



# Active Anomaly Detection for time-domain discoveries

arXiv:astro-ph/1909.13260

*4th Inter-experiment Machine Learning Workshop - CERN  
23 October 2020*

Emille E. O. Ishida

*Laboratoire de Physique de Clermont - Université Clermont-Auvergne  
Clermont Ferrand, France*



# Summary

- Context
- Anomaly detection
- Human in the loop
- Applications to astronomy
  - Simulations
  - Real data
- Conclusions
- Implications

*Context ...*

Astronomy has been,  
traditionally,  
an experience of  
solitude



The old astronomer, by Charlie Bowater <sup>3</sup>

*Context ...*

Astronomy has been,  
traditionally,  
an experience of  
solitude

*Meaning ...*

Discovery happened by chance

or

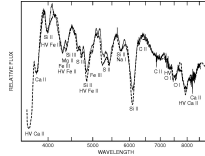
By careful analysis of small volume of highly  
informative data



The old astronomer, by Charlie Bowater <sup>4</sup>

## 2 types of Astronomical data

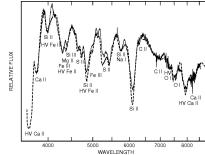
# Spectroscopic



- Pros:
  - High resolution
  - Large information content
  - Enable astrophysical analysis
- Cons:
  - Super expensive
  - Strong requirements on observation conditions

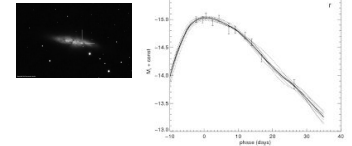
# 2 types of Astronomical data

## Spectroscopic



- Pros:
  - High resolution
  - Large information content
  - Enable astrophysical analysis
- Cons:
  - Super expensive
  - Strong requirements on observation conditions

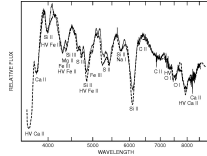
## Photometry



- Pros:
  - Easy to obtain
  - Allows environmental and morphological analysis
  - Enables time domain research of a large number of objects
- Cons:
  - Low resolution (integration over large wavelength range)

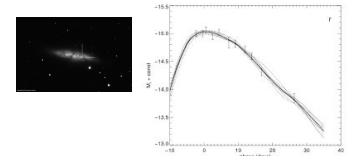
## 2 types of Astronomical data

# Spectroscopic

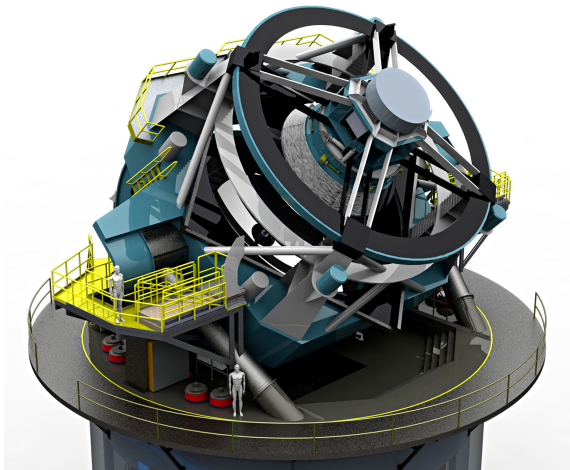


- Pros:
  - High resolution
  - Large information content
  - Enable astrophysical analysis
- Cons:
  - Super expensive
  - Strong requirements on observation conditions

# Photometry



- Pros:
  - Easy to obtain
  - Allows environmental and morphological analysis
  - Enables time domain research of a large number of objects
- Cons:
  - Low resolution (integration over large wavelength range)



# The Vera Rubin Observatory Legacy Survey of Space and Time (LSST)

**~10 million candidates/night**

Over a total life span of 10 years

## Serendipitous discoveries will not happen...





Home

About

Project

Results

Contact

Workshop 2018

Workshop 2019

Workshop 2020

Workshop 2021

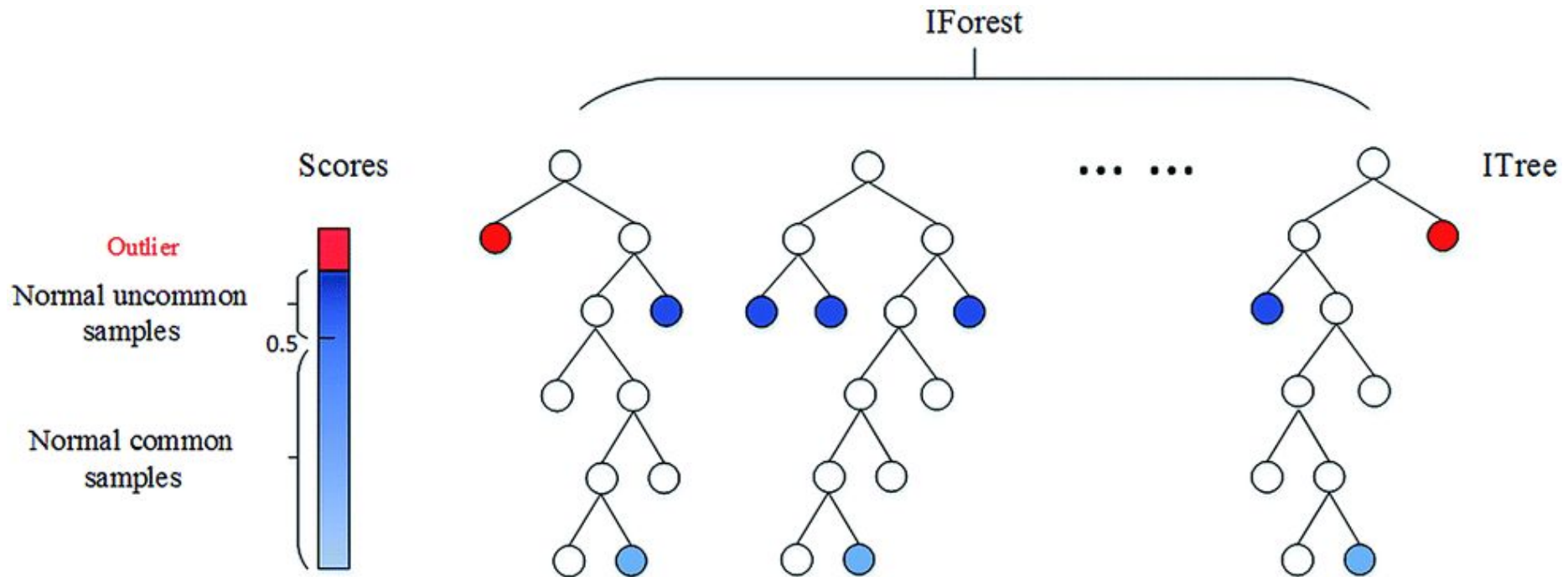
A wide-angle photograph of the night sky showing the Milky Way galaxy stretching across the frame. The galaxy's core is visible as a bright, yellowish-white region. Below the sky, a dark, silhouetted desert landscape is visible, with a small, bright orange light source on the horizon, likely a setting or rising sun or moon.

