

FL EMING : are we able to find SC4K-like LAEs using AI?

Ana Paulino-Afonso

Institute of Astrophysics and Space Sciences (IA-U.Porto)

FCT official team members: Andrew Humphrey (DTx CoLab/IA-U.Porto), Jarle Brinchmann (IA-U.Porto), Israel Matute (IA-U.Lisboa), Rodrigo Pizarro (IA-U.Lisboa/FCUL), Thomas Scott (IA-U.Porto), Pedro Cunha (IA-U.Porto/FCUP), José Fonseca (IA-U.Porto/FCUP), **Afonso do Vale (IA-U.Porto/FCUP)**, and **Bruno Cerqueira (IA-U.Porto/FCUP)**

Unofficial team members, but people always available to discuss and help: Bruno Ribeiro (GBL/Celfocus), David Sobral (BNP Paribas), **XGAL team**, and **CRISP team**

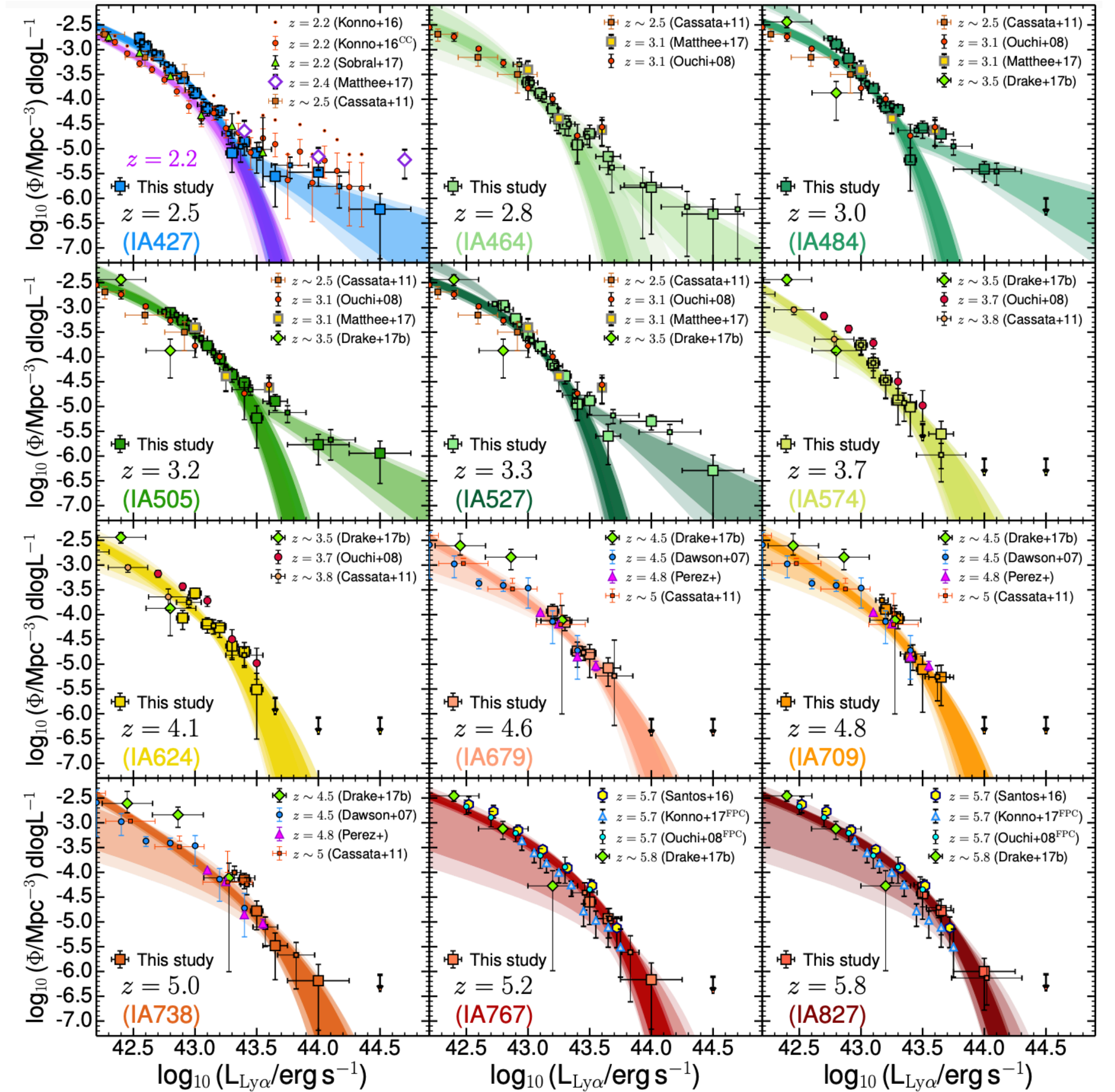
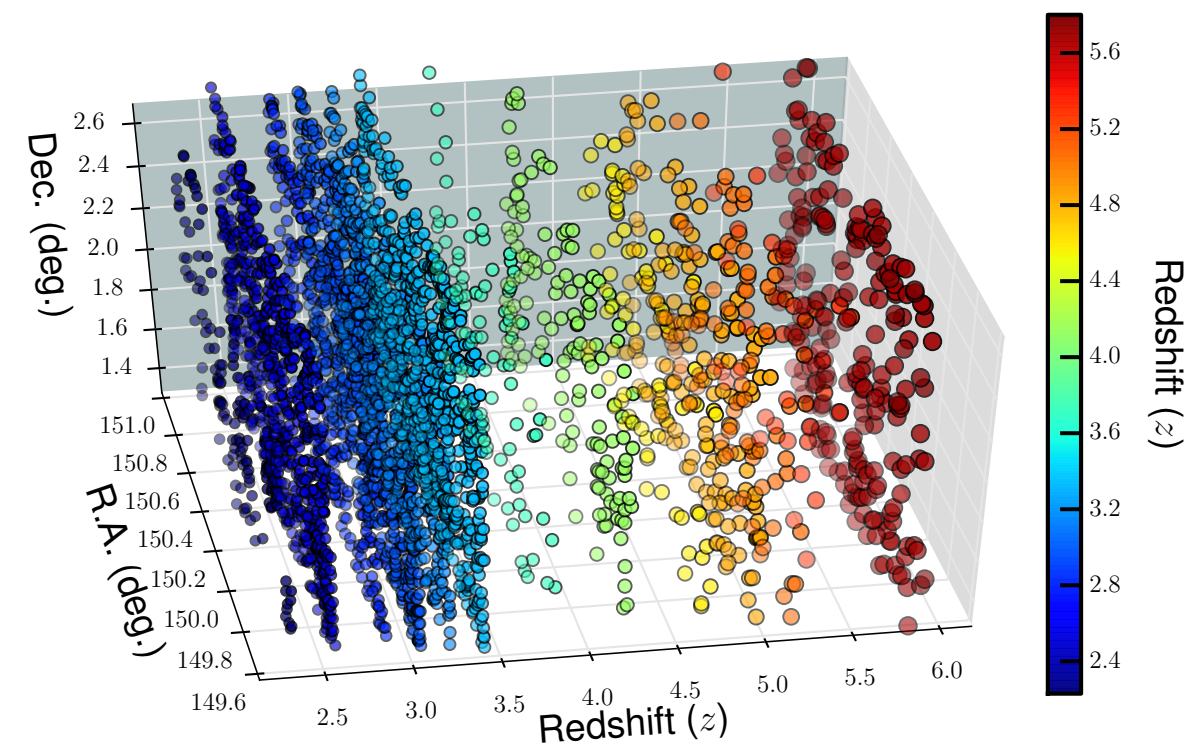
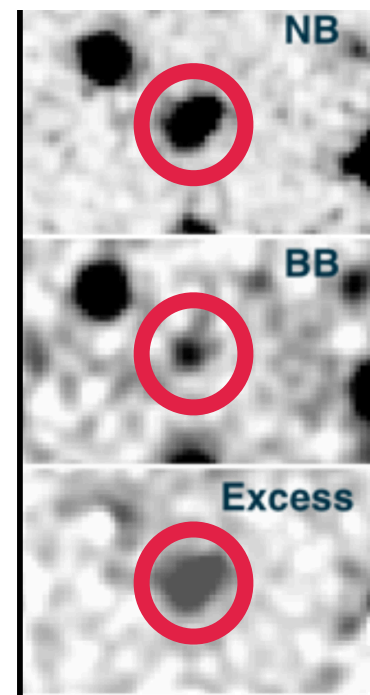
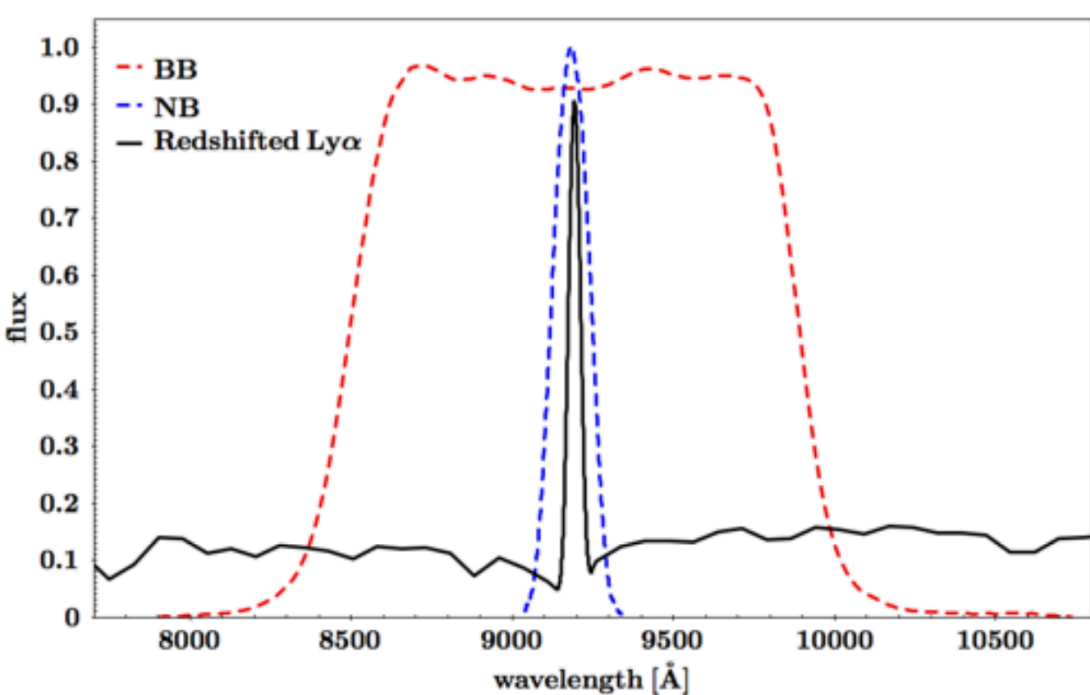
This work was supported by Fundação para a Ciência e a Tecnologia (FCT) through the research grants UIDB/04434/2020 and UIDP/04434/2020, the FCT Investigador FCT Contract No. 2020.03946.CEECIND, and in the form of an exploratory project with the EXPL/FIS-AST/1085/2021 reference (PI: Paulino-Afonso).



Introducing & Motivation: SC4K

Sobral et al. 2018a

- A good, well-understood selection that can be applied with current instrumentation
- Well calibrated + sensitive + resulting in a **representative population of galaxies**
- Able to uniformly select **large samples**
- **Different epochs + large areas + best-studied fields**

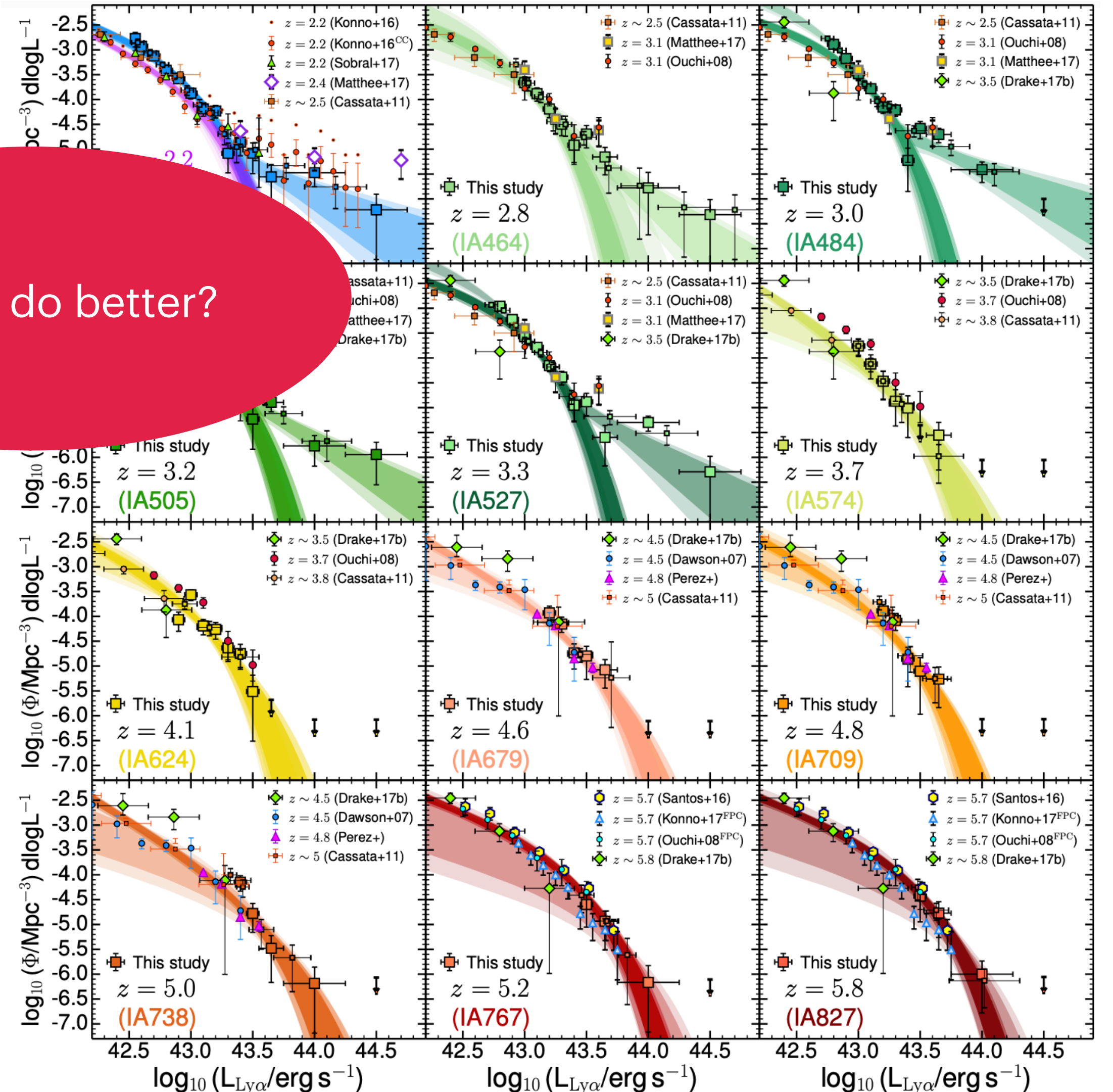
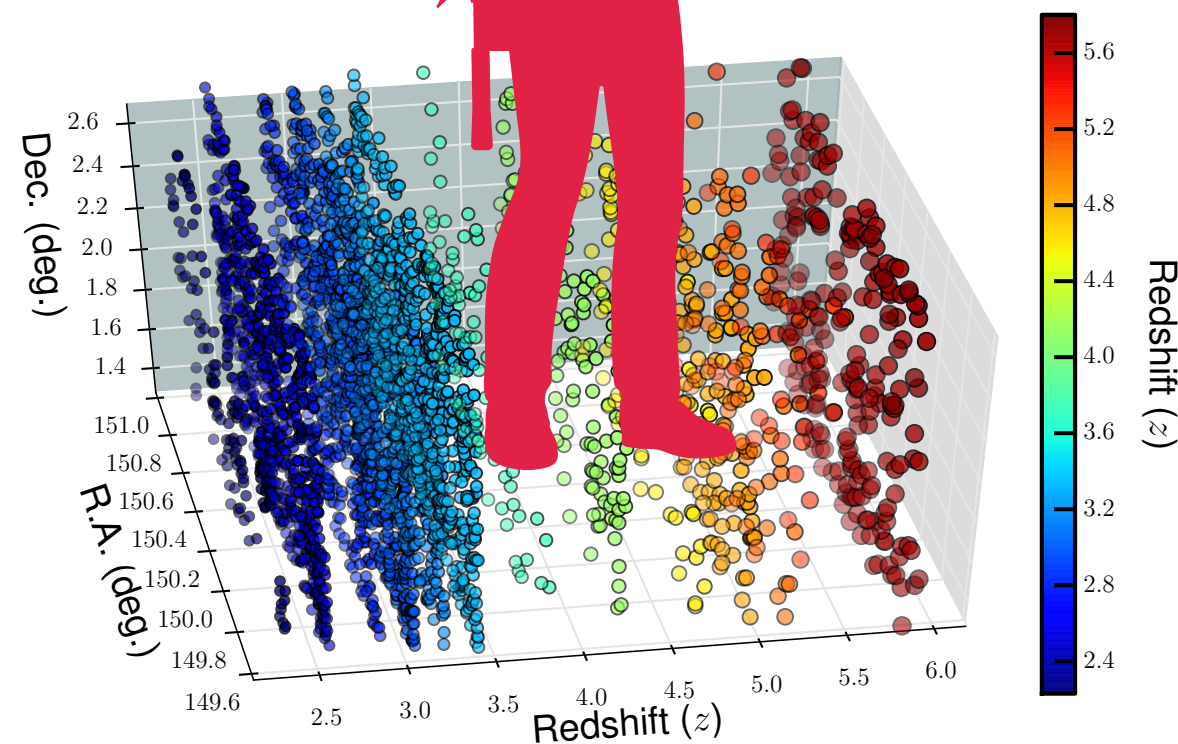
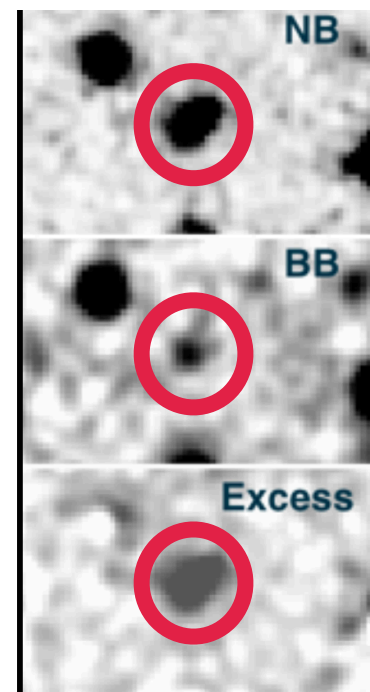
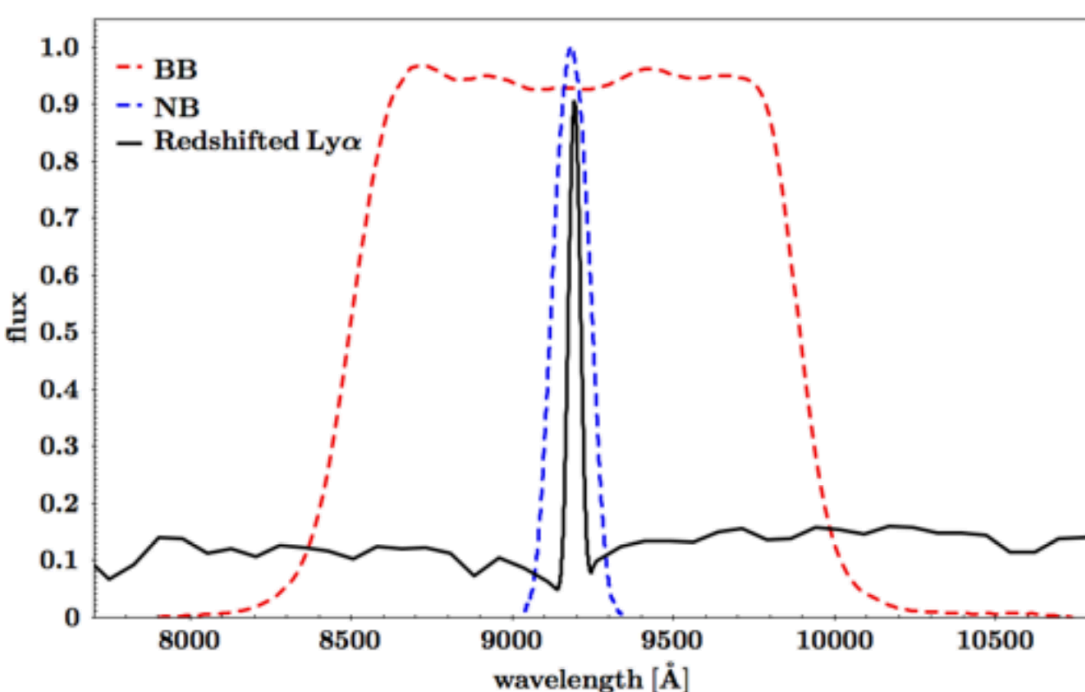


Introducing & Motivation: SC4K

Sobral et al. 2018a

- A good, well-understood selection that can be applied with current instrumentation
- Well calibrated + sensitive + resulting in a representative population of galaxies
- Able to uniformly select large samples
- Different epochs + large areas + best-studied fields

Can we do better?



