

## The LSST AGN Data Challenge: Selection methods

DORĐE V. Savić,<sup>1,2</sup> ISIDORA JANKOV,<sup>3</sup> WEIXIANG YU,<sup>4</sup> VINCENZO PETRECCA,<sup>5,6</sup> MATTHEW J. TEMPLE,<sup>7,\*</sup>  
QINGLING NI,<sup>8</sup> RAPHAEL SHIRLEY,<sup>9,10</sup> ANDJELKA B. KOVAČEVIĆ,<sup>3,11</sup> MLADEN NIKOLIĆ,<sup>3</sup> DRAGANA ILIĆ,<sup>3,12</sup>  
LUKA Č. POPOVIĆ,<sup>2,3</sup> MAURIZIO PAOLILLO,<sup>5,6</sup> SWAYAMTRUPTA PANDA,<sup>13,14,†</sup> ALEKSANDRA ČIPRIJANOVIĆ,<sup>15</sup> AND  
GORDON T. RICHARDS<sup>4</sup>

<sup>1</sup>*Institut d’Astrophysique et de Géophysique, Université de Liège  
Allée du 6 Août 19c, 4000 Liège, Belgium*

<sup>2</sup>*Astronomical Observatory, Volgina 7, 11000 Belgrade, Serbia*

<sup>3</sup>*University of Belgrade - Faculty of Mathematics, Department of astronomy, Studentski trg 16 Belgrade, Serbia*

<sup>4</sup>*Drexel University, Department of Physics, 32 S. 32nd Street, Philadelphia, PA 19104, USA*

<sup>5</sup>*Department of Physics, University of Napoli “Federico II”, via Cinthia 9, 80126 Napoli, Italy*

<sup>6</sup>*INAF - Osservatorio Astronomico di Capodimonte, via Moiariello 16, 80131 Napoli, Italy*

<sup>7</sup>*Instituto de Estudios Astrofísicos, Universidad Diego Portales, Av. Ejército Libertador 441, Santiago, Chile*

<sup>8</sup>*Max-Planck-Institut für extraterrestrische Physik (MPE), Gießenbachstraße 1, D-85748 Garching bei München, Germany*

<sup>9</sup>*Astronomy Centre, Department of Physics and Astronomy, University of Southampton, Southampton SO17 1BJ, UK*

<sup>10</sup>*Institute of Astronomy, University of Cambridge, Madingley Road, Cambridge CB3 0HA, UK*

<sup>11</sup>*PIFI Research Fellow, Key Laboratory for Particle Astrophysics, Institute of High Energy Physics, Chinese Academy of Sciences, 19B  
Yuquan Road, 100049 Beijing, China*

<sup>12</sup>*Humboldt Research Fellow, Hamburger Sternwarte, Universität Hamburg, Gojenbergsweg 112, 21029 Hamburg, Germany*

<sup>13</sup>*Laboratório Nacional de Astrofísica - MCTI, R. dos Estados Unidos, 154 - Nações, Itajubá - MG, 37504-364, Brazil*

<sup>14</sup>*Center for Theoretical Physics, Polish Academy of Sciences, Al. Lotników 32/46, 02-668 Warsaw, Poland*

<sup>15</sup>*Fermi National Accelerator Laboratory, P.O. Box 500, Batavia, IL 60510, USA*

### ABSTRACT

Development of the Rubin Observatory Legacy Survey of Space and Time (LSST) includes a series of Data Challenges (DC) arranged by various LSST Scientific Collaborations (SC) that are taking place during the project’s preoperational phase. The AGN Science Collaboration Data Challenge (AGNSC-DC) is a partial prototype of the expected LSST AGN data, aimed at validating machine learning approaches for AGN selection and characterization in large surveys like LSST. The AGNSC-DC took part in 2021 focusing on accuracy, robustness, and scalability. The training and the blinded datasets were constructed to mimic the future LSST release catalogs using the data from the Sloan Digital Sky Survey Stripe 82 region and the XMM-Newton Large Scale Structure Survey region. Data features were divided into astrometry, photometry, color, morphology, redshift and class label with the addition of variability features and images. We present the results of four DC submitted solutions using both classical and machine learning methods. We systematically test the performance of supervised (support vector machine, random forest, extreme gradient boosting, artificial neural network, convolutional neural network) and unsupervised (deep embedding clustering) models when applied to the problem of classifying/clustering sources as stars, galaxies or AGNs. We obtained classification accuracy 97.5% for supervised and clustering accuracy 96.0% for unsupervised models and 95.0% with a classic approach for a blinded dataset. We find that variability features significantly improve the accuracy of the trained models and correlation analysis among different bands enables a fast and inexpensive first order selection of quasar candidates.

