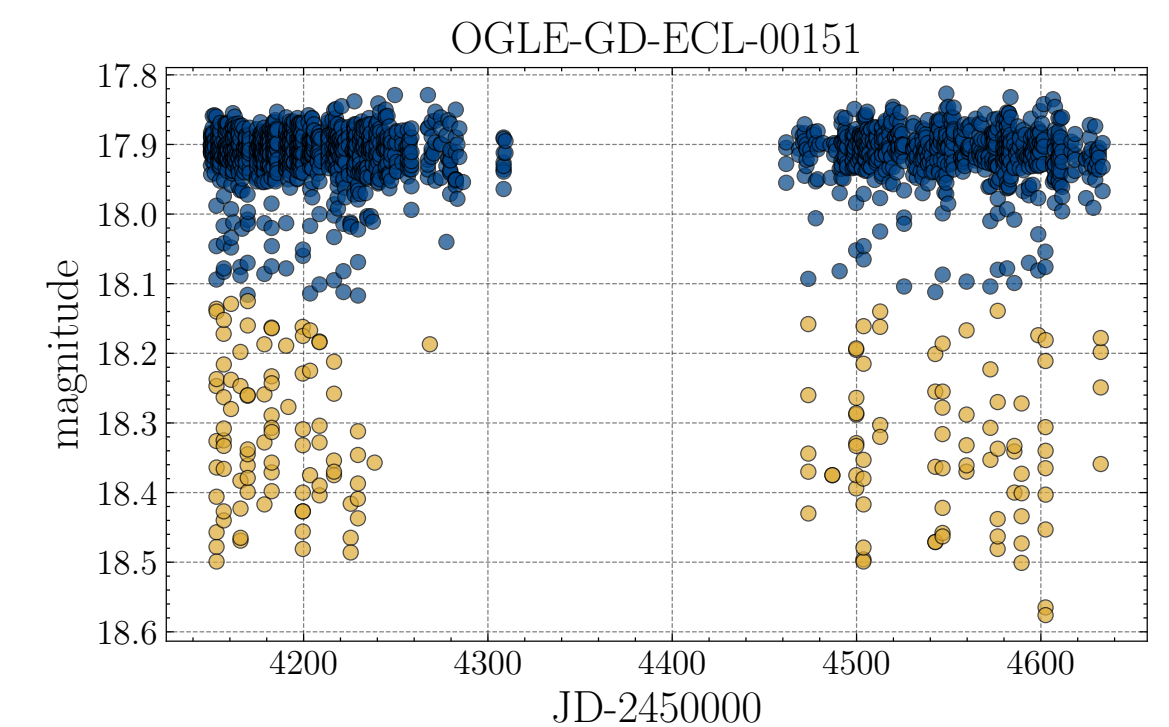- `Session` class which handles interactive pipeline

# Astronomical data: time-series features

`light-curve` package in Python and Rust, https://github.com/light-curve

- Rich feature set

  - Magnitude statistics: mean-, median-, momentum- quartile-ba[...]

  - Shape-based: Stetson (1996) K, $\eta^e$ (Kim+ 2014)

  - "Fast" Lomb–Scargle periodogram peaks and other derivatives

  - Parametric fits: linear, SN-like functions: Bazin+ 2009,
    Villar+ 2019, **new Rainbow approach Russeil+2024**

  - New Otsu-split extractor: powerful features to classify recurren[...]
    outbursts, eclipsing binaries, etc (Lavrukhina+2023)