Introducing & Motivation: SILVERRUSH

Ono et al. 2022

ACCEPTED FOR PUBLICATION IN APJ
Preprint typeset using L^AT_EX style emulateapj v. 12/16/11

9,318 new LAE candidates @ $2.2 < z < 7.0$

177 of them with spectroscopic confirmation



SILVERRUSH X: MACHINE LEARNING-AIDED SELECTION OF **9,318 LAES**
AT $z = 2.2, 3.3, 4.9, 5.7, 6.6, \text{ AND } 7.0$ FROM THE HSC SSP AND CHORUS SURVEY DATA

YOSHIAKI ONO¹, RYOHEI ITOH^{1,2}, TAKATOSHI SHIBUYA³, MASAMI OUCHI^{4,1,5}, YUICHI HARIKANE^{1,6}, SATOSHI YAMANAKA^{7,8},
AKIO K. INOUE^{9,7}, TOSHIYUKI AMAGASA^{10,11}, DAICHI MIURA¹⁰, MAIKI OKURA¹⁰, KAZUHIRO SHIMASAKU^{12,13},
IKURU IWATA⁴, YOSHIAKI TANIGUCHI¹⁴, SEIJI FUJIMOTO¹⁵, MASANORI IYE⁴, ANTON T. JAELANI^{16,17},
NOBUNARI KASHIKAWA^{12,13}, SHOTARO KIKUCHIHARA^{1,12}, SATOSHI KIKUTA¹¹, MASAKAZU A.R. KOBAYASHI¹⁸,
HARUKA KUSAKABE¹⁹, CHIEN-HSIU LEE²⁰, YONGMING LIANG⁴, YOSHIKI MATSUOKA⁸, RIEKO MOMOSE¹²,
TOHRU NAGAO⁸, KIMIHIKO NAKAJIMA⁴, AND KEN-ICHI TADAKI⁴

Accepted for publication in ApJ

ABSTRACT

We present a new catalog of 9,318 Ly α emitter (LAE) candidates at $z = 2.2, 3.3, 4.9, 5.7, 6.6,$ and 7.0 that are photometrically selected by the SILVERRUSH program with a machine learning technique from large area (up to 25.0 deg^2) imaging data with six narrowband filters taken by the Subaru Strategic Program with Hyper Suprime-Cam (HSC SSP) and a Subaru intensive program, Cosmic Hydrogen Reionization Unveiled with Subaru (CHORUS). We construct a convolutional neural network that distinguishes between real LAEs and contaminants with a completeness of 94% and a contamination rate of 1%, enabling us to efficiently remove contaminants from the photometrically selected LAE candidates. We confirm that our LAE catalogs include 177 LAEs that have been spectroscopically identified in our SILVERRUSH programs and previous studies, ensuring the validity of our machine learning selection. In addition, we find that the object-matching rates between our LAE catalogs and our previous results are $\simeq 80\text{--}100\%$ at bright NB magnitudes of $\lesssim 24 \text{ mag}$. We also confirm that the surface number densities of our LAE candidates are consistent with previous results. Our LAE catalogs will be made public on our project webpage.

Keywords: galaxies: formation — galaxies: evolution — galaxies: high-redshift

Introducing & Motivation: SILVERRUSH

Ono et al. 2022

ACCEPTED FOR PUBLICATION IN APJ
Preprint typeset using L^AT_EX style emulateapj v. 12/16/11

9,318 new LAE candidates @ $2.2 < z < 7.0$

177 of them with spectroscopic confirmation



SILVERRUSH X: MACHINE LEARNING-AIDED SELECTION OF **9,318 LAES**
AT $z = 2.2, 3.3, 4.9, 5.7, 6.6,$ AND 7.0 FROM THE HSC SSP AND CHORUS SURVEY DATA

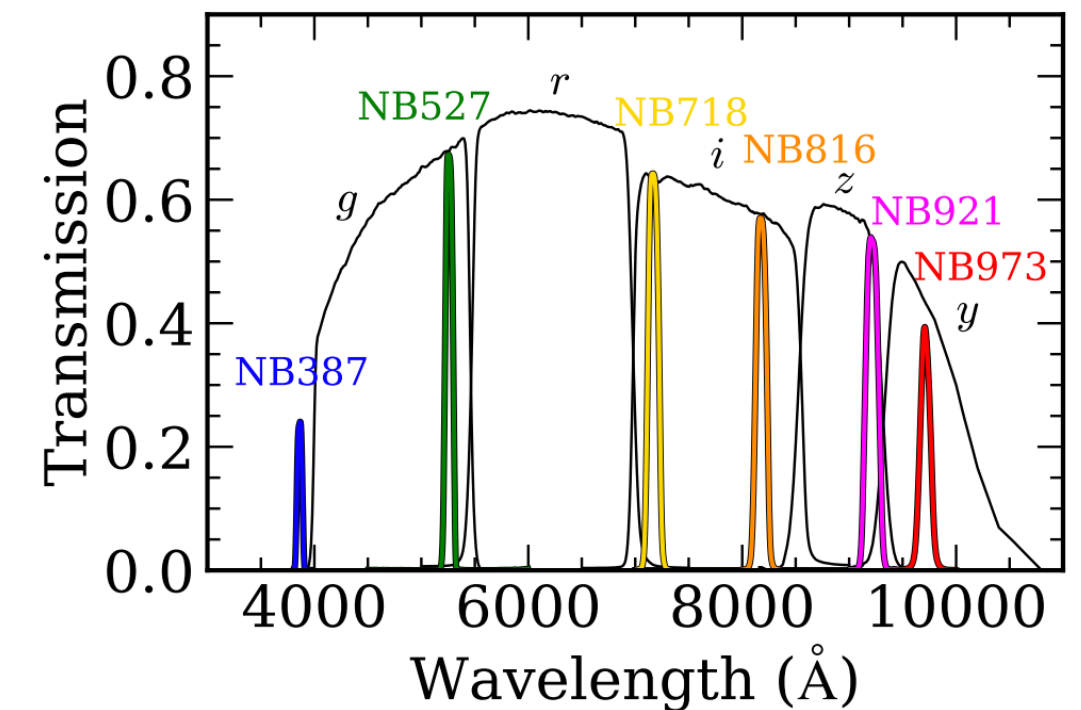
YOSHIAKI ONO¹, RYOHEI ITOH^{1,2}, TAKATOSHI SHIBUYA³, MASAMI OUCHI^{4,1,5}, YUICHI HARIKANE^{1,6}, SATOSHI YAMANAKA^{7,8}, AKIO K. INOUE^{9,7}, TOSHIYUKI AMAGASA^{10,11}, DAICHI MIURA¹⁰, MAIKI OKURA¹⁰, KAZUHIRO SHIMASAKU^{12,13}, IKURU IWATA⁴, YOSHIAKI TANIGUCHI¹⁴, SEIJI FUJIMOTO¹⁵, MASANORI IYE⁴, ANTON T. JANELANI^{16,17}, NOBUNARI KASHIKAWA^{12,13}, SHOTARO KIKUCHIHARA^{1,12}, SATOSHI KIKUTA¹¹, MASAKAZU A.R. KOBAYASHI¹⁸, HARUKA KUSAKABE¹⁹, CHIEN-HSIU LEE²⁰, YONGMING LIANG⁴, YOSHIKI MATSUOKA⁸, RIEKO MOMOSE¹², TOHRU NAGAO⁸, KIMIHIKO NAKAJIMA⁴, AND KEN-ICHI TADAKI⁴

Accepted for publication in ApJ

ABSTRACT

We present a new catalog of 9,318 Ly α emitter (LAE) candidates at $z = 2.2, 3.3, 4.9, 5.7, 6.6,$ and 7.0 that are photometrically selected by the SILVERRUSH program with a machine learning technique from large area (up to 25.0 deg^2) imaging data with six narrowband filters taken by the Subaru Strategic Program with Hyper Suprime-Cam (HSC SSP) and a Subaru intensive program, Cosmic Hydrogen Reionization Unveiled with Subaru (CHORUS). We construct a convolutional neural network that distinguishes between real LAEs and contaminants with a completeness of 94% and a contamination rate of 1%, enabling us to efficiently remove contaminants from the photometrically selected LAE candidates. We confirm that our LAE catalogs include 177 LAEs that have been spectroscopically identified in our SILVERRUSH programs and previous studies, ensuring the validity of our machine learning selection. In addition, we find that the object-matching rates between our LAE catalogs and our previous results are $\simeq 80\text{--}100\%$ at bright NB magnitudes of $\lesssim 24$ mag. We also confirm that the surface number densities of our LAE candidates are consistent with previous results. Our LAE catalogs will be made public on our project webpage.

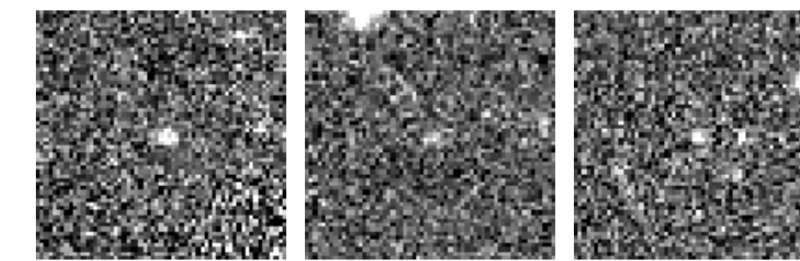
Keywords: galaxies: formation — galaxies: evolution — galaxies: high-redshift



Simulated LAEs

Satellite trails & Noise

Class 1



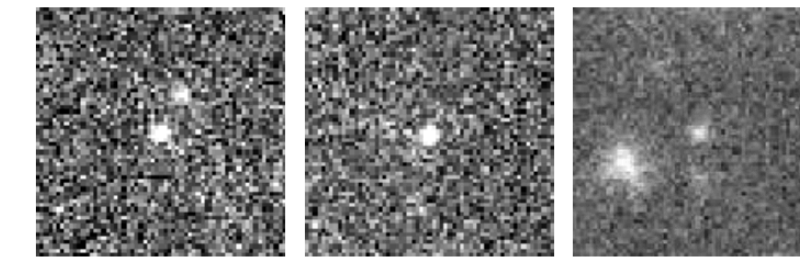
S/N: 3-4

Class 5



Bright

Class 2



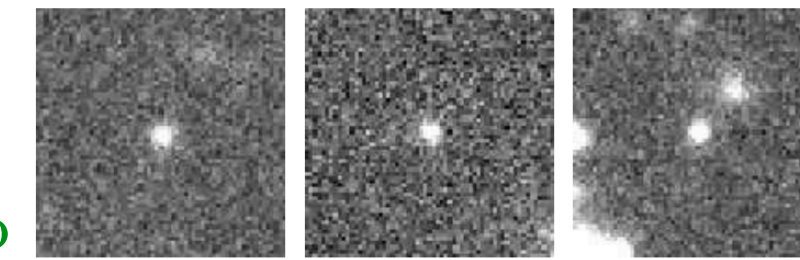
S/N: 4-10

Class 6



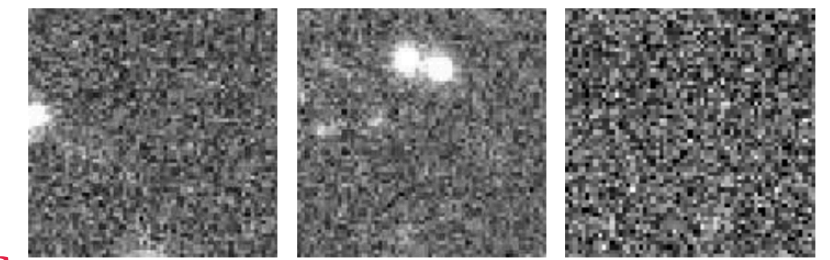
Faint

Class 3



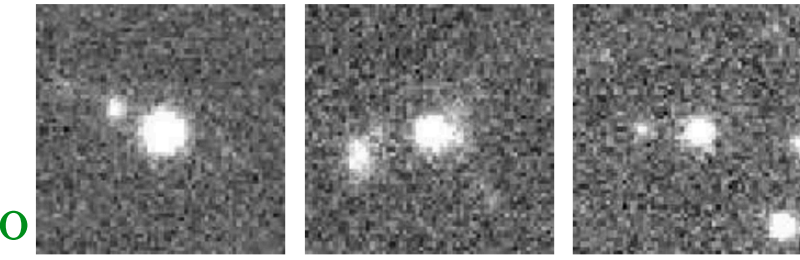
S/N: 10-30

Class 7



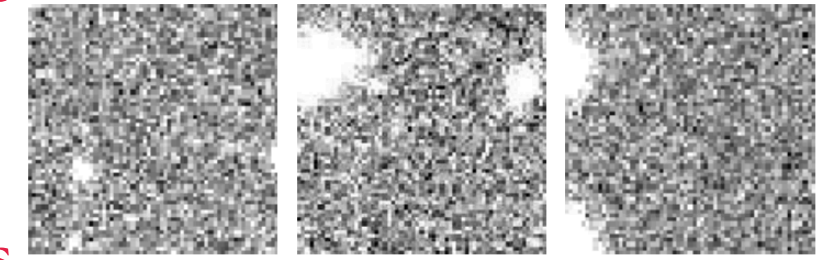
W/o sky residuals

Class 4



S/N: 30-200

Class 8



W/ sky residuals

Introducing & Motivation: SILVERRUSH

Ono et al. 2022

ACCEPTED FOR PUBLICATION IN APJ
Preprint typeset using L^AT_EX style emulateapj v. 12/16/11

9,318 new LAE candidates @ $2.2 < z < 7.0$

177 of them with spectroscopic confirmation



SILVERRUSH X: MACHINE LEARNING-AIDED SELECTION OF 9,318 LAES
AT $z = 2.2, 3.3, 4.9, 5.7, 6.6,$ AND 7.0 FROM THE HSC SSP AND CHORUS SURVEY DATA

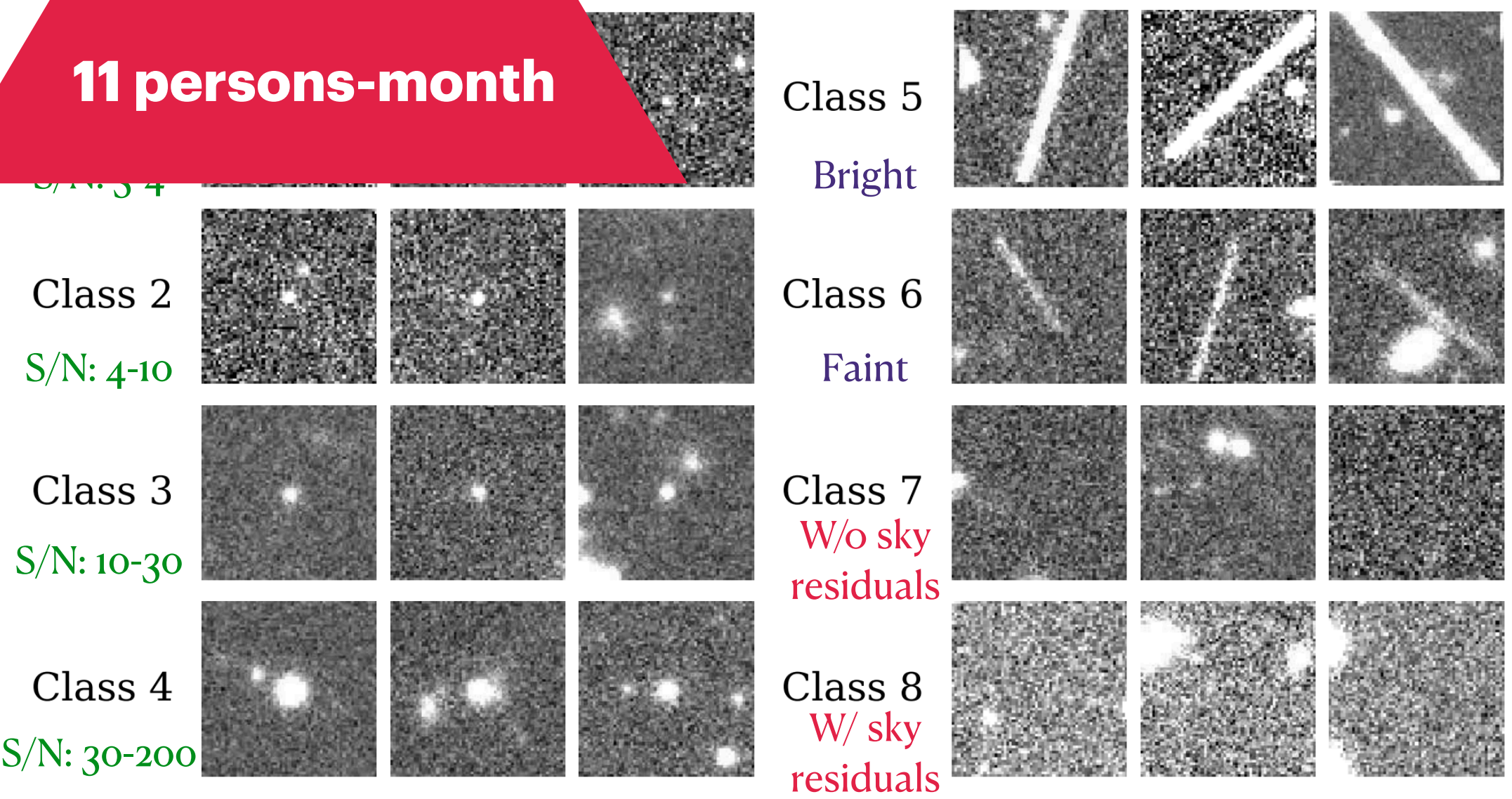
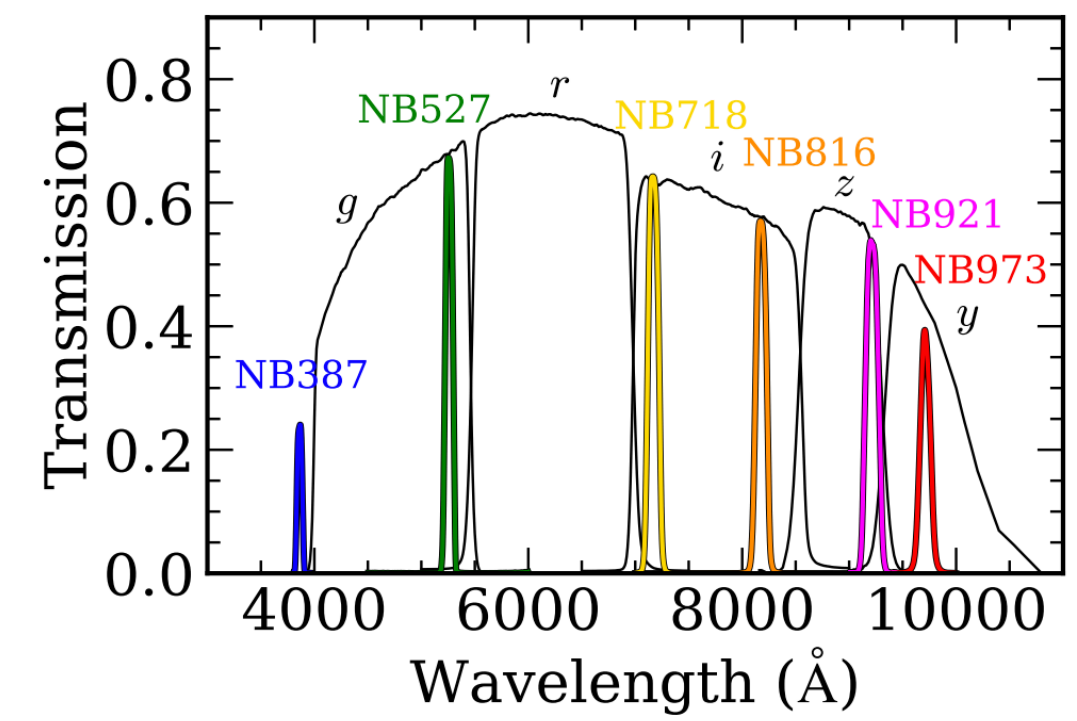
YOSHIAKI ONO¹, RYOHEI ITOH^{1,2}, TAKATOSHI SHIBUYA³, MASAMI OUCHI^{4,1,5}, YUICHI HARIKANE^{1,6}, SATOSHI YAMANAKA^{7,8}, AKIO K. INOUE^{9,7}, TOSHIYUKI AMAGASA^{10,11}, DAICHI MIURA¹⁰, MAIKI OKURA¹⁰, KAZUHIRO SHIMASAKU^{12,13}, IKURU IWATA⁴, YOSHIAKI TANIGUCHI¹⁴, SEIJI FUJIMOTO¹⁵, MASANORI IYE⁴, ANTON T. JAELANI^{16,17}, NOBUNARI KASHIKAWA^{12,13}, SHOTARO KIKUCHIHARA^{1,12}, SATOSHI KIKUTA¹¹, MASAKAZU A.R. KOBAYASHI¹⁸, HARUKA KUSAKABE¹⁹, CHIEN-HSIU LEE²⁰, YONGMING LIANG⁴, YOSHIKI MATSUOKA⁸, RIEKO MOMOSE¹², TOHRU NAGAO⁸, KIMIHIKO NAKAJIMA⁴, AND KEN-ICHI TADAKI⁴

Accepted for publication in ApJ

ABSTRACT

We present a new catalog of 9,318 Ly α emitter (LAE) candidates at $z = 2.2, 3.3, 4.9, 5.7, 6.6,$ and 7.0 that are photometrically selected by the SILVERRUSH program with a machine learning technique from large area (up to 25.0 deg^2) imaging data with six narrowband filters taken by the Subaru Strategic Program with Hyper Suprime-Cam (HSC SSP) and a Subaru intensive program, Cosmic Hydrogen Reionization Unveiled with Subaru (CHORUS). We construct a convolutional neural network that distinguishes between real LAEs and contaminants with a completeness of 94% and a contamination rate of 1%, enabling us to efficiently remove contaminants from the photometrically selected LAE candidates. We confirm that our LAE catalogs include 177 LAEs that have been spectroscopically identified in our SILVERRUSH programs and previous studies, ensuring the validity of our machine learning selection. In addition, we find that the object-matching rates between our LAE catalogs and our previous results are $\simeq 80\text{--}100\%$ at bright NB magnitudes of $\lesssim 24$ mag. We also confirm that the surface number densities of our LAE candidates are consistent with previous results. Our LAE catalogs will be made public on our project webpage.