# Supernova Anomaly Detection

<https://snad.space/>

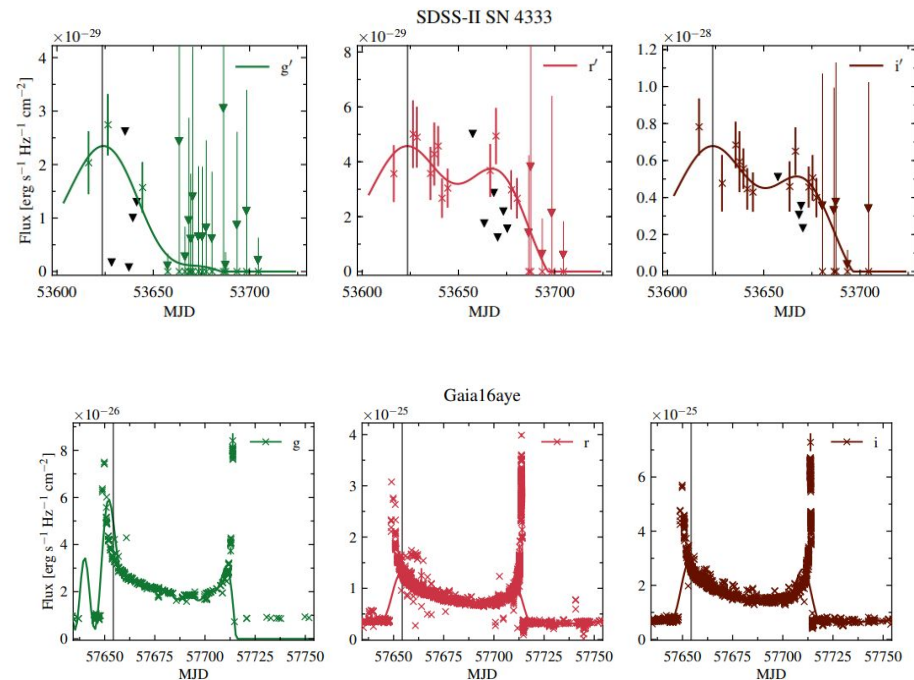
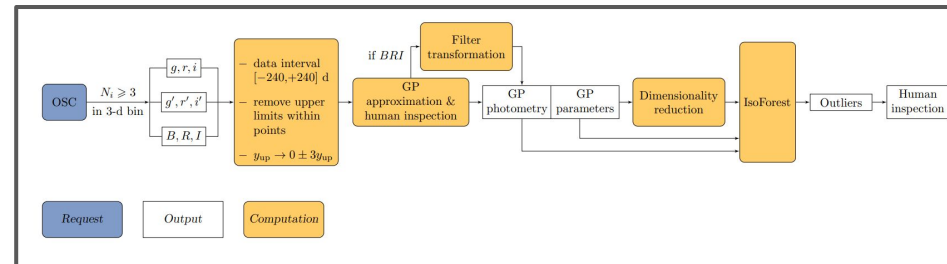


# Isolation Forest



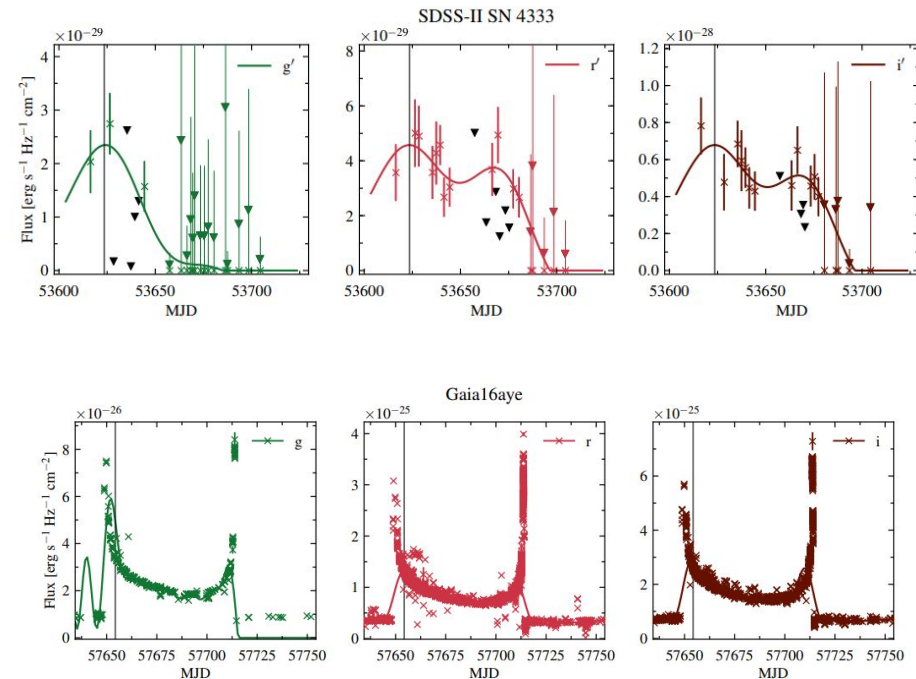
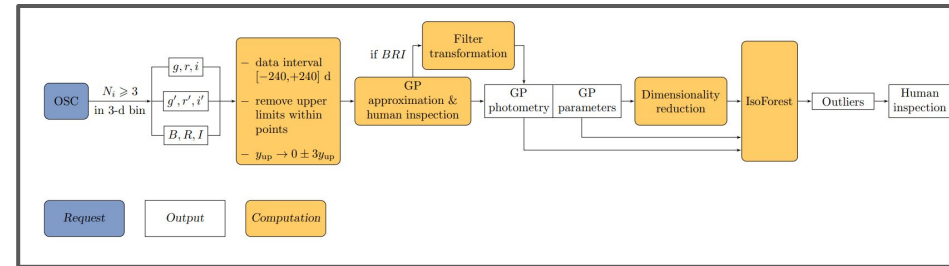
# First try: the Open Supernova Catalog