- Hundreds of unit tests, pre-built wheels for Linux and macOS

- Serves **three ZTF/LSST brokers**: AMPEL, ANTARES, Fink
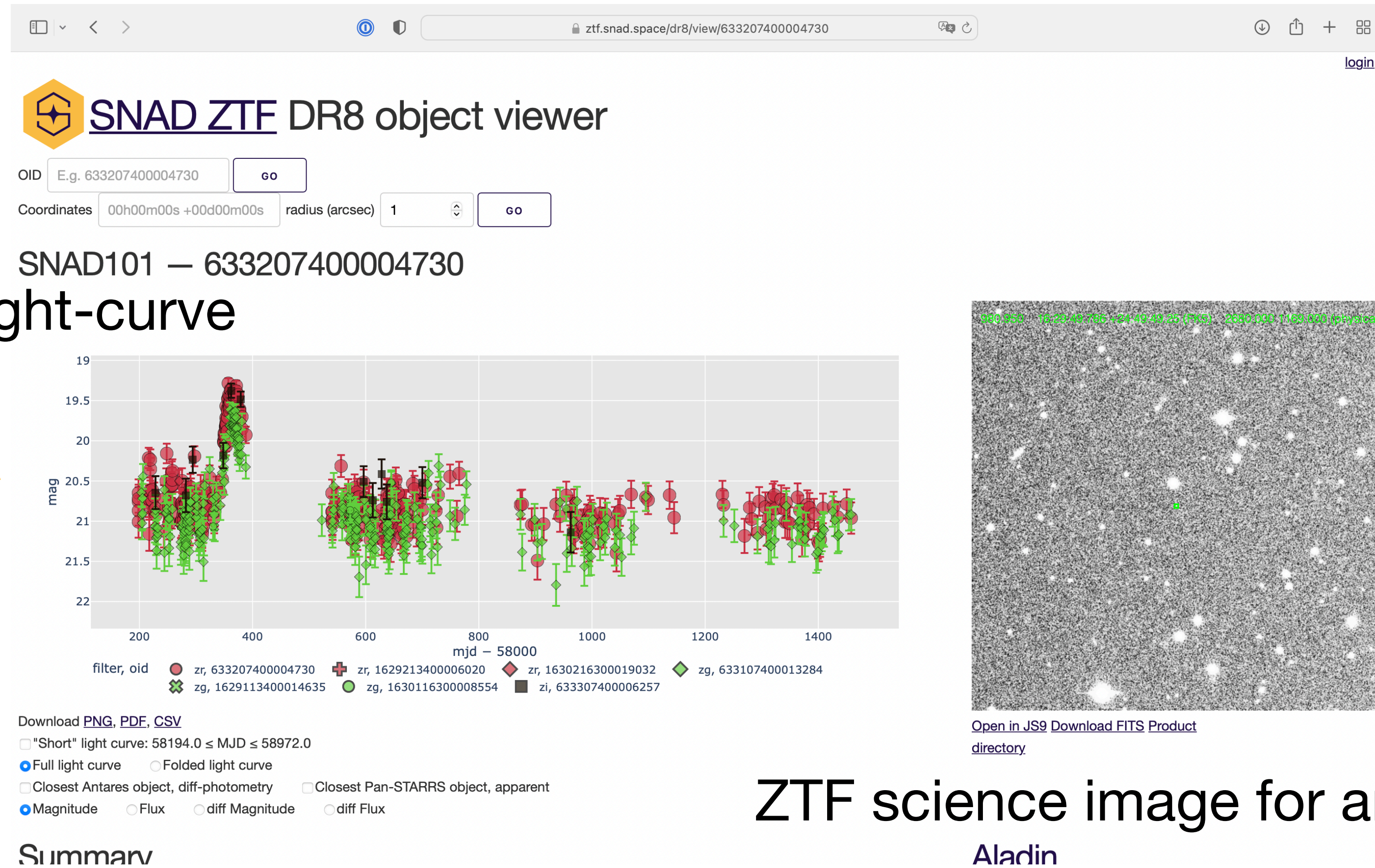
- `python3 -m pip install light-curve`



Villar+19, arXiv:1905.07422



VanderPlas & Ivezić 15, arXiv:1502.01344

# Astronomical data: Expert Portal
## SNAD VIewer, htts://ztf.snad.space



Self-matched ZTF light-curve

ZTF science image for any detection

# Astronomical data: Expert Portal
## SNAD VIewer, htts://ztf.snad.space



Name, type, period, distance & extension from other catalogs and our periodogram

# Tutorial Notebooks



Basic tutorial

- "Static" anomaly detection
- Toy data
- Light-curve features



US names time series



MNIST digit images