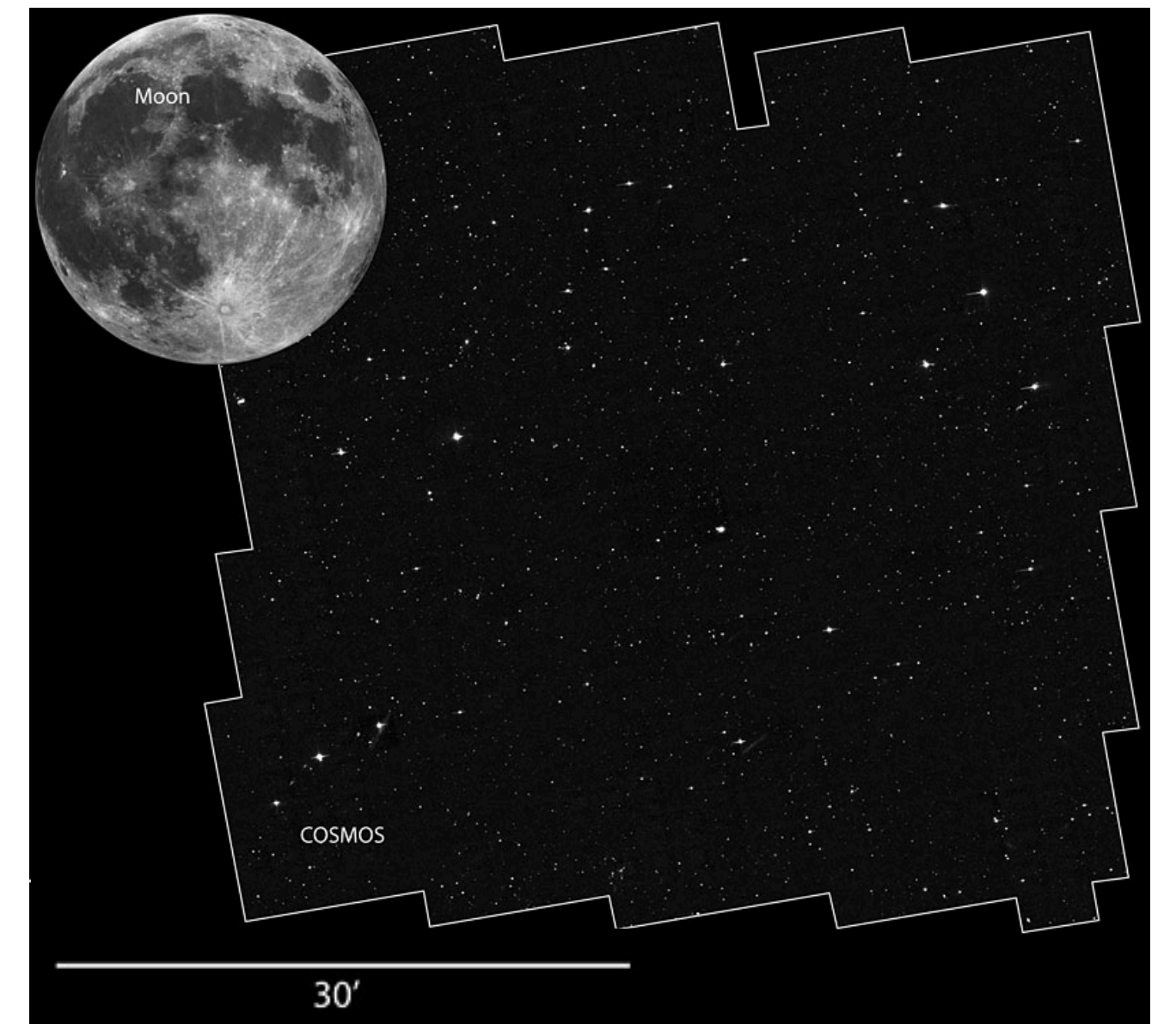
Keywords: galaxies: formation — galaxies: evolution — galaxies: high-redshift



[astro-ph.GA] 5 Apr 2021

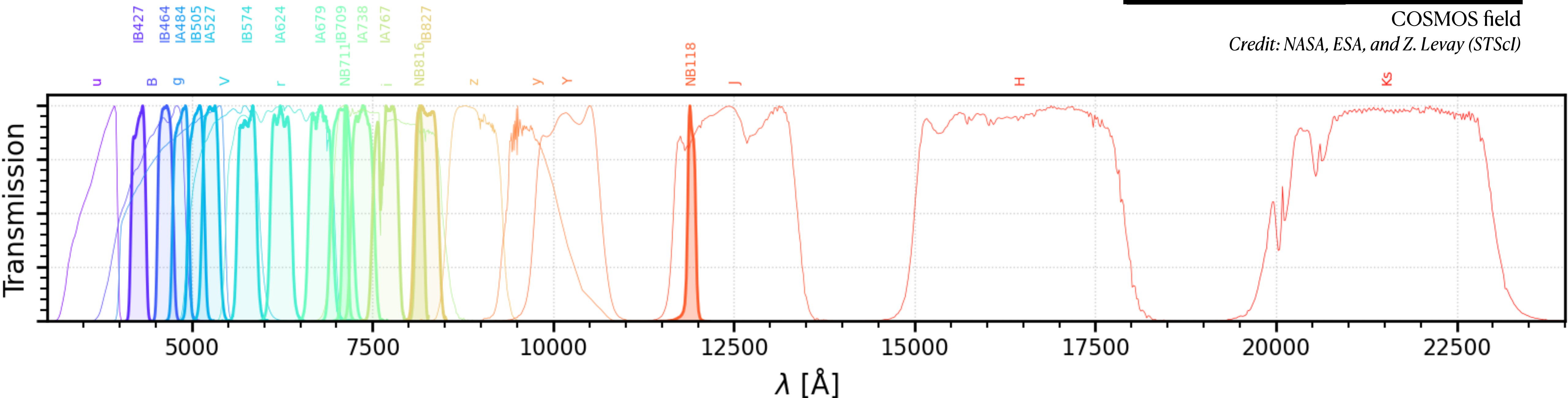
Data & Feature Engineering

- COSMOS 2020 ([Weaver et al. 2022a](#))
- SC4K 2018 ([Sobral et al. 2018a](#))



COSMOS field
Credit: NASA, ESA, and Z. Levay (STScI)

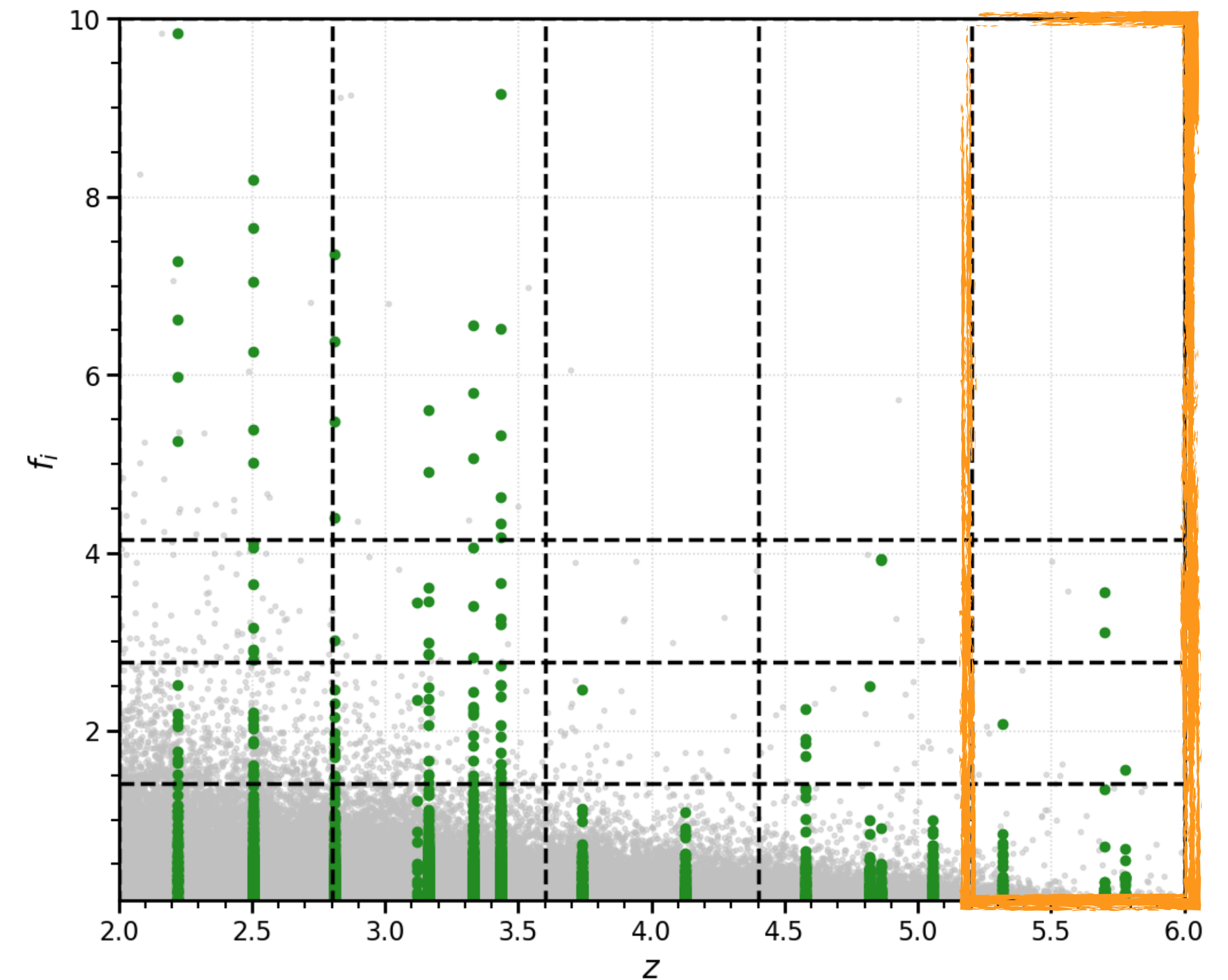
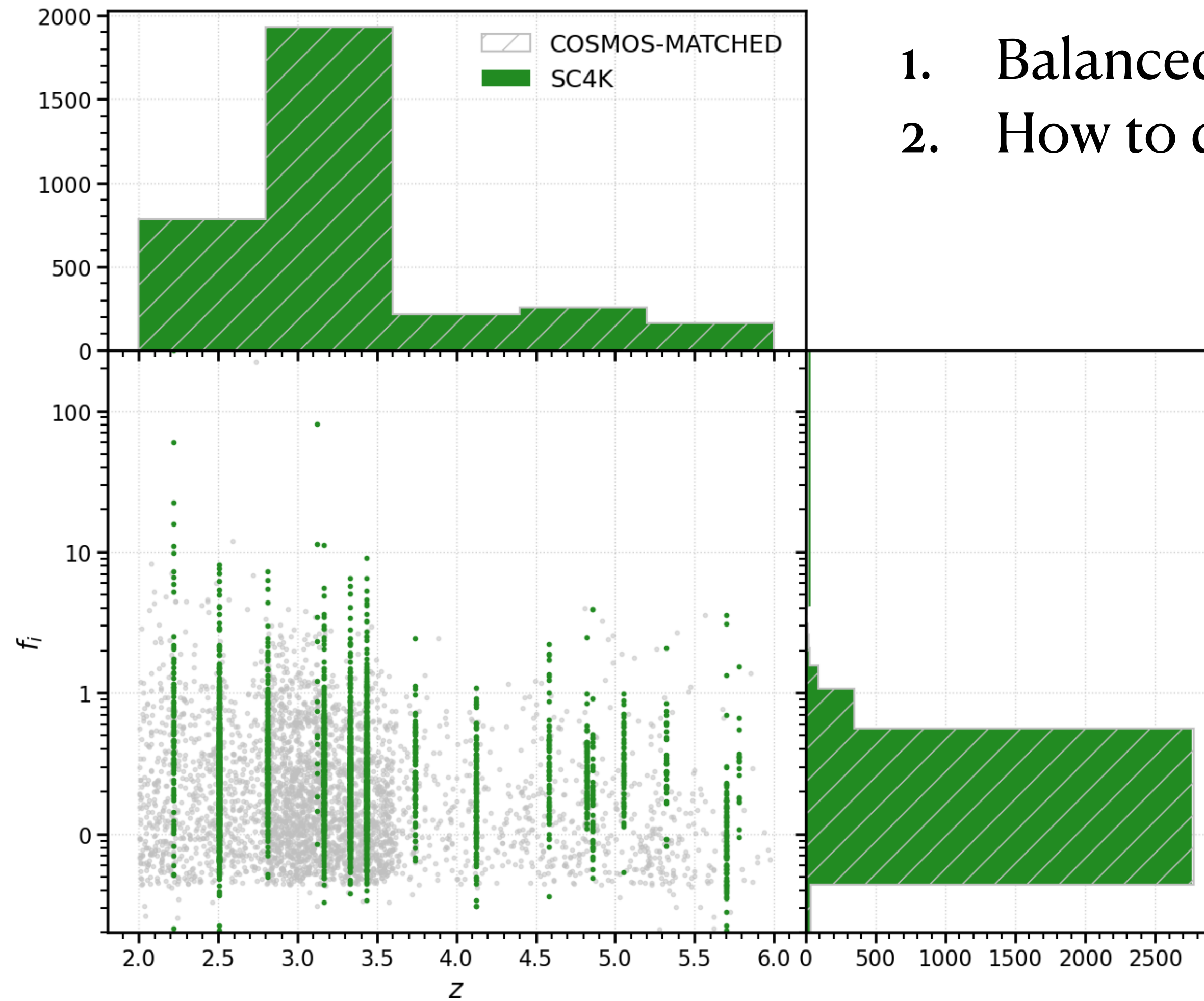
← Optical to Near-Infrared →
Fluxes + Uncertainties & Normalisation (unit-norm)



Data & Feature Engineering

1. Balanced or unbalanced data?
2. How to define a proper sample of non-LAEs?

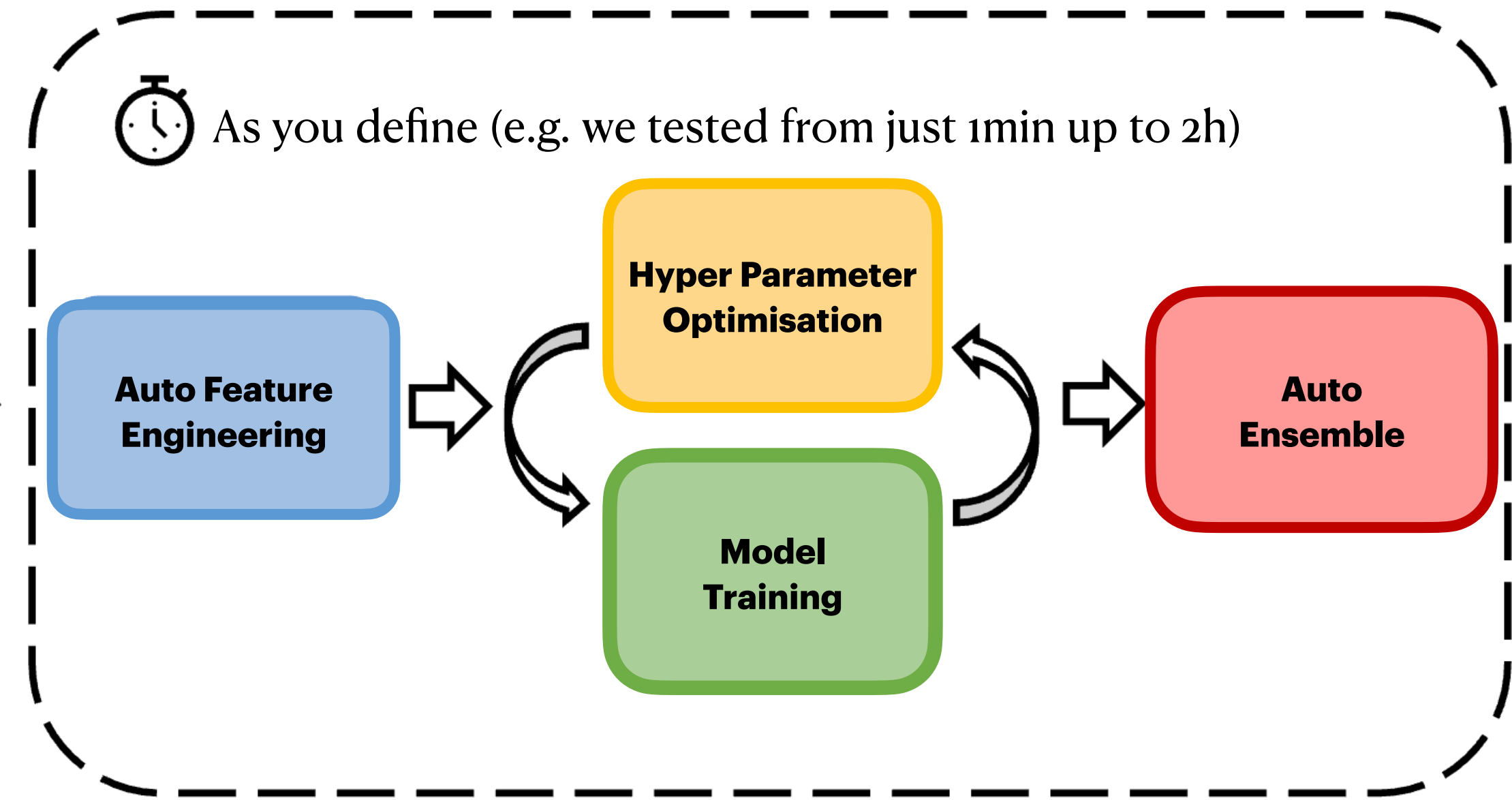
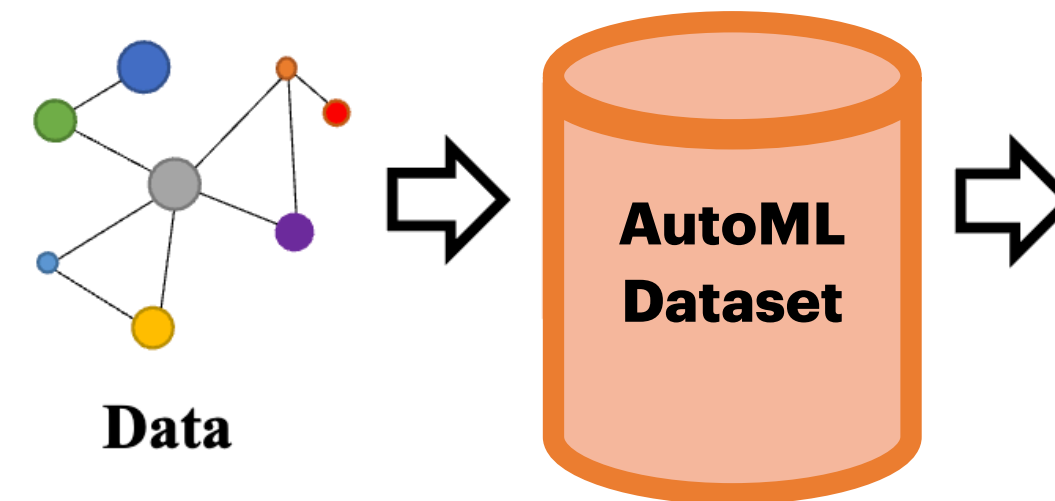
7 samples