- Real data:
  - Uncertainties
  - Upper limits
  - Different filter sets
- Pre-processing:
  - Filter translation
  - Selection cuts
  - 2D Gaussian Process
  - 3 sets: photo, photo + GP param, tSNE
- Data and analysis:
  - Initial data: 2000 light curves
  - Anomaly detection via Isolation Forest
  - Visually inspected 2% in each set ( $\sim 100$  objs)
- Results:
  - 81 identified anomalies
  - SLSN, peculiar SNe, miss-classified stars
  - 1 AGN and 1 binary micro-lensing



# First try: the Open Supernova Catalog

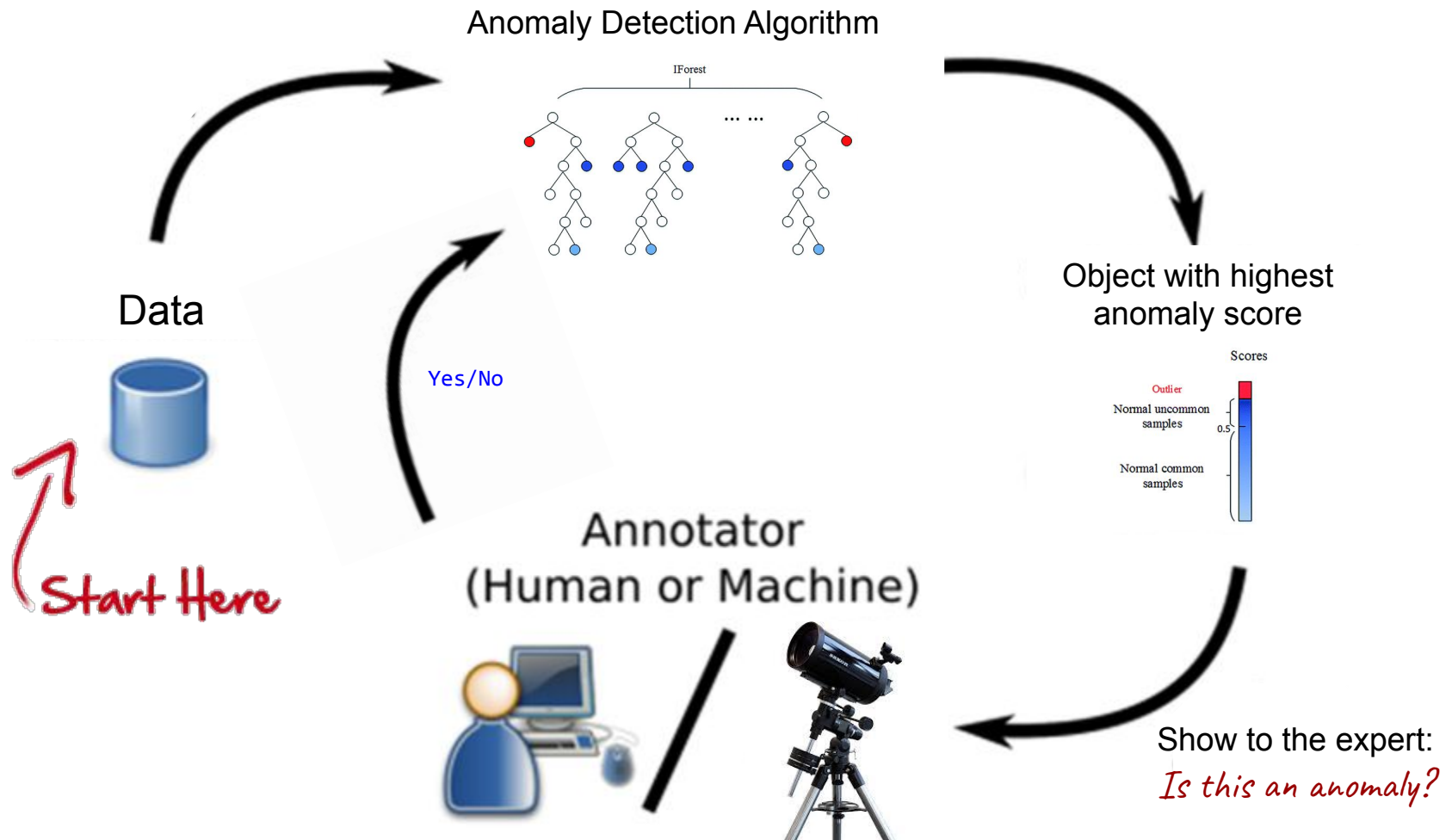
- Real data:
  - Uncertainties
  - Upper limits
  - Different filter sets
- Pre-processing:
  - Filter translation
  - Selection cuts
  - 2D Gaussian Process
  - 3 sets: photo, photo + GP param, tSNE
- Data and analysis:
  - Initial data: 2000 light curves
  - Anomaly detection via Isolation Forest
  - Visually inspected 2% in each set ( $\sim 100$  objs)
- Results:
  - 81 identified anomalies
  - SLSN, peculiar SNe, miss-classified stars
  - 1 AGN and 1 binary micro-lensing