Methods: H2O AutoML

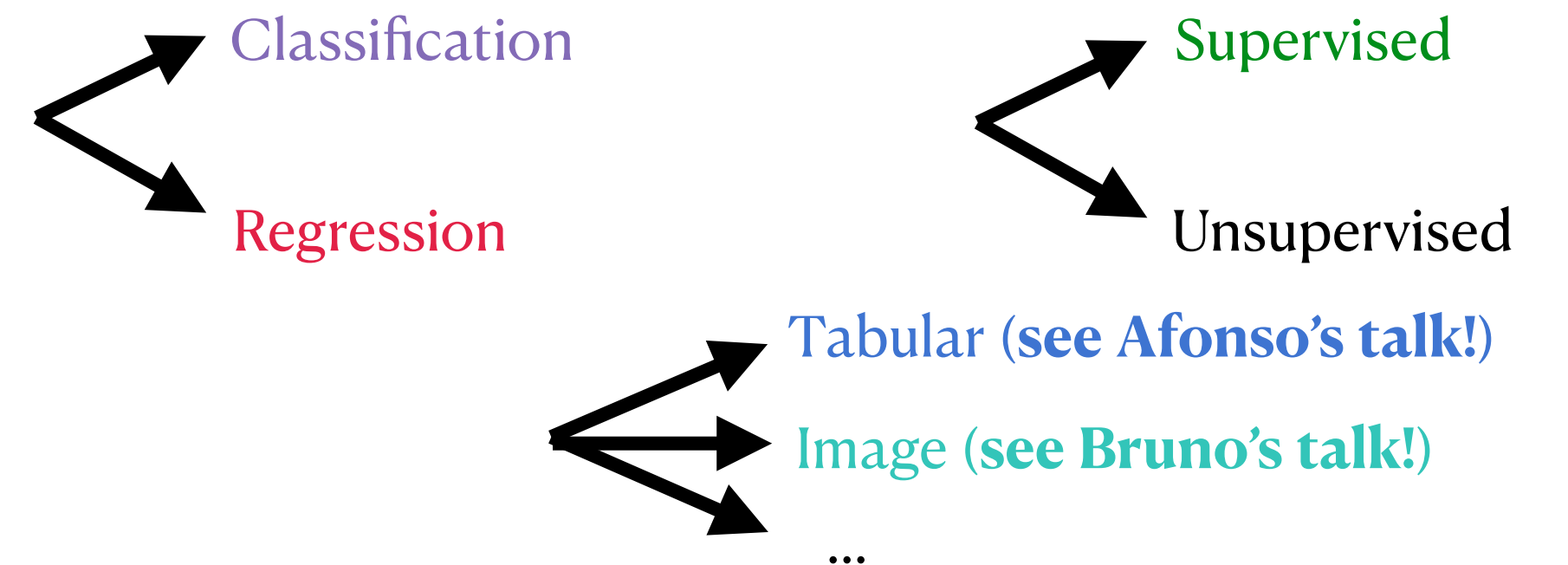


LeDell & Poirier 2020



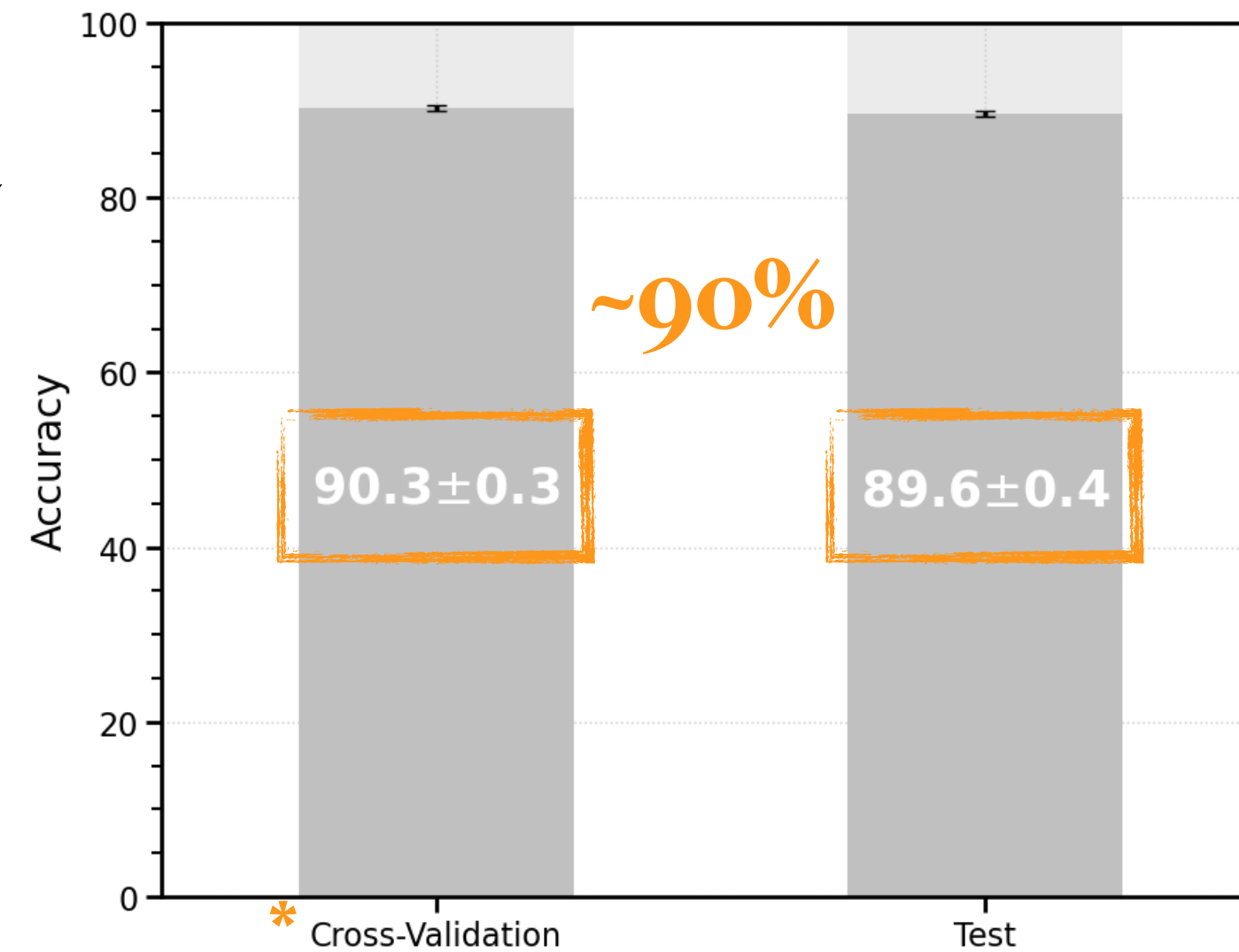
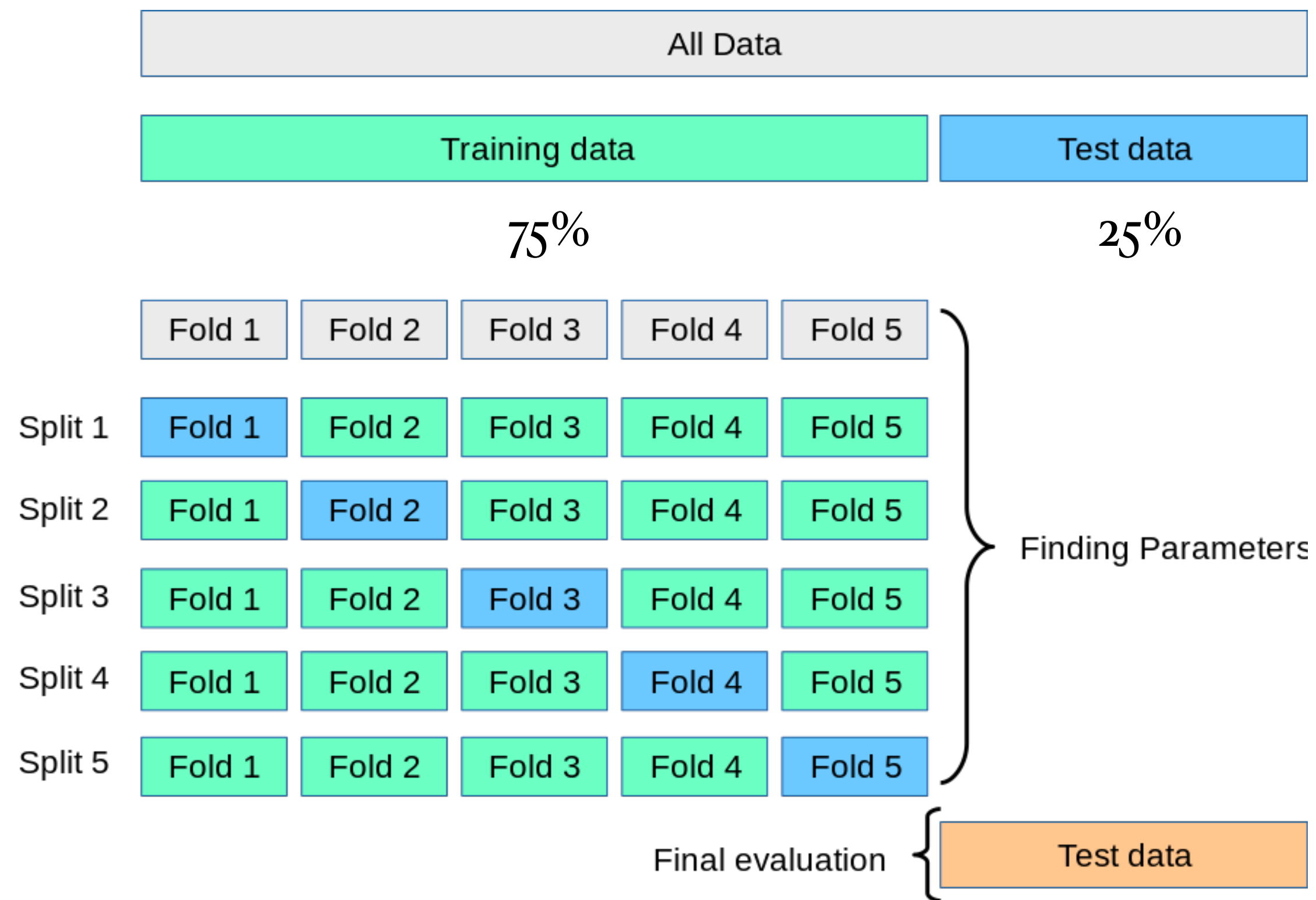
Algorithms:

- Distributed Random Forest (DRF, including both Random Forest and Extremely Randomised Trees (XRT) models)
- Generalised Linear Model (GLM)
- XGBoost
- **Gradient Boosting Machines (GBM)**
- DeepLearning
- Stacked Ensemble



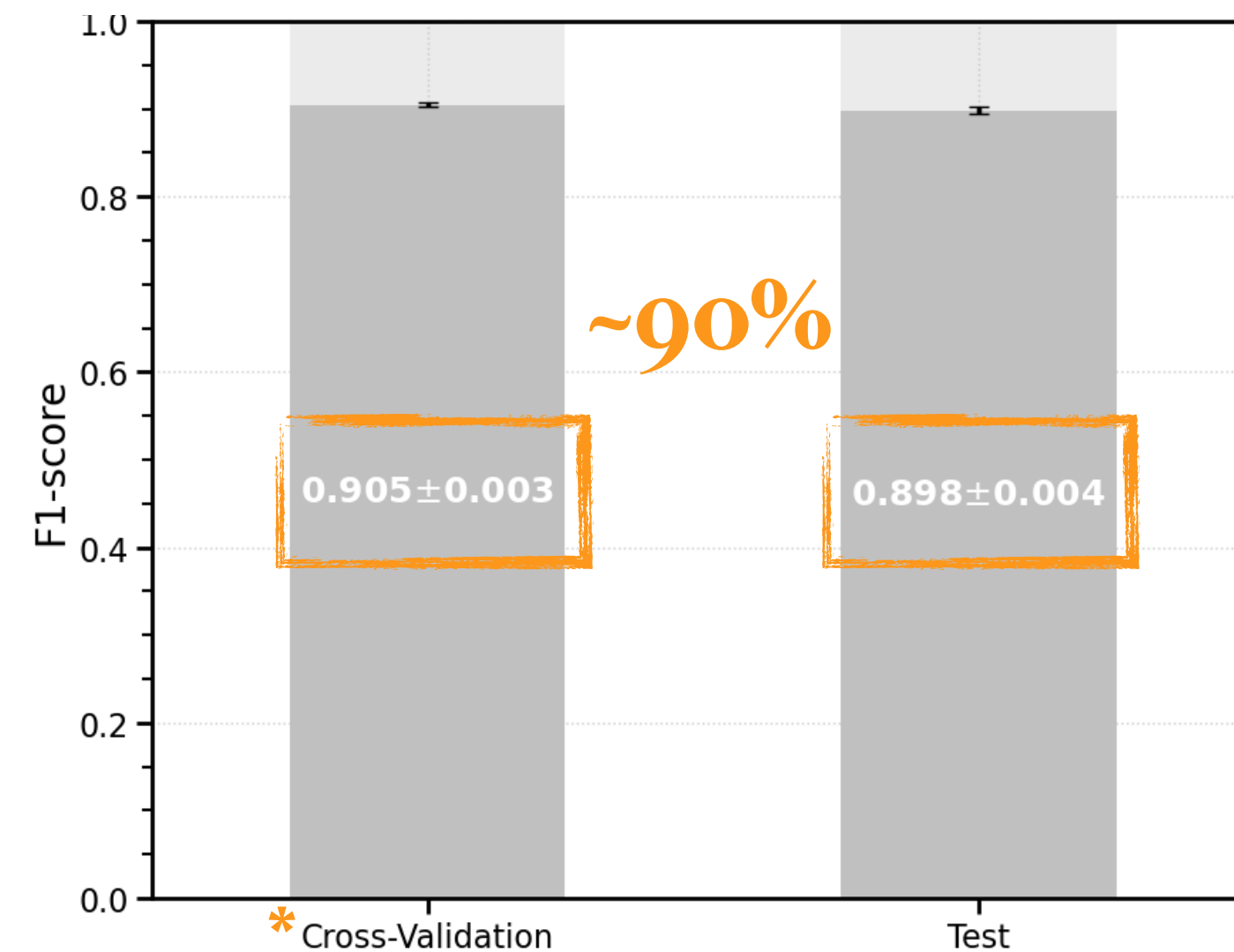
AutoML: evaluating the classification model

* **Cross-validation**: a resampling method that uses different portions of the data to test and train a model on different iterations. **Important** to protect against overfitting, particularly when the amount of data is limited.



Accuracy is the ratio of predictions that exactly match the true class labels.

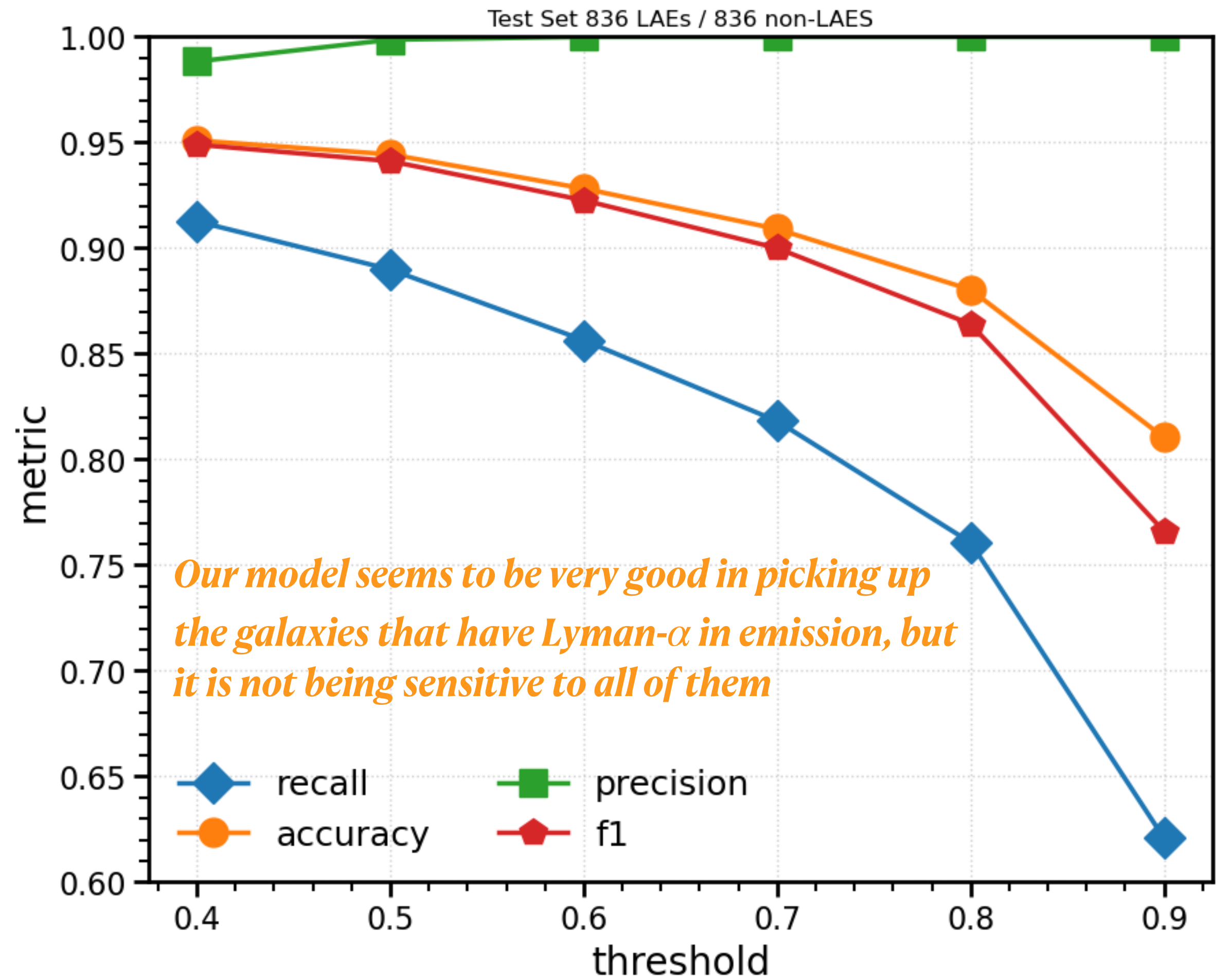
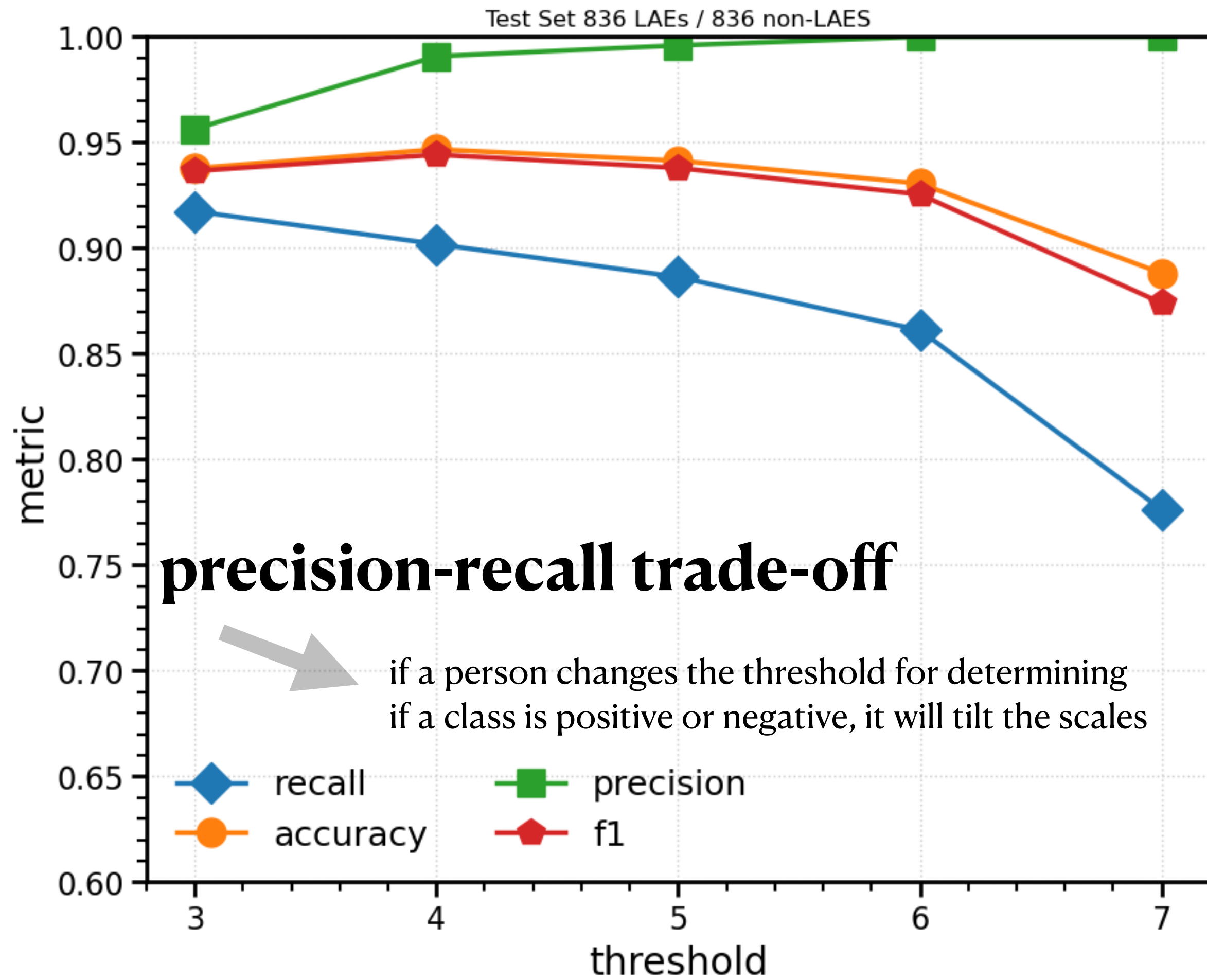
$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$



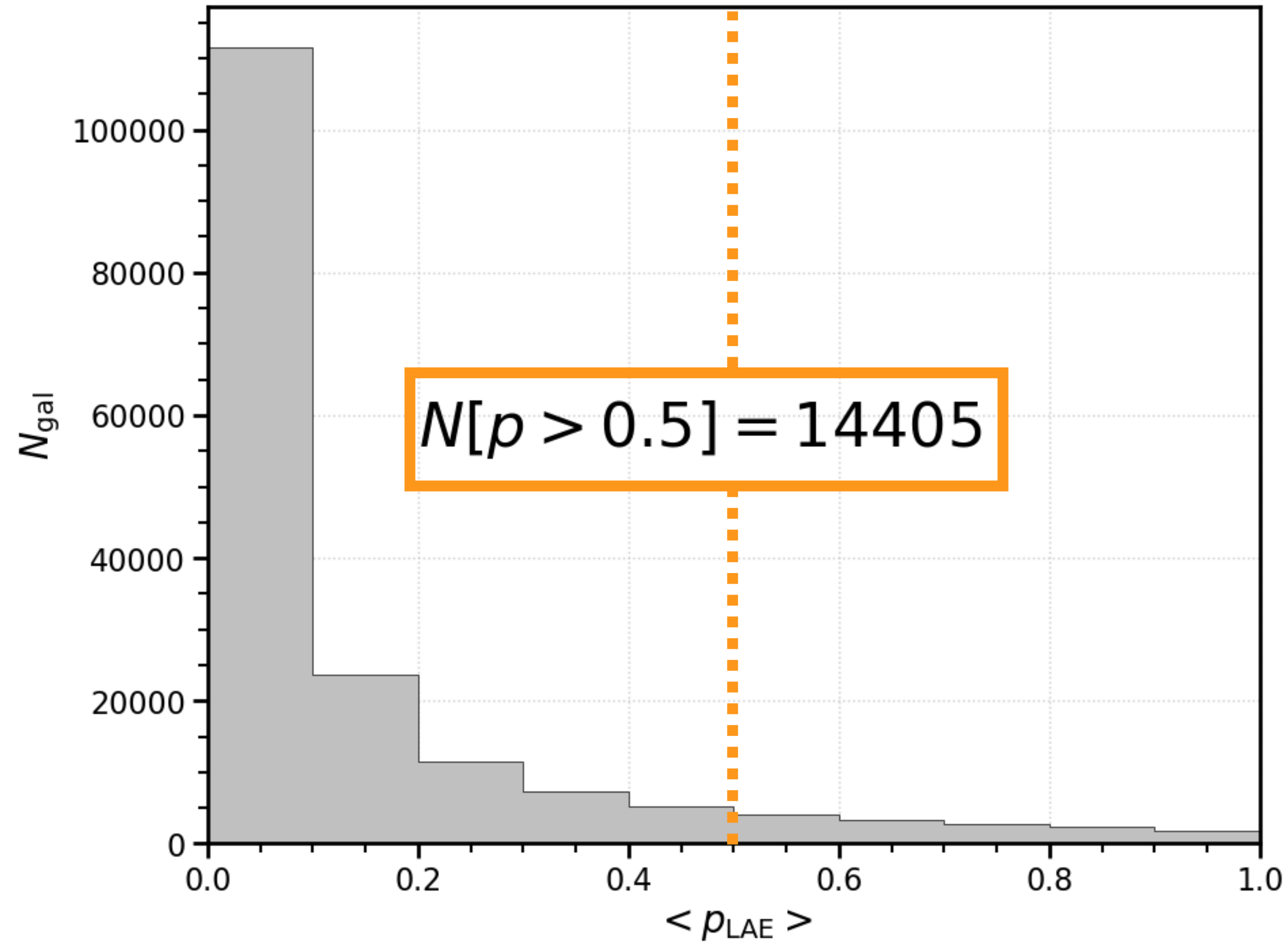
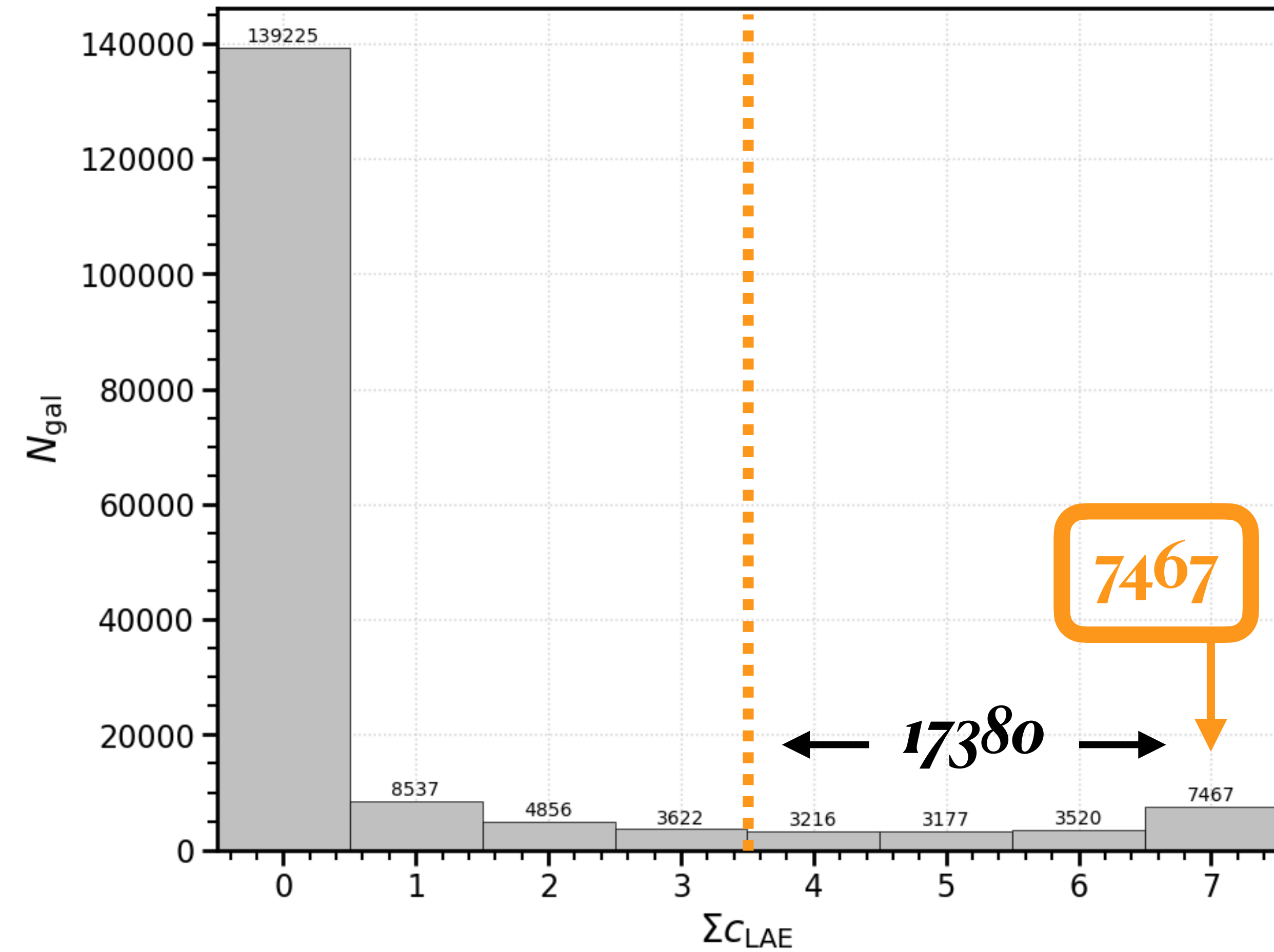
F1-score is a good balanced measure of both false positives and false negatives.

$$F1 \text{ Score} = \frac{2TP}{2TP + FP + FN}$$

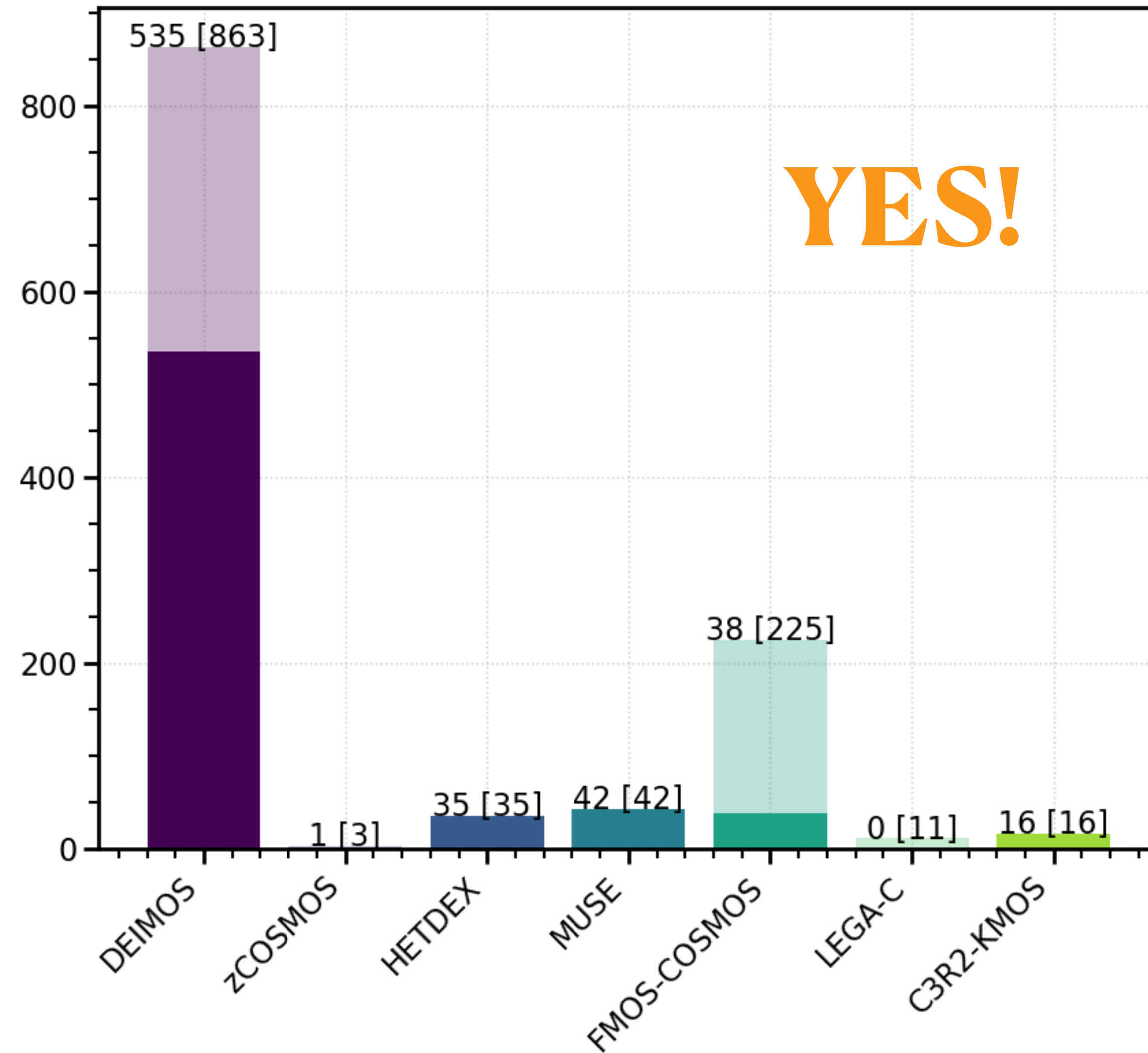
AutoML: evaluating the classification model



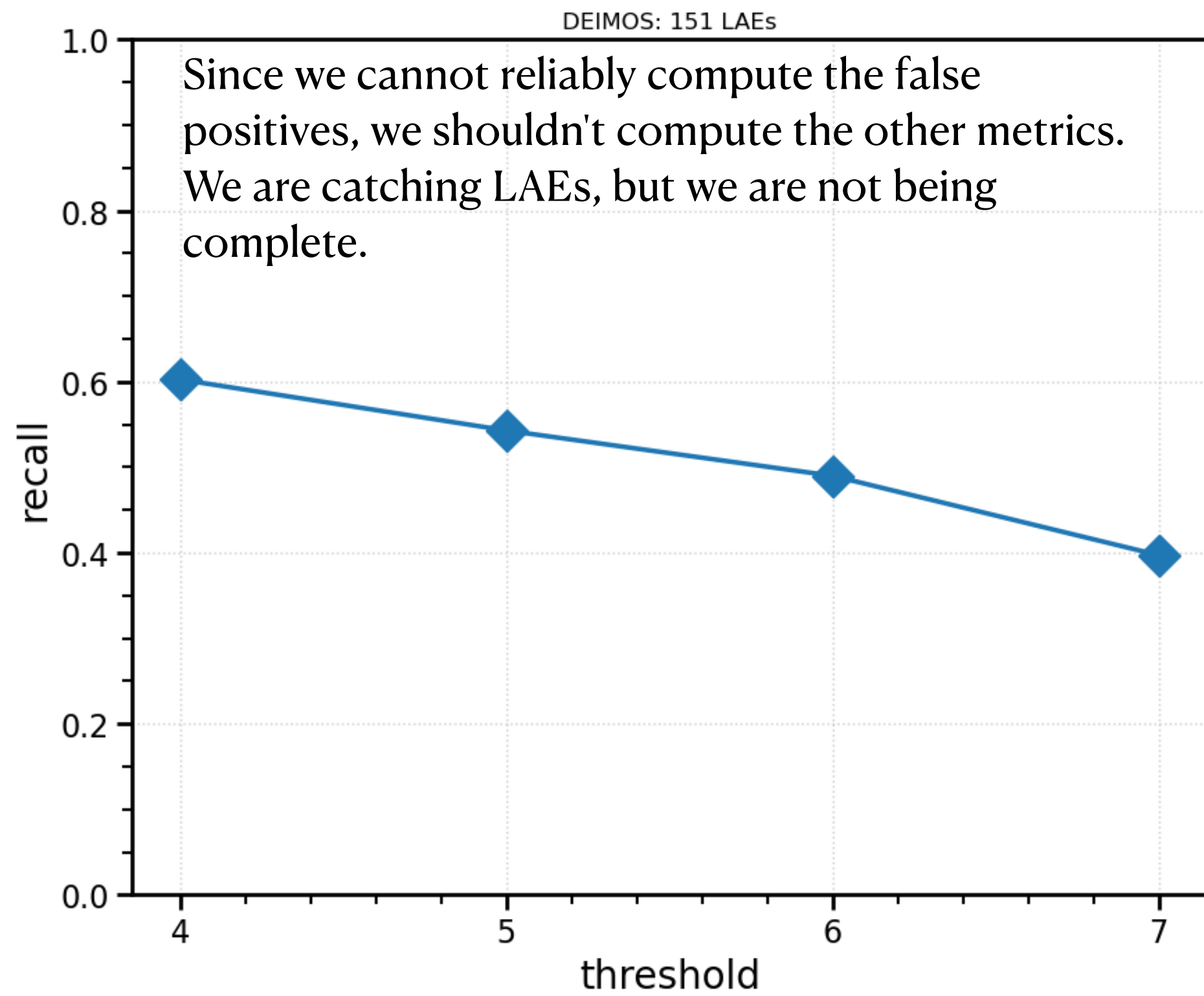
Results: how many LAEs are probably escaping us?