Problem:

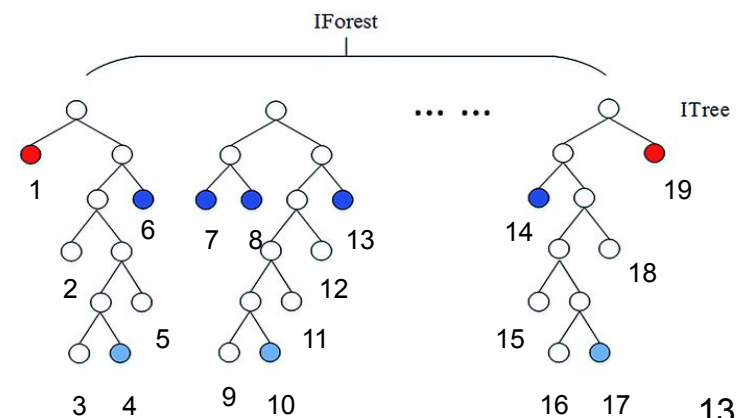
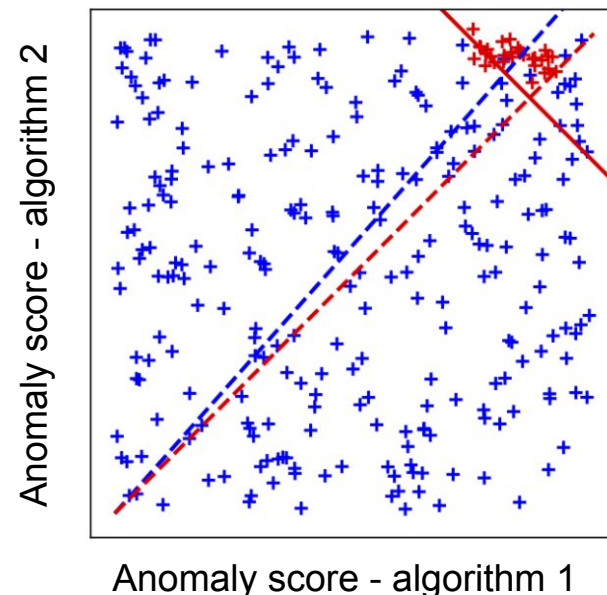
*What is scientifically interesting is in the eye of the beholder.*

# Active Learning for Anomaly Detection



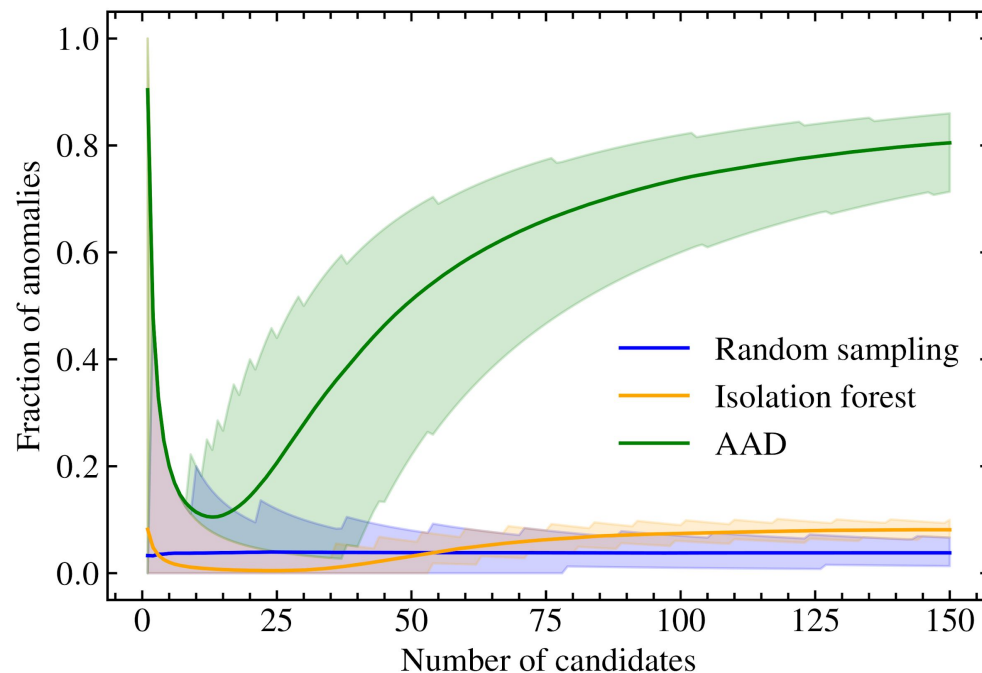
# Ensemble learning and expert feedback

- Ensemble learning: *when you do not know, ask around!*
- 2 perfectly accurate AD algorithms will agree on the scores for true anomalies
- In the real case one will be more accurate than the other, so we need to assign **weights**
- **Active Anomaly Discovery:**
  - Start with a normal Isolation Forest
  - Consider each decision path leading to a leaf node as an weak AD algorithm (ensemble member)
  - Assign an equal weight to each ensemble member
  - Show the most anomalous obj to the expert
  - If `expert_answer == yes`:
    - Show next obj with highest anomaly scoreelse:
    - Update weights