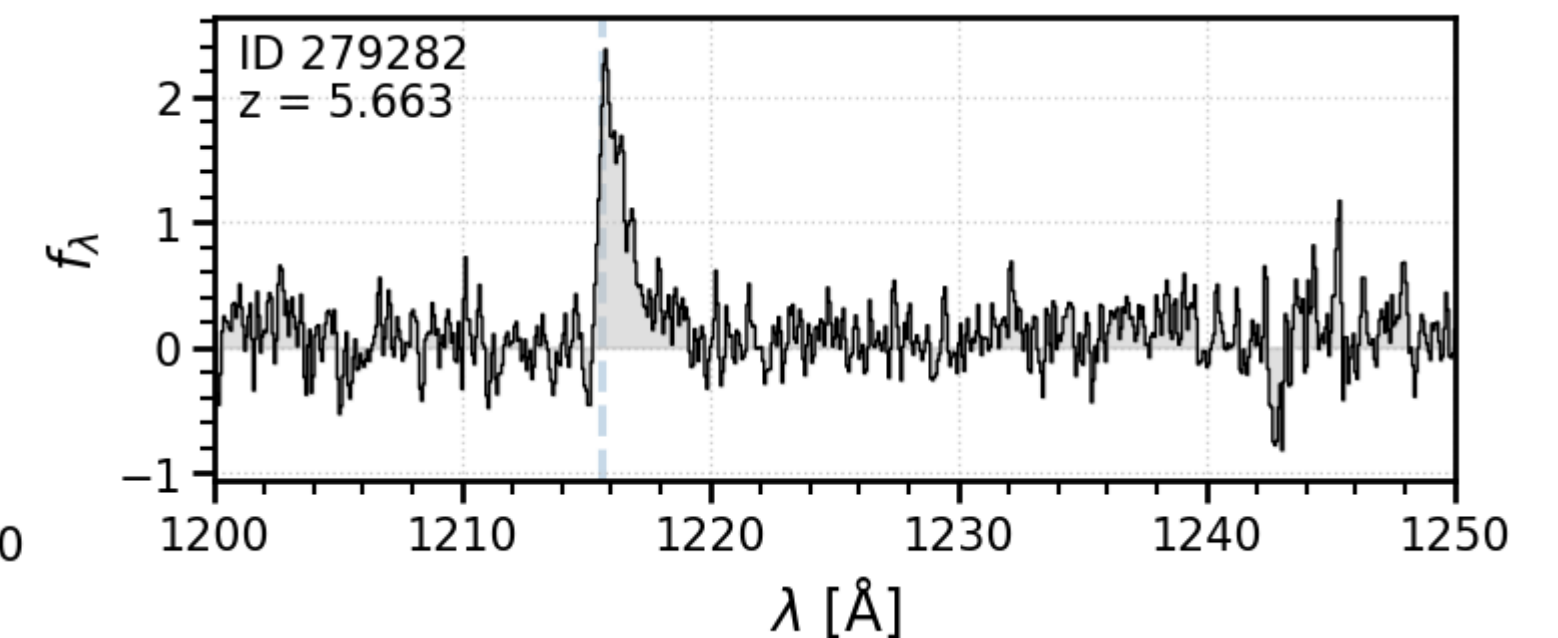
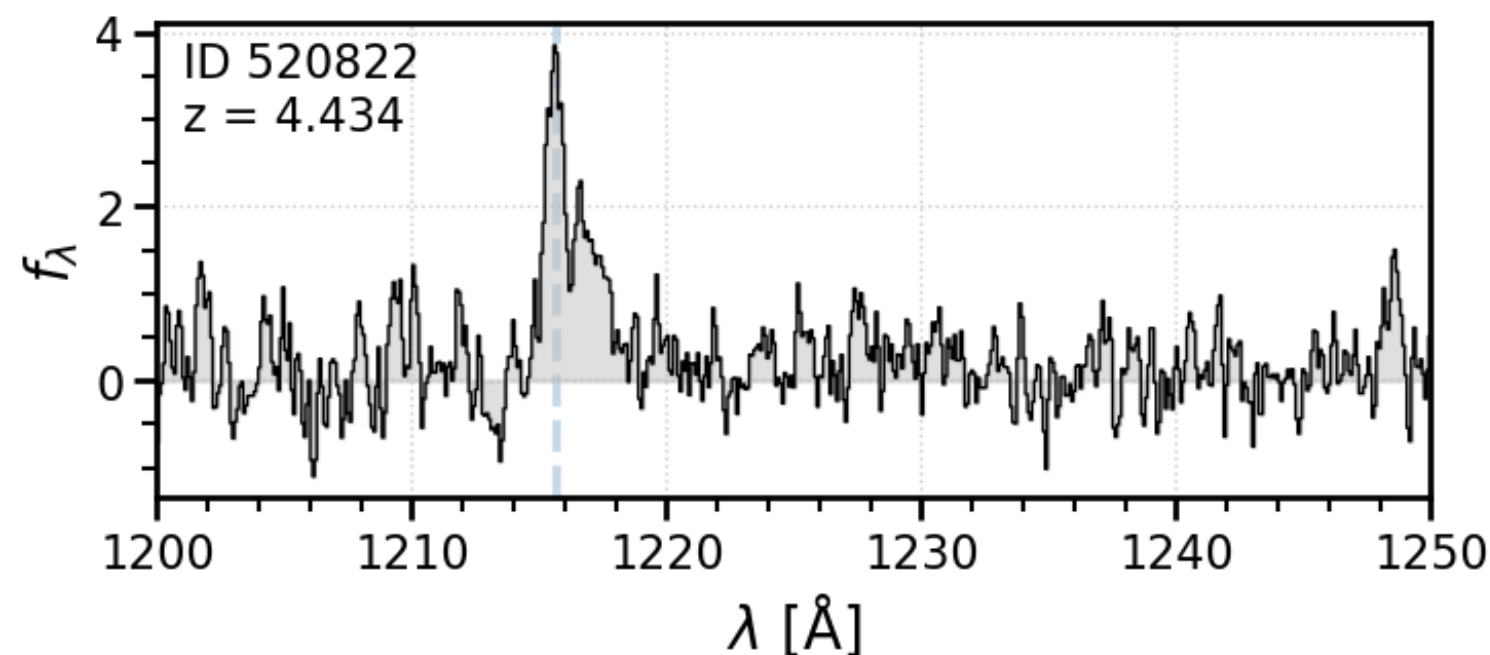
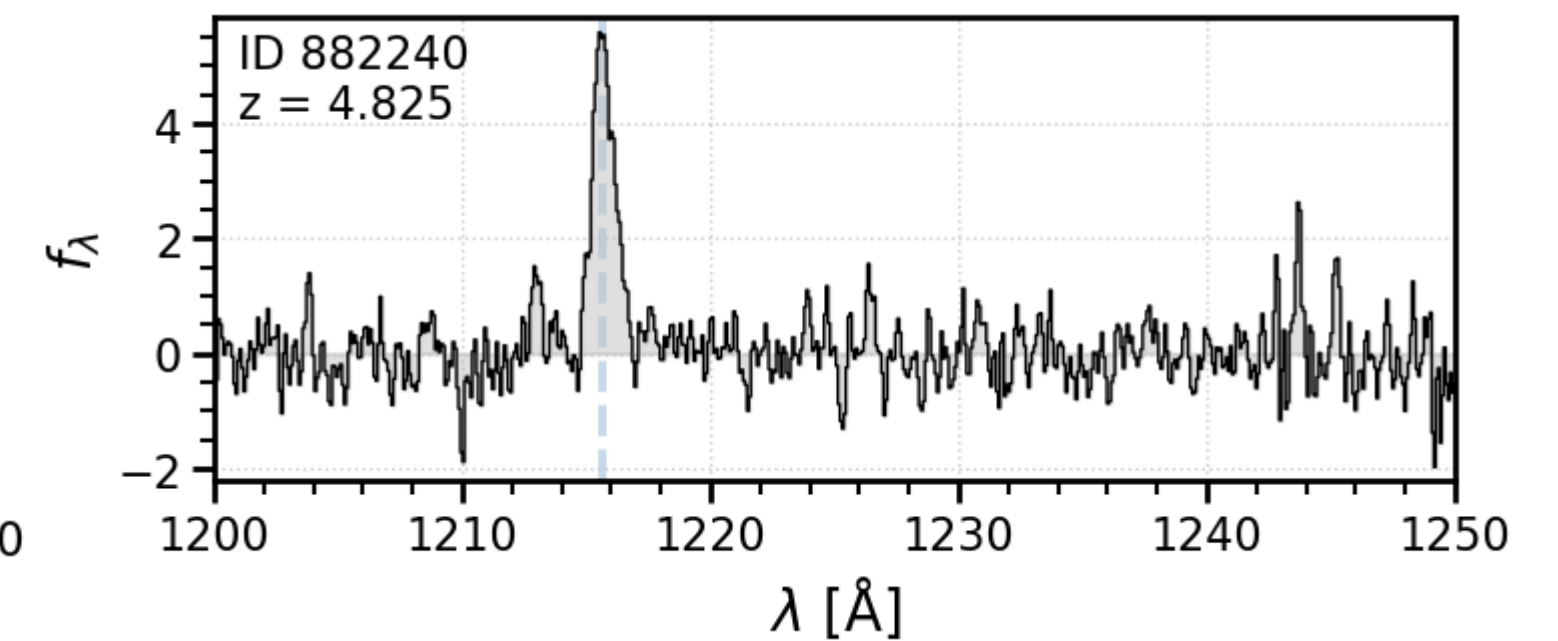
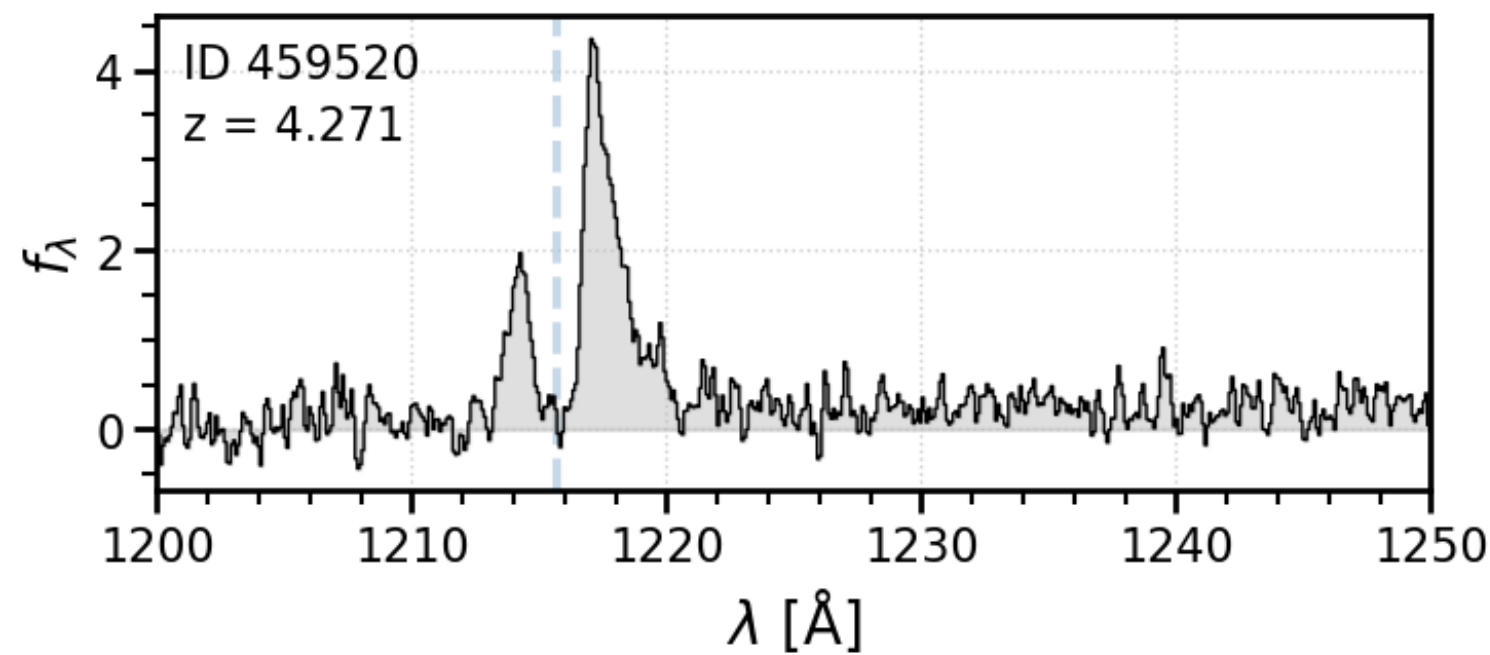
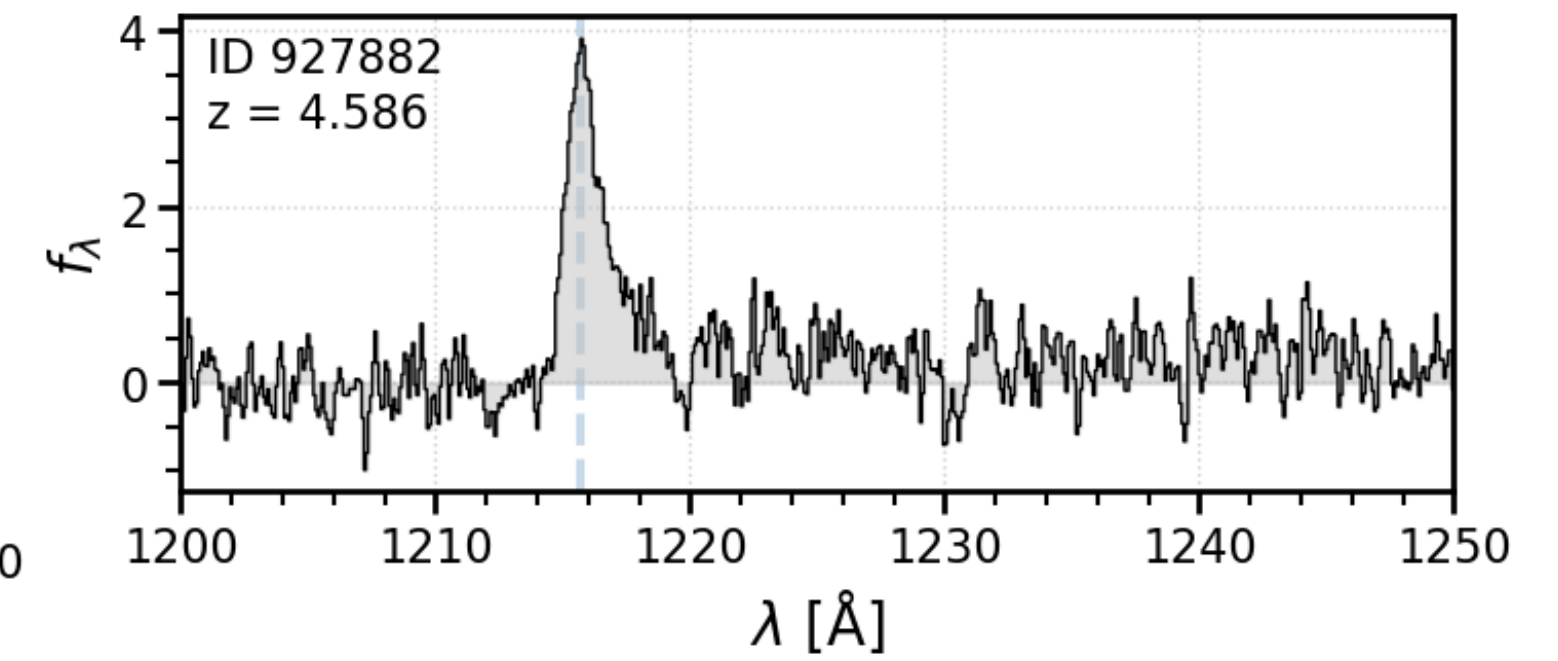
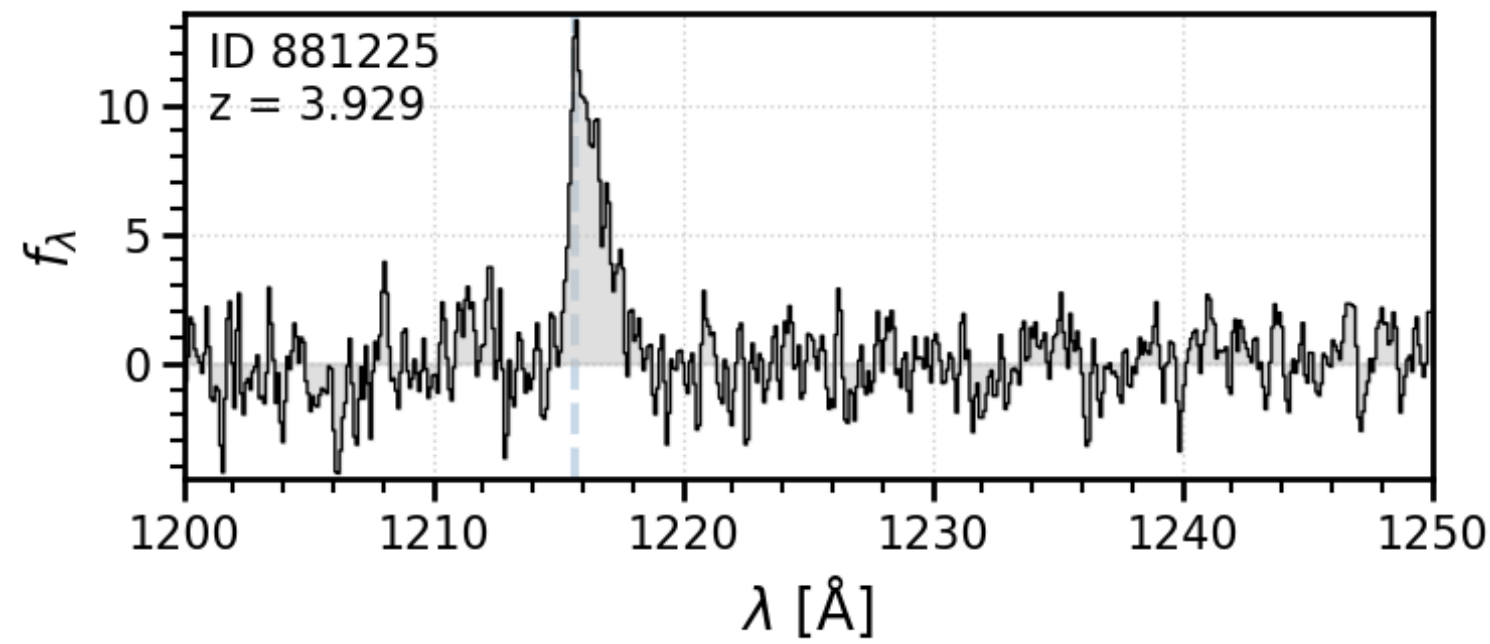
Results: can we spectroscopically confirm some of them?



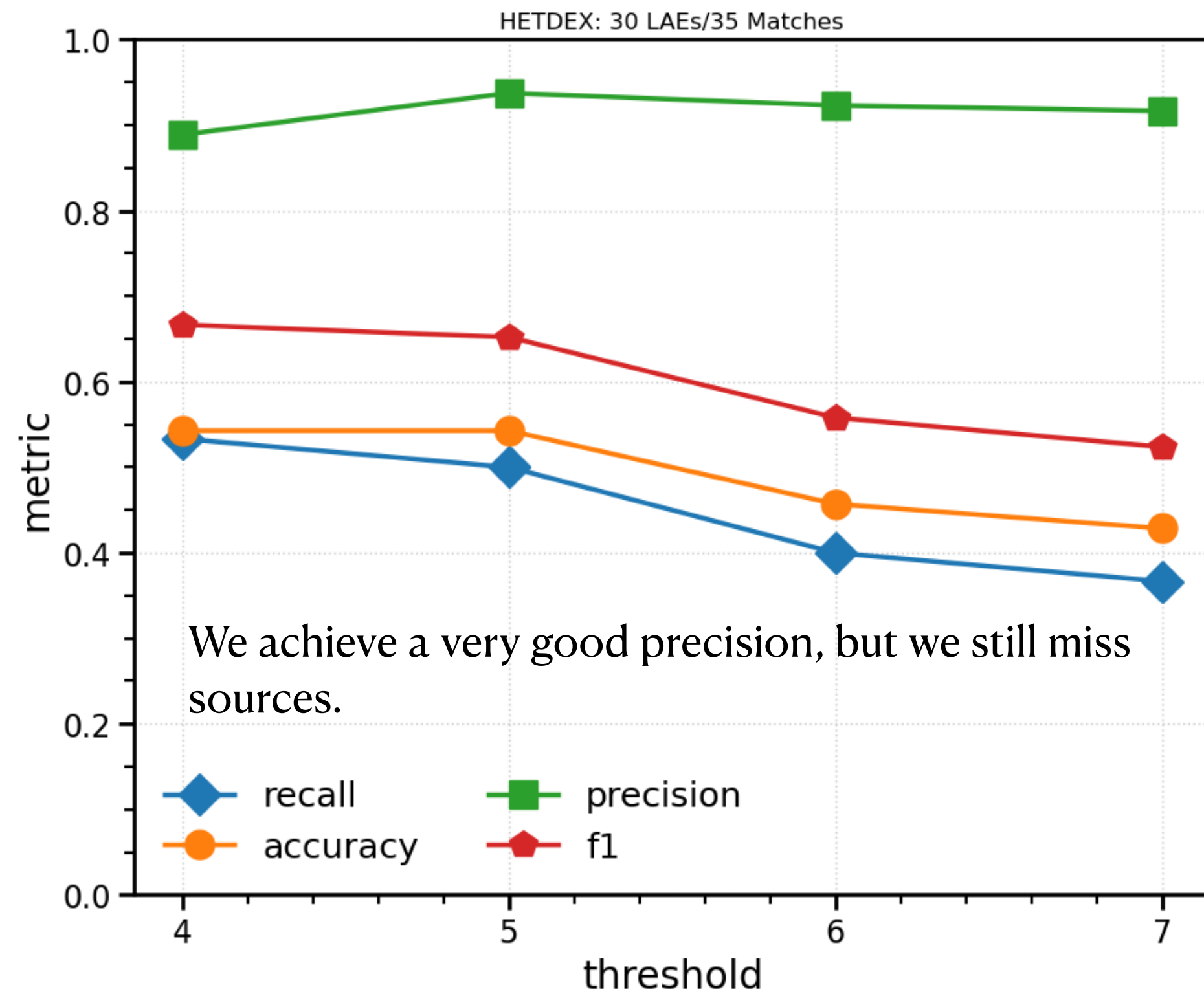
Results: DEIMOS spectroscopic confirmations