# Simulations: the PLAsTiCC data set

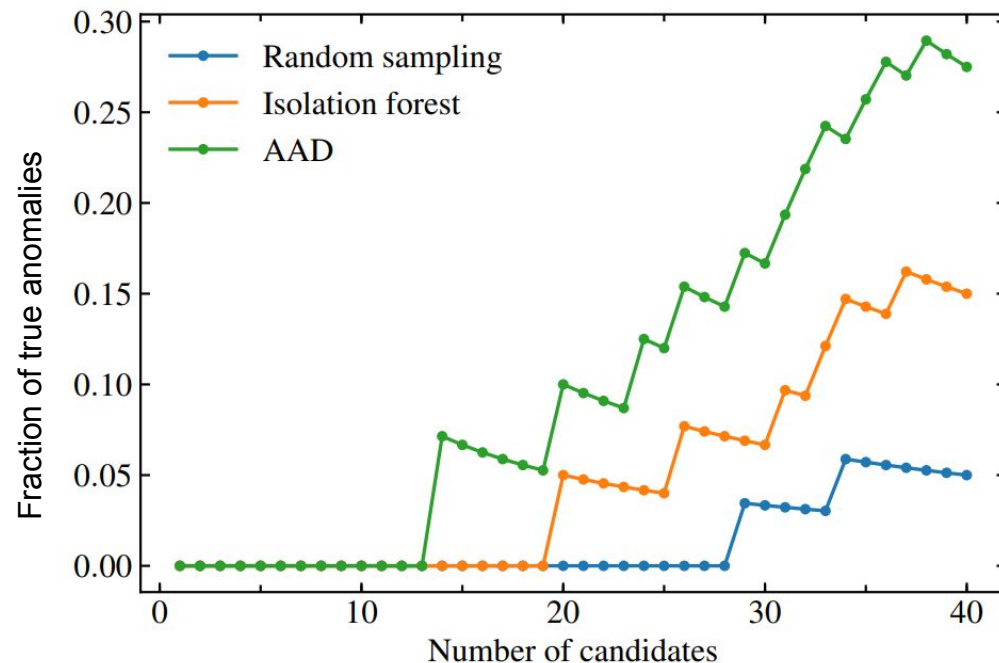
- Data from the Kaggle PLAsTiCC data set restricted to Supernova-like events
- Initial sample:  $\sim 7000$  light curves, 3 known classes, 3 peculiar classes (277 anomalies, 4%)
- 145 objects scrutinized ( $\sim 2\%$ ), on average:
  - **Random Sampling:** 5 real anomalies
  - **Isolation Forest:** 12 real anomalies
  - **Active Anomaly Discovery:** 120 real anomalies





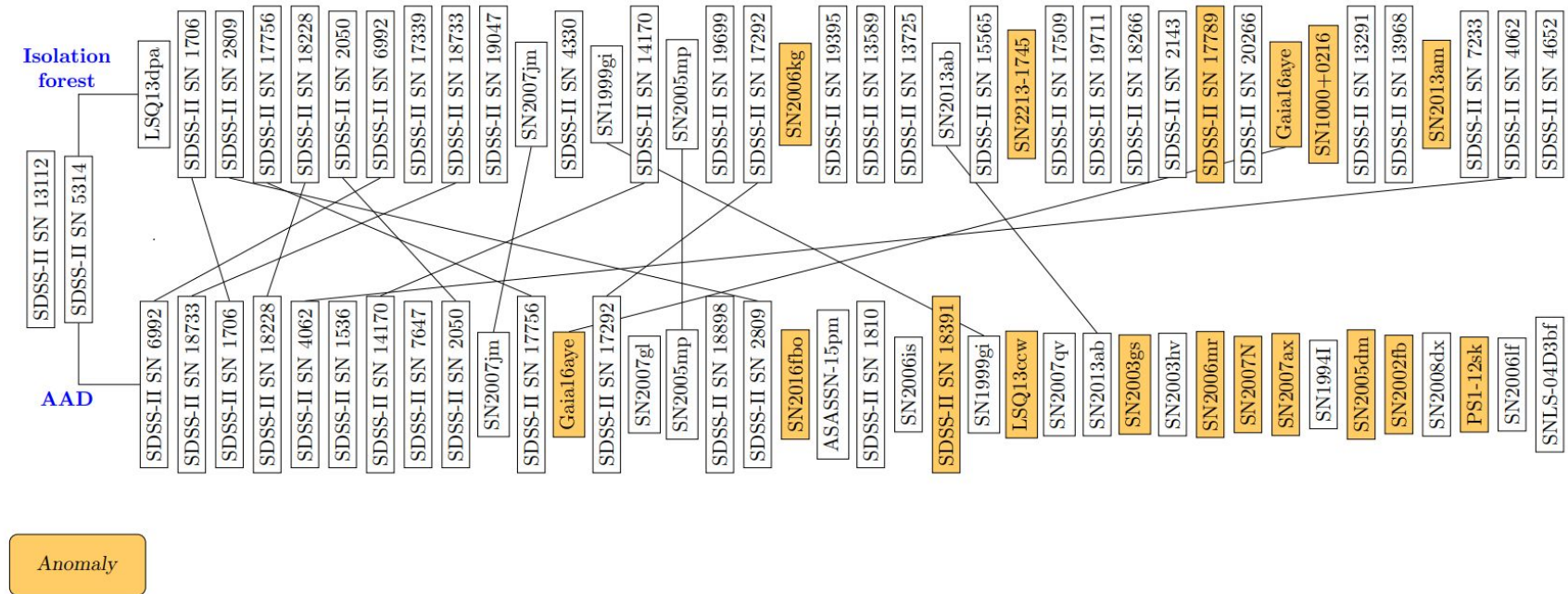
# Real data: The Open Supernova Catalog

- **Anomaly:**
  - Miss-classification (non-SNe)
  - Unusual light curve behavior
  - Previously known 91bg-like and 91T-like
- **Not-anomaly:**
  - Bad Gaussian process fitting
  - Not enough signal
  - Identified artifacts
- **Results within 2% contamination (40 objs):**
  - Random sampling: 2 (5%)
  - Isolation Forest: 5 (15%)
  - Active Anomaly Discovery: 11 (27%)



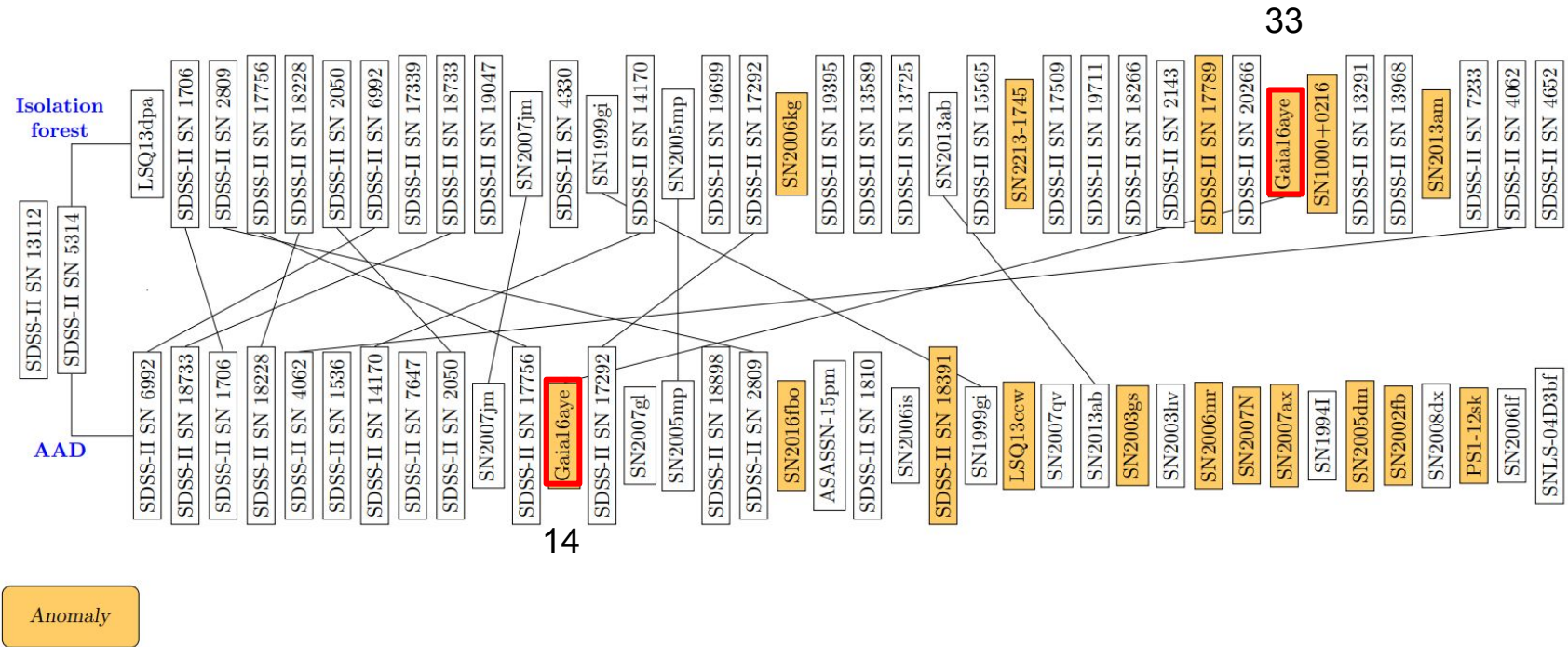
*AAD was able to increase the incidence of true anomalies presented to the expert in 80%*

# Real data: The Open Supernova Catalog

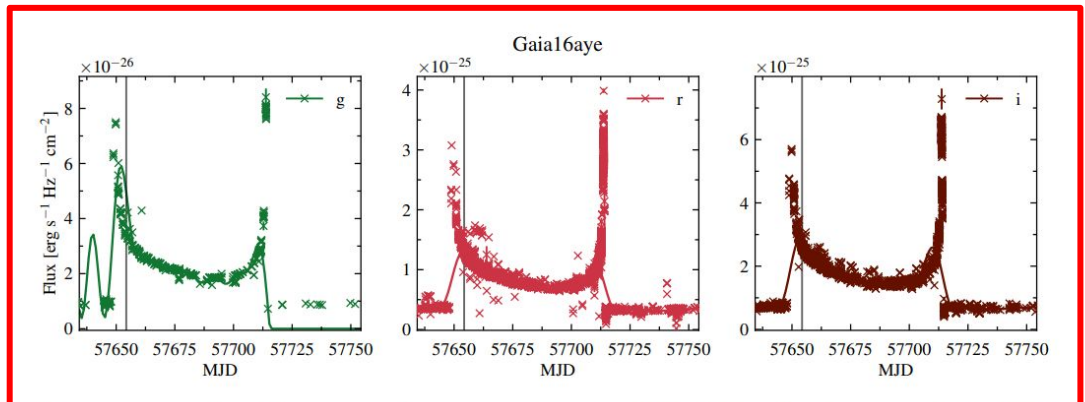


- It requires some time for changes to be effectively incorporated
- Late queries:
  - Objects which were not found in the static case
  - Higher concentration of true anomalies