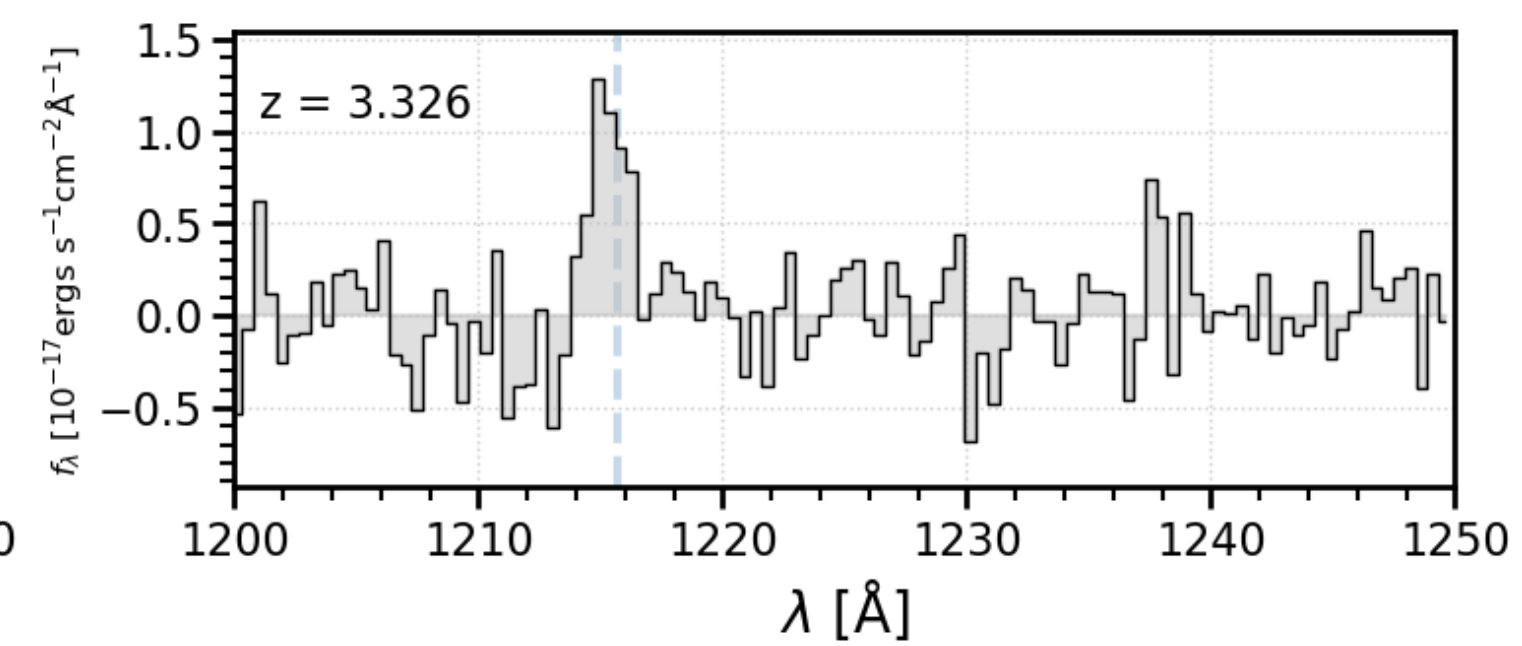
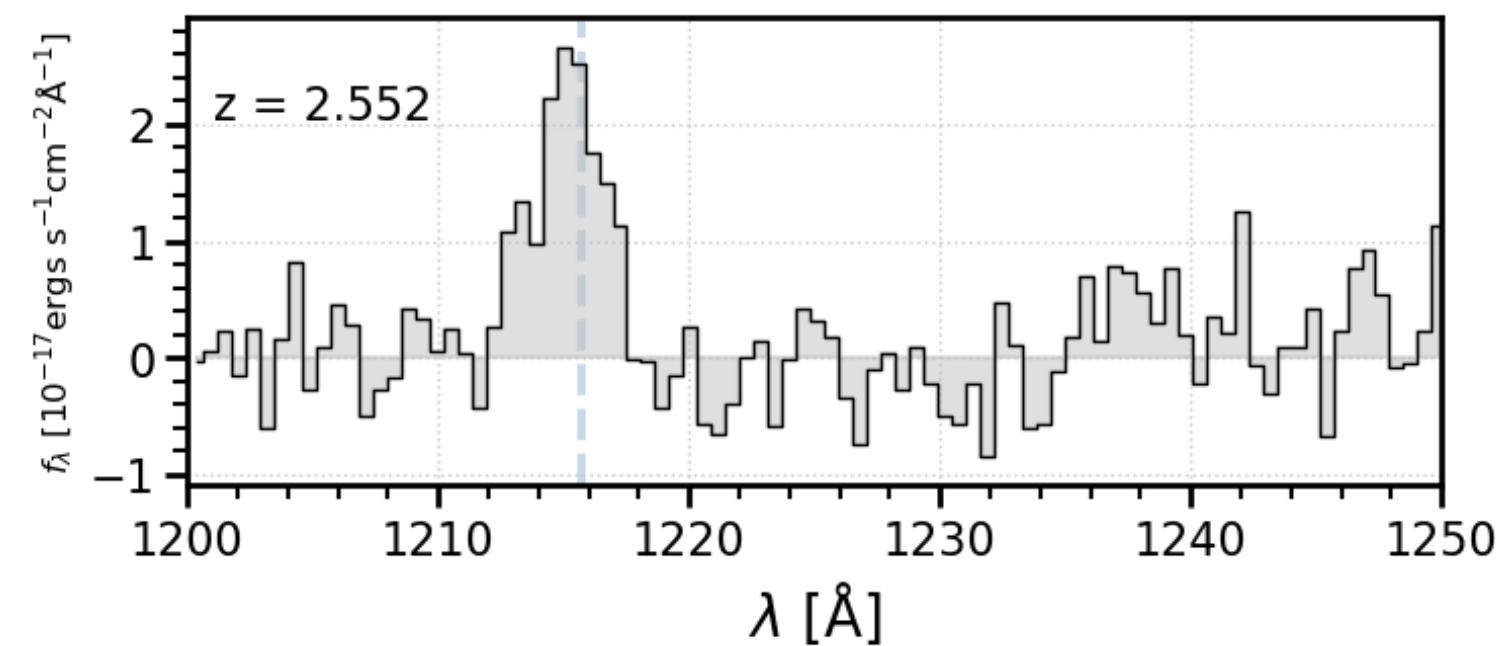
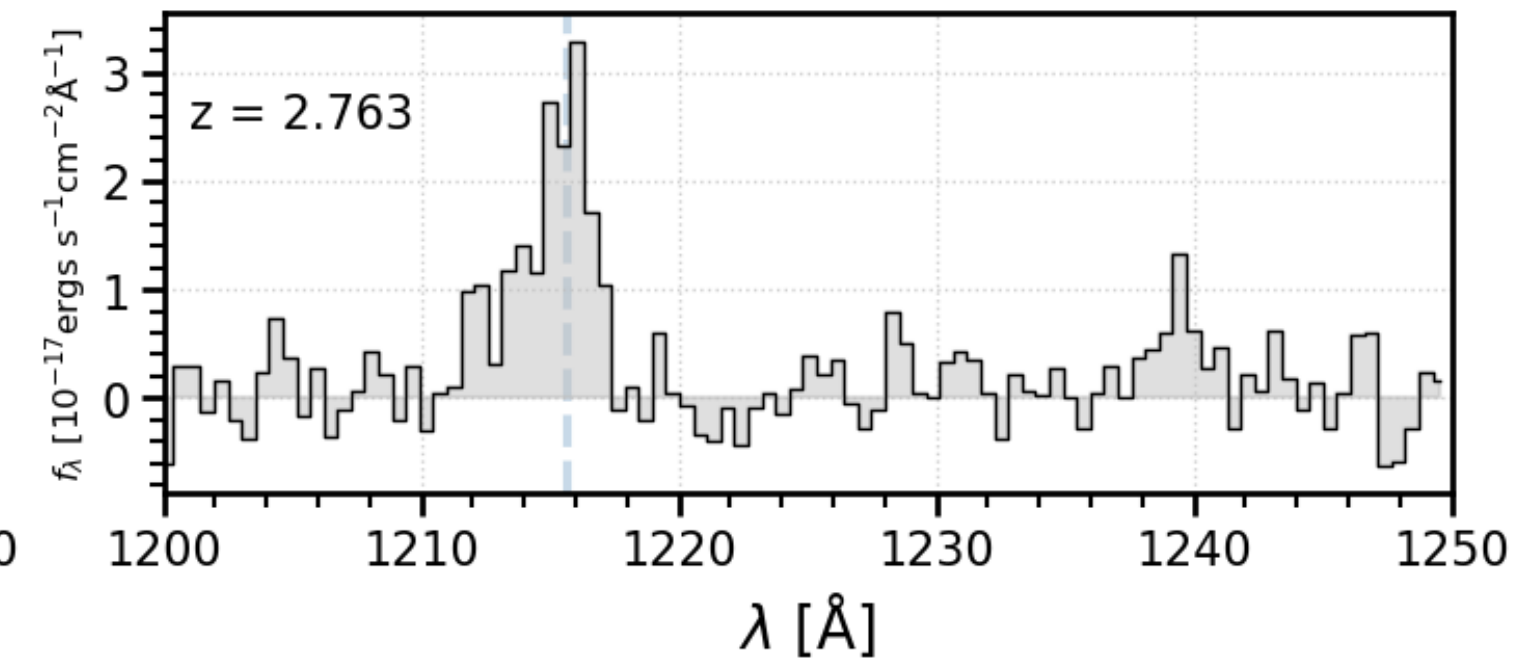
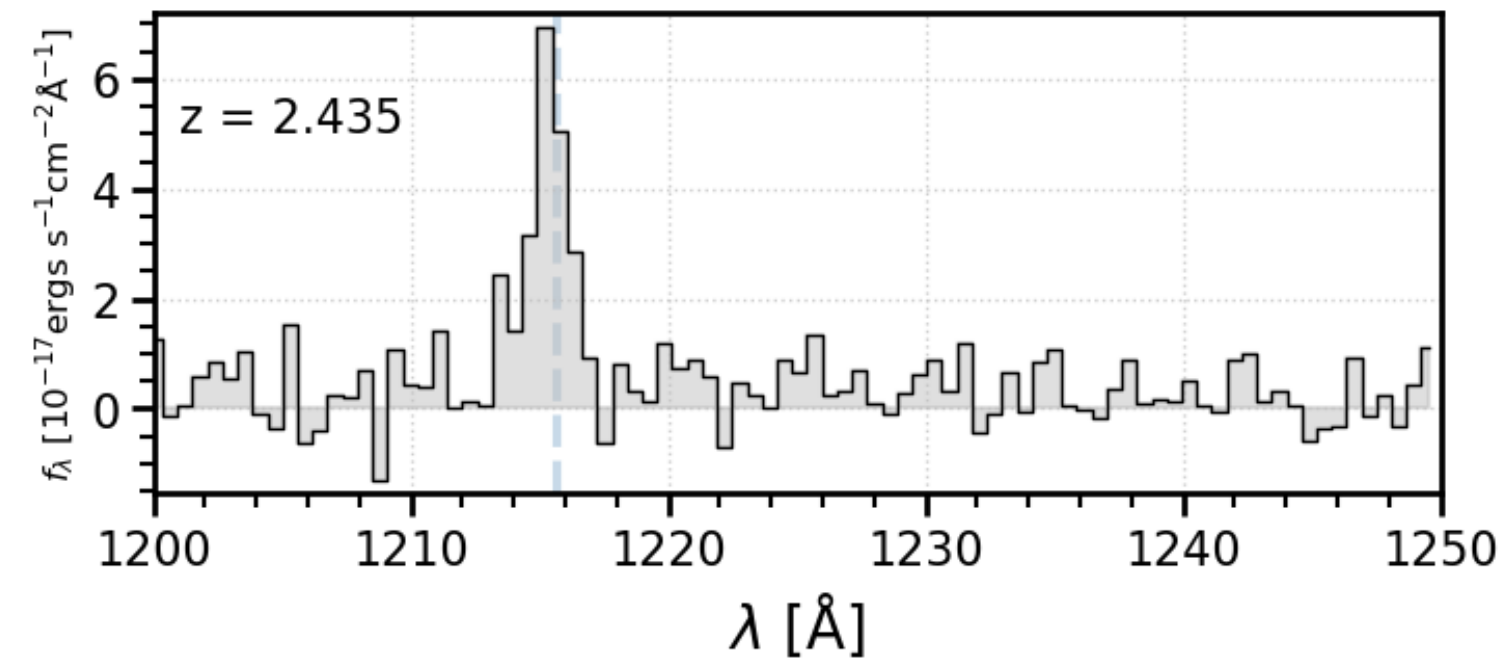
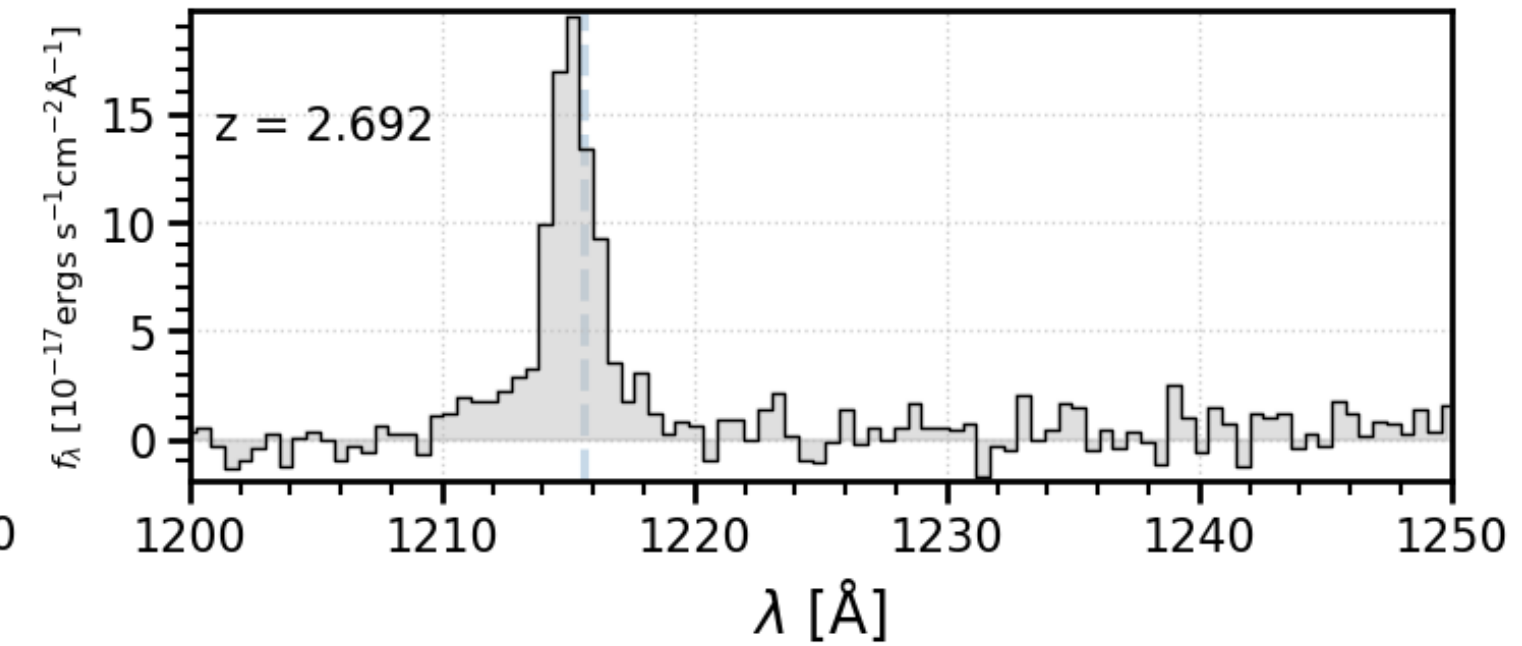
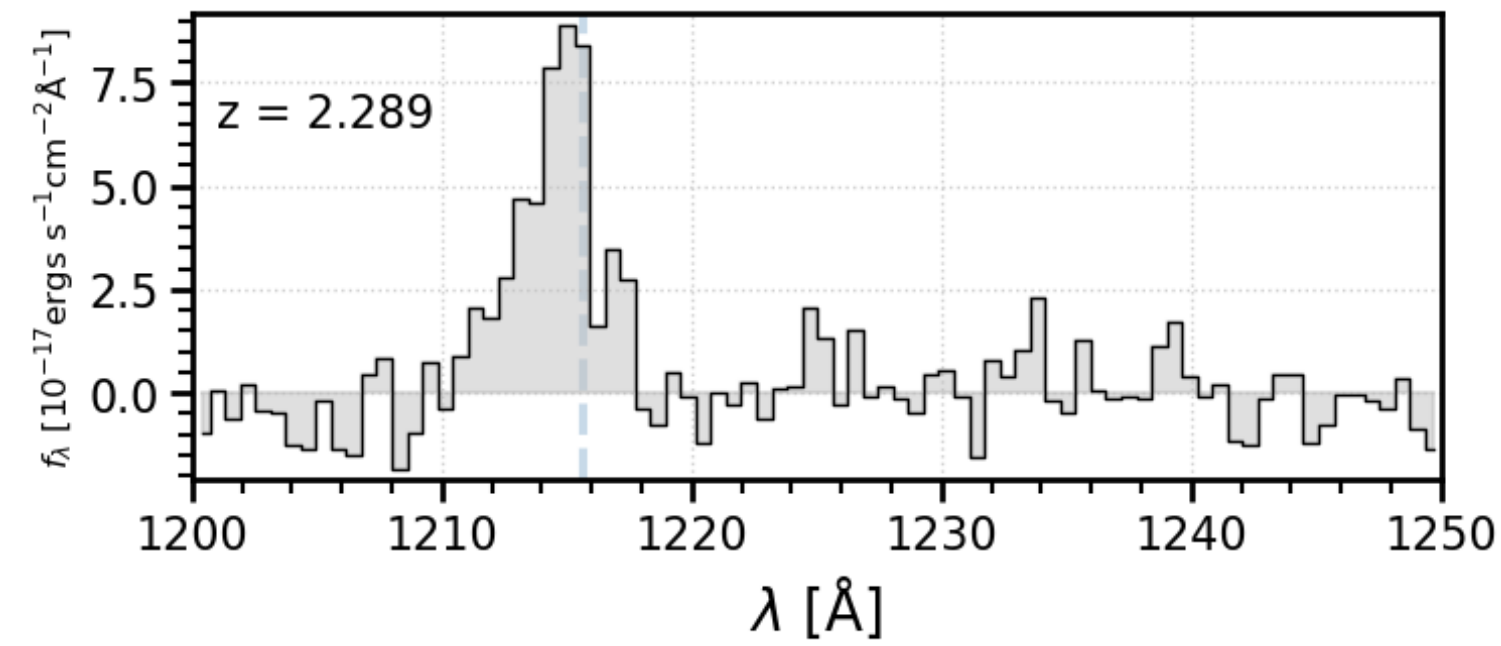
Hasinger et al. 2018



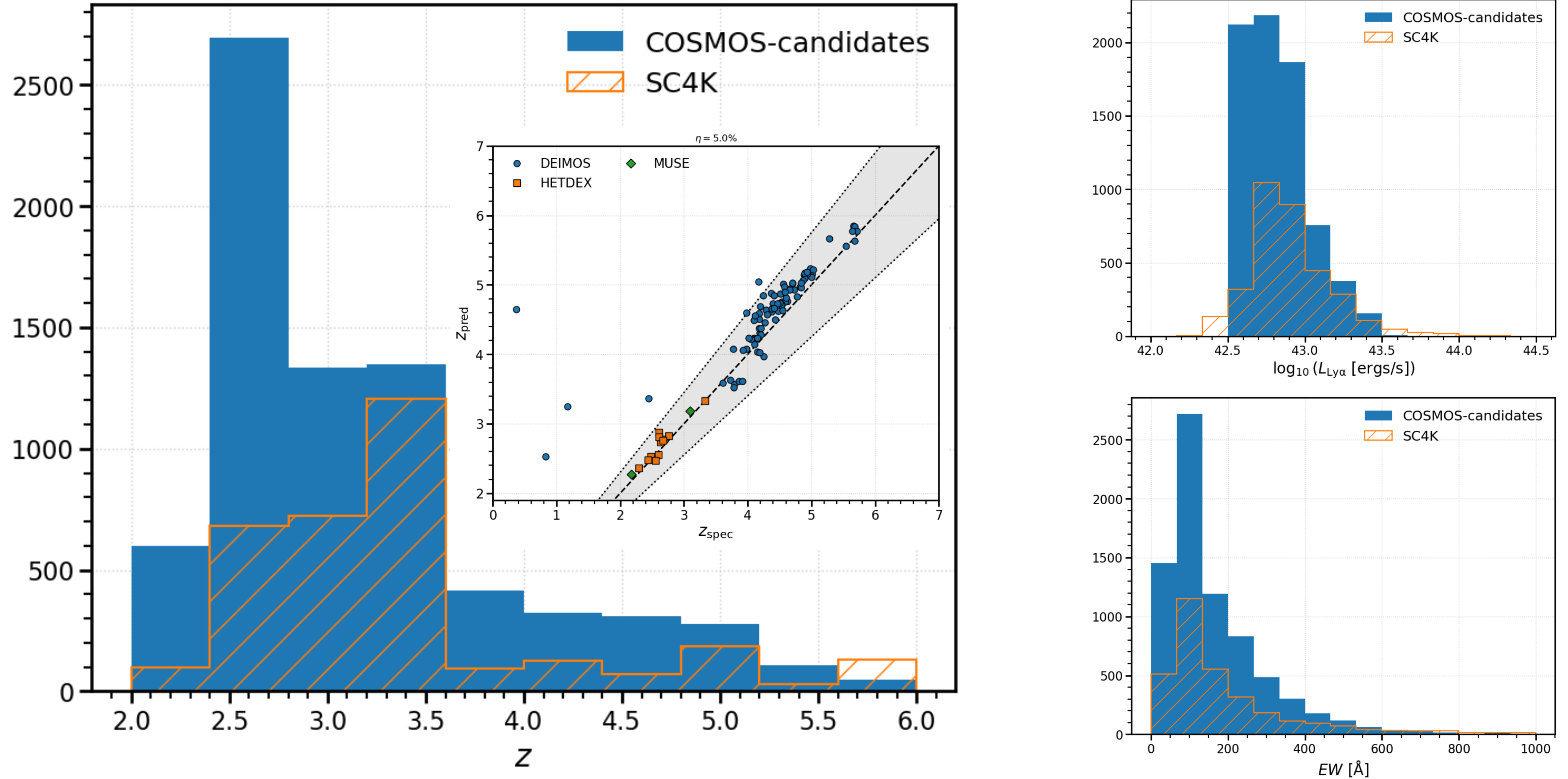
Results: HETDEX spectroscopic confirmations