# Active Anomaly Detection for time-domain discoveries



*Fast identification of binary microlensing event*



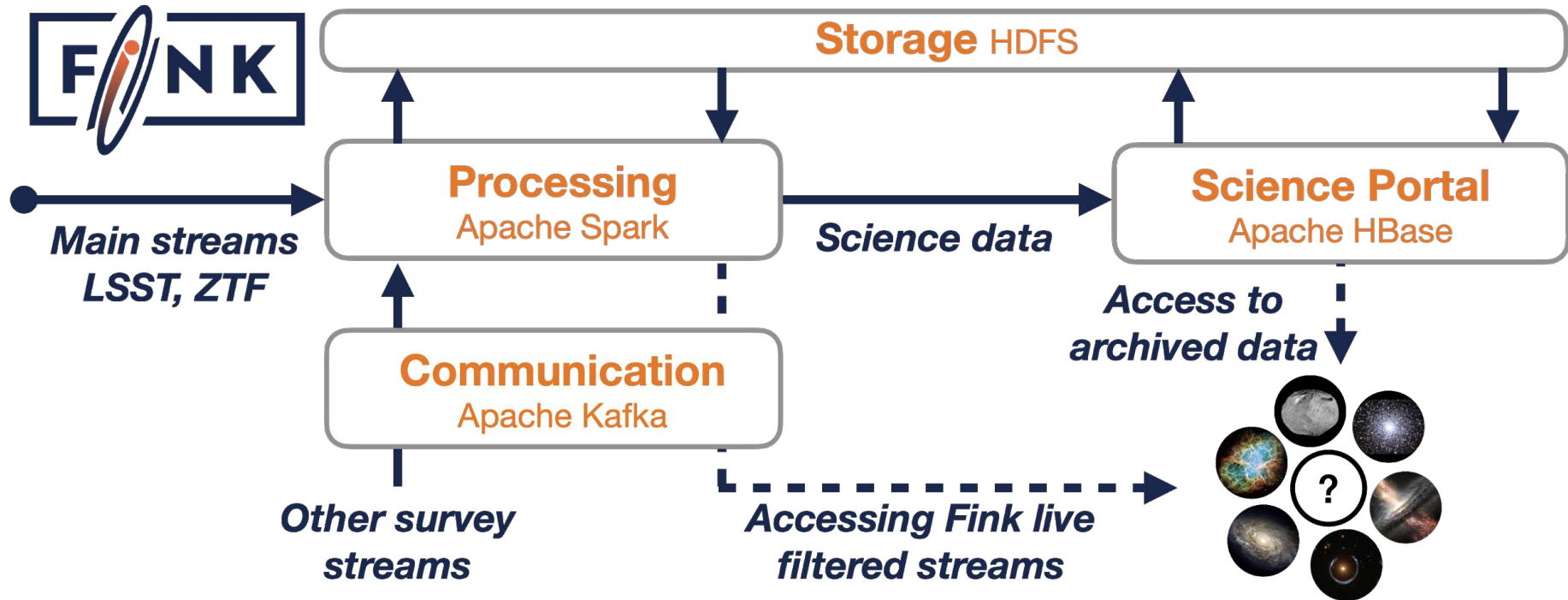
# Conclusions

- Active Anomaly Detection can be a powerful tool to boost discoveries
- Approach is still under development in other fields
  - Opportunity to develop astronomy-oriented strategies
- Astronomical data has many caveats which are not necessarily taken into account by off-the-shelf algorithms
- Collaboration is essential

<https://snad.space/>

# Implications

*A French-born broker to digest LSST alerts and search for interesting astrophysical objects*



Community-driven project with important elements on  
Adaptive Learning and Bayesian Deep Learning



# Thank you, Merci, Спасибо

*From the SNAD team!*



**SNAD**

<https://snad.space/>



Back-up slides

# Active Anomaly Discovery

$$l(q, \mathbf{w}; z_i, y_i) =$$

$$\begin{cases} 0 & \text{if } \mathbf{w} \cdot \mathbf{z}_i \geq q \text{ and } y_i = \text{anomaly} \\ 0 & \text{if } \mathbf{w} \cdot \mathbf{z}_i < q \text{ and } y_i = \text{normal} \\ q - \mathbf{w} \cdot \mathbf{z}_i & \text{if } \mathbf{w} \cdot \mathbf{z}_i < q \text{ and } y_i = \text{anomaly} \\ \mathbf{w} \cdot \mathbf{z}_i - q & \text{if } \mathbf{w} \cdot \mathbf{z}_i \geq q \text{ and } y_i = \text{normal} \end{cases},$$

---

## Algorithm 2 Active Anomaly Discovery (AAD)

---

**Input:** Dataset  $\mathbf{H}$ , budget  $B$

Initialize the weights  $\mathbf{w}^{(0)} = \{\frac{1}{\sqrt{m}}, \dots, \frac{1}{\sqrt{m}}\}$

Set  $t = 0$

Set  $\mathbf{H}_A = \mathbf{H}_N = \emptyset$

**while**  $t \leq B$  **do**

$t = t + 1$

    Set  $\mathbf{a} = \mathbf{H} \cdot \mathbf{w}$  (i.e.,  $\mathbf{a}$  is the vector of anomaly scores)

    Let  $\mathbf{z}_i$  = instance with highest anomaly score (where  $i = \arg \max_i (a_i)$ )

    Get feedback  $\{\text{'anomaly'}/\text{'nominal'}\}$  on  $\mathbf{z}_i$

**if**  $\mathbf{z}_i$  is *anomaly* **then**

$\mathbf{H}_A = \{\mathbf{z}_i\} \cup \mathbf{H}_A$

**else**

$\mathbf{H}_N = \{\mathbf{z}_i\} \cup \mathbf{H}_N$

**end if**

15:  $\mathbf{w}^{(t)}$  = compute new weights; normalize  $\|\mathbf{w}^{(t)}\| = 1$

**end while**

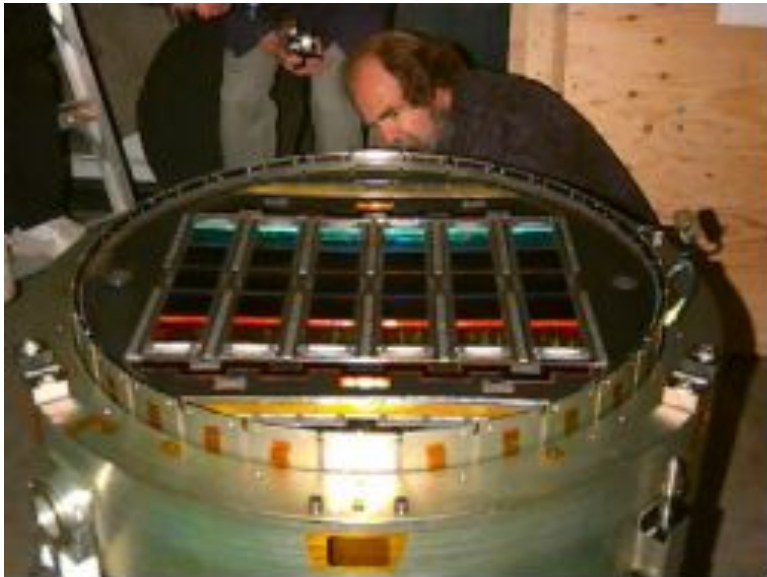
---

$$\begin{aligned} \mathbf{w}^{(t)} = \arg \min_{\mathbf{w}, \xi} \frac{C_A}{|\mathbf{H}_A|} & \left( \sum_{\mathbf{z}_i \in \mathbf{H}_A} \ell(\hat{q}_\tau(\mathbf{w}^{(t-1)}), \mathbf{w}; (\mathbf{z}_i, y_i)) \right) \\ & + \frac{1}{|\mathbf{H}_N|} \left( \sum_{\mathbf{z}_i \in \mathbf{H}_N} \ell(\hat{q}_\tau(\mathbf{w}^{(t-1)}), \mathbf{w}; (\mathbf{z}_i, y_i)) \right) \\ & + \frac{C_\xi}{|\mathbf{H}_A|} \left( \sum_{\mathbf{z}_i \in \mathbf{H}_A} \ell(\mathbf{z}_\tau^{(t-1)} \cdot \mathbf{w}, \mathbf{w}; (\mathbf{z}_i, y_i)) \right) \\ & + \frac{C_\xi}{|\mathbf{H}_N|} \left( \sum_{\mathbf{z}_i \in \mathbf{H}_N} \ell(\mathbf{z}_\tau^{(t-1)} \cdot \mathbf{w}, \mathbf{w}; (\mathbf{z}_i, y_i)) \right) \\ & + \|\mathbf{w} - \mathbf{w}_p\|^2 \end{aligned} \quad (2)$$

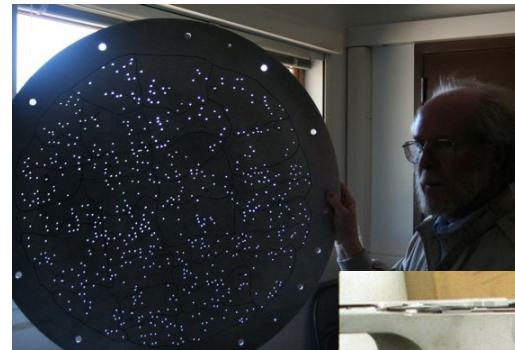
where,  $\mathbf{w}_p = \frac{\mathbf{w}_U}{\|\mathbf{w}_U\|} = [\frac{1}{\sqrt{m}}, \dots, \frac{1}{\sqrt{m}}]^T$ ,  $\mathbf{z}_\tau^{(t-1)}$  and  $\hat{q}_\tau(\mathbf{w}^{(t-1)})$  are

# Photometry x Spectroscopy

*An example from SDSS*



Exposure time 2 x 54s



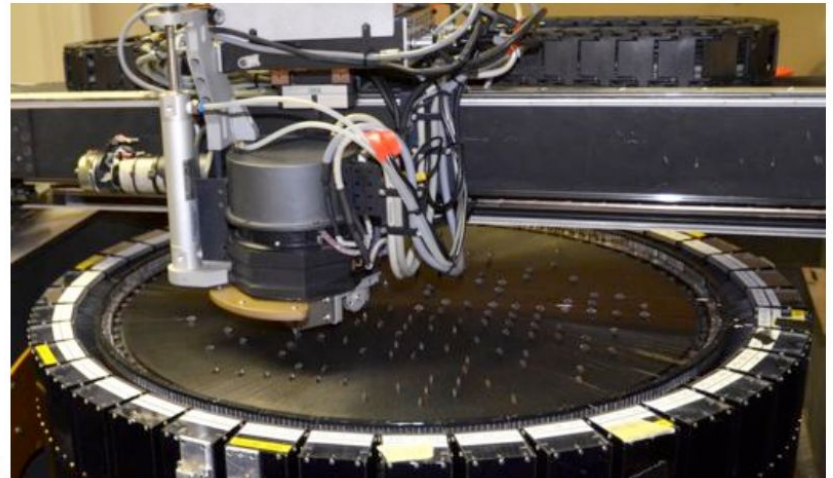
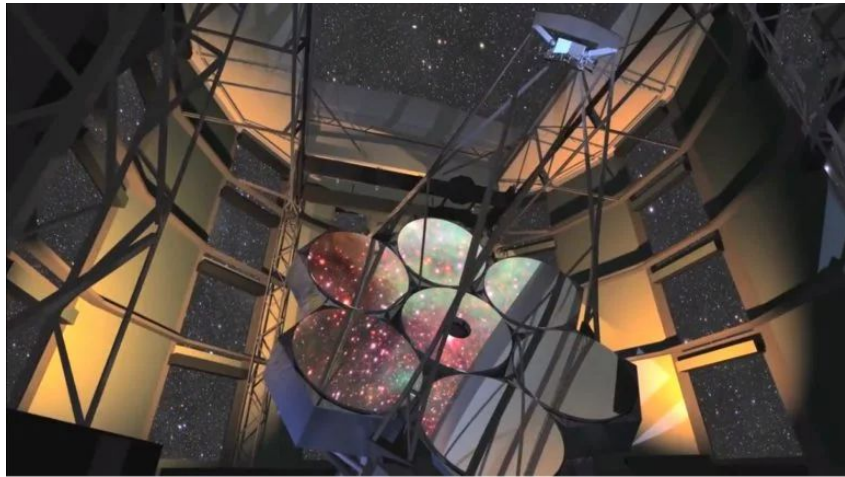
Integration time of at least

45 minutes<sub>23</sub>

# Photometry x Spectroscopy

*An example from the Australian Astronomical Observatory*

For the Giant Magellan Telescope (GMT)  
First light 2025



Integration time much larger...