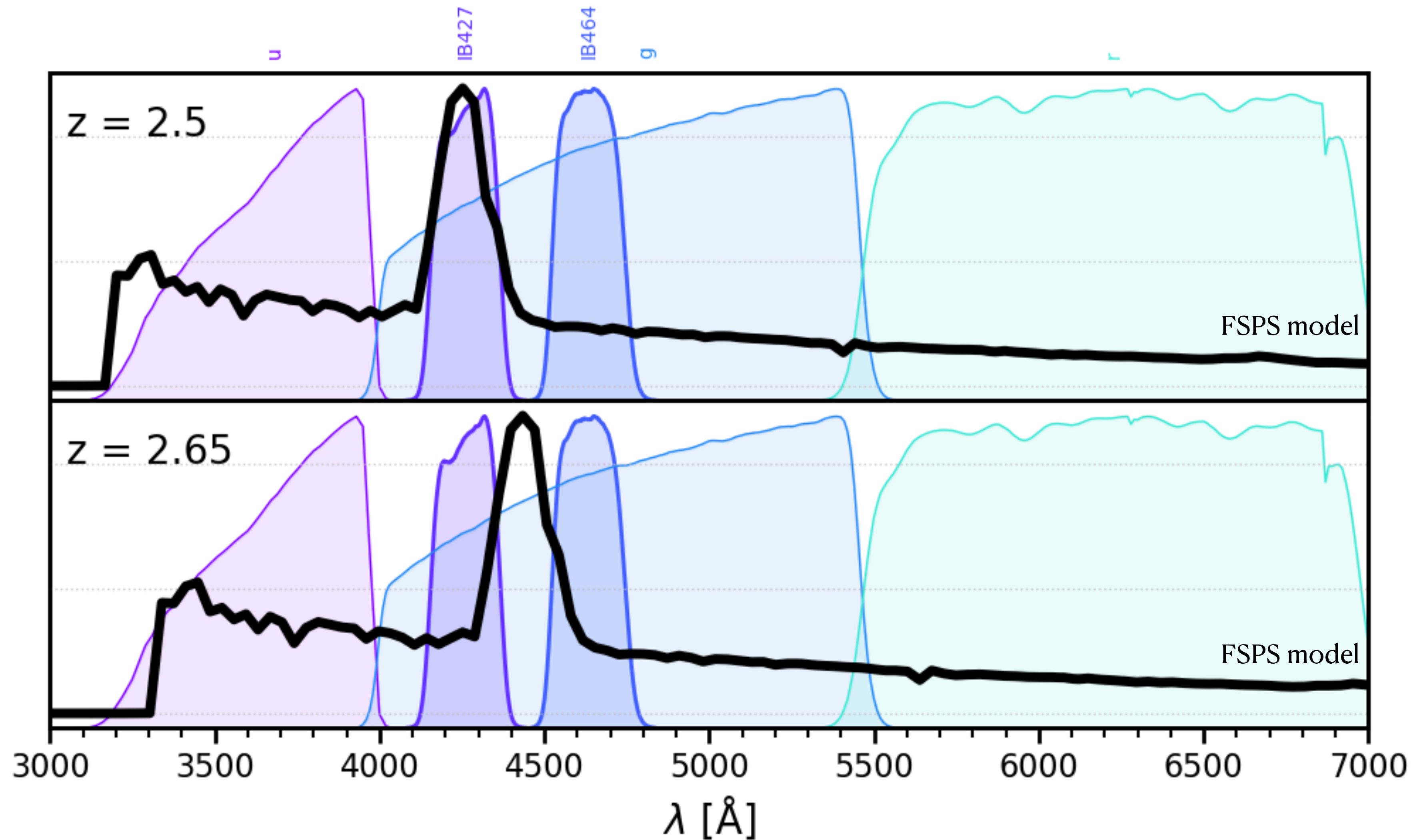
Cooper et al. 2023



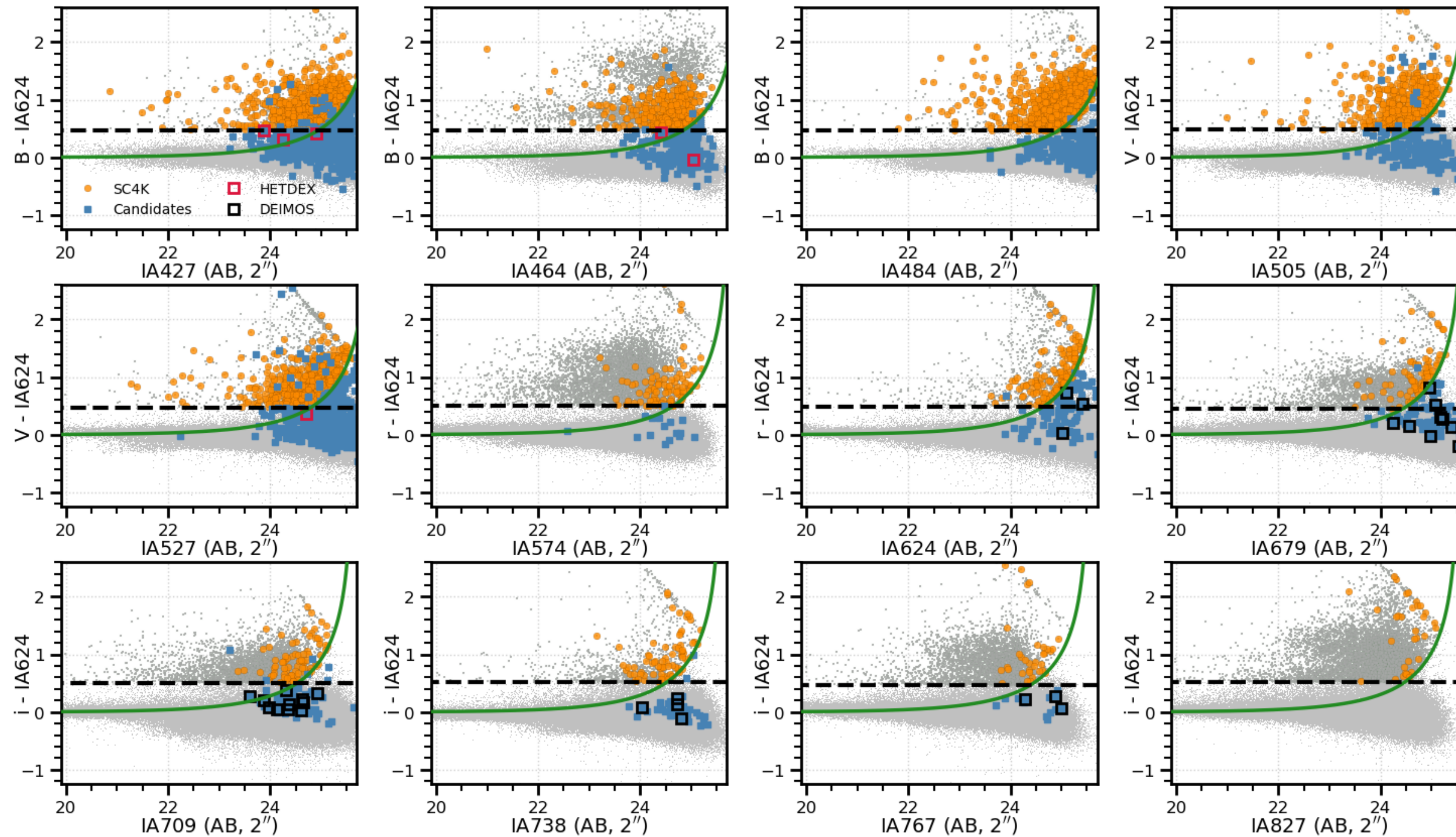
Results: predicting their z_s , $L_{Ly\alpha}$, and EWs



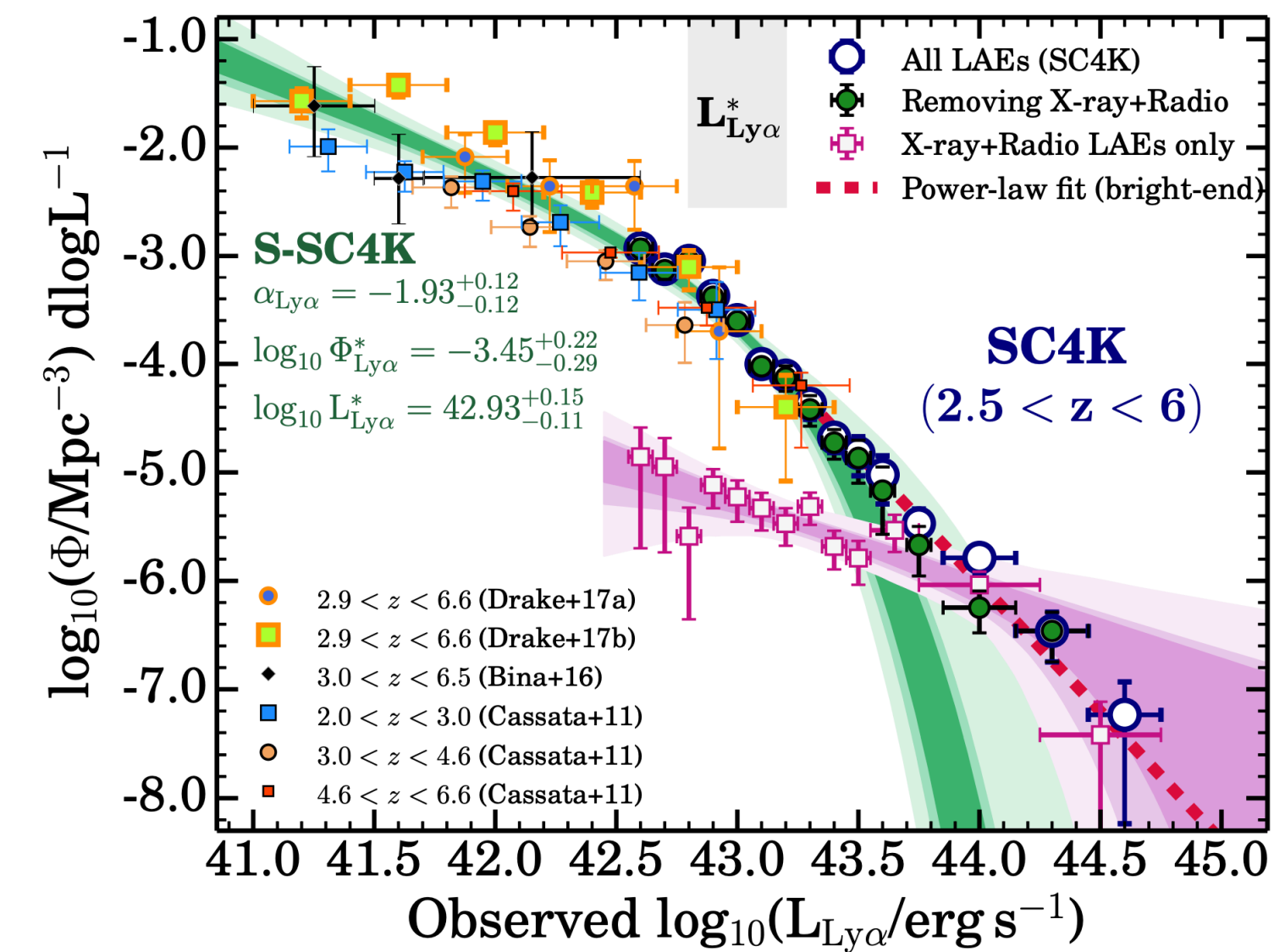
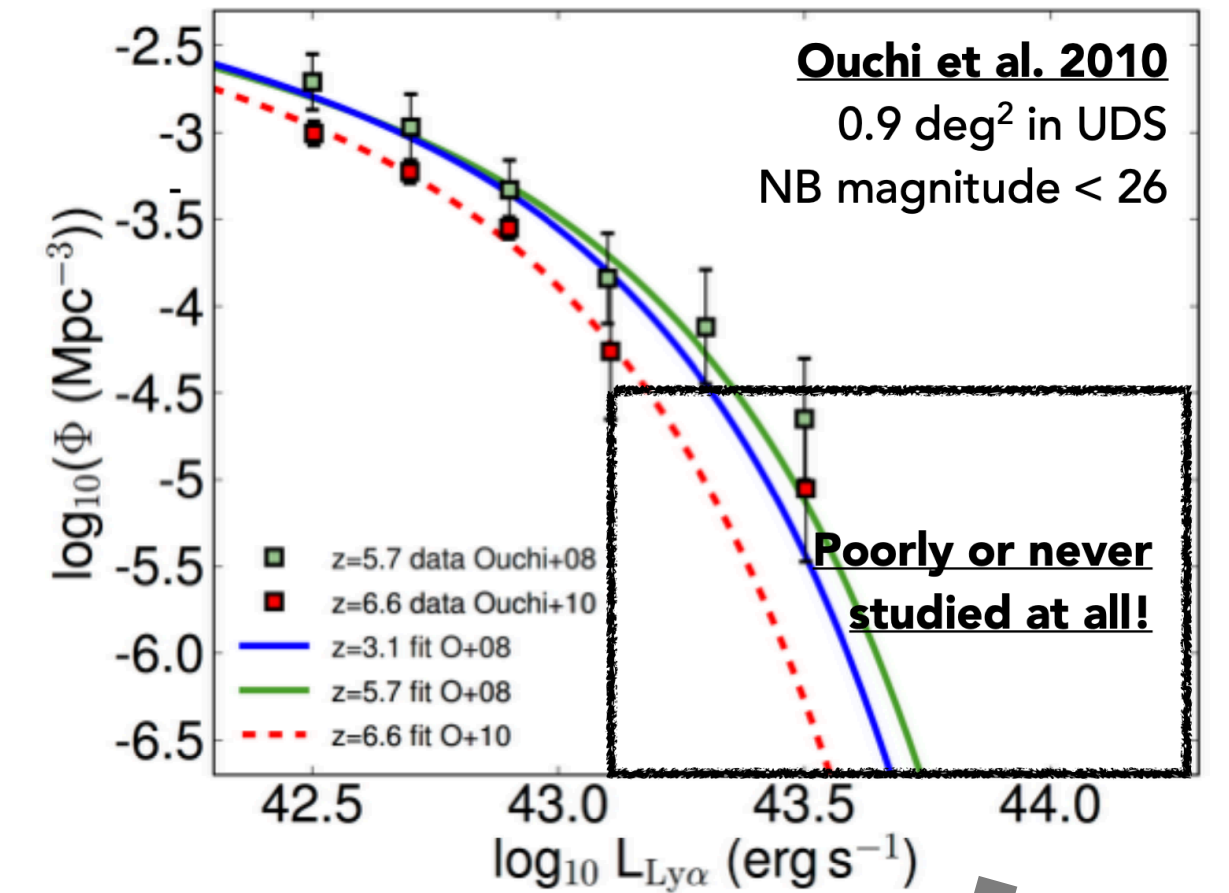
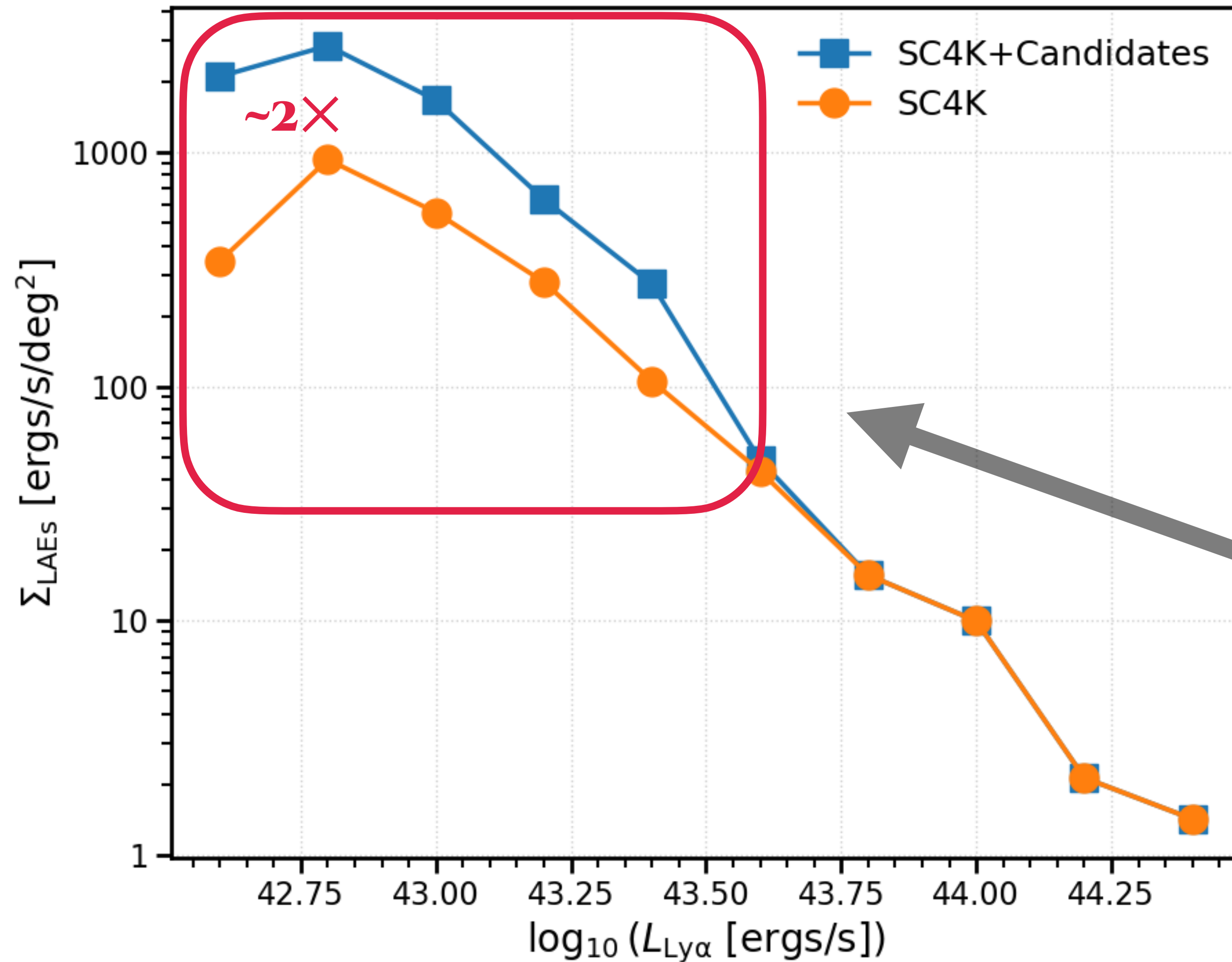
Discussion: why are we missing them?



Discussion: why are we missing them?



Discussion: what are the implications?



Take-home messages

- AI can be used to increase samples of LAEs.
- We are getting pure sample, but we are still far from having complete samples.
- We can also predict basic properties of LAEs which is useful for effective target selection.
- We still need to prove generalisation capabilities and discuss uncertainties.

Future directions

FLAEMING: Finding Lyman- α emitters through machine learning

T-FLAEMING: Transcoding the way of finding LAEs in the era of large surveys and AI