*Keywords:* galaxies: active; methods: statistical; surveys: catalogs; astrostatistics techniques: classification

\* Fondecyt fellow

† CNPq fellow

## 1. INTRODUCTION

A few percent of galaxies show enhanced emission from the nucleus that typically surpasses the stellar emission from the rest of the galaxy (e.g., Macuga et al. 2019); such sources are known as active galactic nuclei (AGNs). Emission from AGNs is produced by an accretion disk and ionized clouds surrounding a central super-massive black hole (Salpeter 1964; Zel'dovich & Novikov 1964; Antonucci 1993; Netzer 2015). AGNs emit across the whole electromagnetic spectrum (Padovani et al. 2017) and are readily observed at large distances due to their high luminosity, with potential to be used as probes of cosmology (Risaliti & Lusso 2019; Panda et al. 2019; Czerny et al. 2022). AGNs have profound effects on the life and evolution of their entire host galaxy (Ferrarese & Merritt 2000; Gebhardt et al. 2000; Kormendy & Ho 2013). Outflows and jets interact with the local environment and release a large amount of energy capable of driving away the nearby gas, hence terminating star formation (Fabian 2012). Moreover, AGNs also have an impact on the surrounding hot intergalactic medium and play an active role in the evolution of the host galaxy clusters (Eckert et al. 2021). Therefore, each successful detection and observation of AGNs and measuring their physical properties is crucial for many areas of modern astrophysics and cosmology.

Ongoing and forthcoming large-scale photometric surveys (e.g. Zwicky Transient Facility - Bellm 2014; Pan-STARRS - Chambers et al. 2016; Gaia - Gaia Collaboration et al. 2016) will produce catalogs for a vast number of sources which brings the astronomy to the new era of “big data”. The Vera C. Rubin Observatory Legacy Survey of Space and Time (LSST) is designed to address the main challenges for probing dark energy and dark matter, solar system exploration, exploring the transient optical sky and mapping the Milky Way (LSST Science Collaboration et al. 2017; Ivezić et al. 2019). With a state-of-the-art 3.2 Gigapixel flat-focal array camera mounted on an 8.4m telescope, LSST will cover the whole observable sky every  $\sim 4$  nights in the optical/near-infrared *ugrizy*-bands. The expected data volume of LSST is  $\sim 300$  PB of raw data and  $\sim 4 \times 10^{10}$  objects after 10 years of planned survey (Ivezić et al. 2019). Every night, LSST will monitor tens of millions of AGNs over  $\sim 18\,000$  deg<sup>2</sup> area (Luo et al. 2017; De Cicco et al. 2021). Although the actual number of AGNs that will be detected is a subject of the optimal observing strategy (Bianco et al. 2022), LSST will produce an AGN sample that supersedes the largest current AGN samples by more than an order of magnitude. This present work is a preparatory step towards producing a high-purity AGN sample with LSST.

To identify AGNs within LSST, the main challenge is separating AGNs from normal galaxies and stars. Construction of LSST’s AGN census will build upon a considerable volume of past work, making use of colors, proper motion, variability and image morphologies. The idea of performing a multi-faceted quasar selection (ie., combining information from multiple observables) has long been proposed (Koo et al. 1986). However, the quality, quantity, and type of data of LSST will allow for a more complete AGN selection and thus these approaches should be considered from scratch.

Color selection has been widely used as the gold standard for identification of unobscured AGNs since their discovery (Koo & Kron 1982; Warren et al. 1991; Richards et al. 2002), but we expect AGN colors to change as a function of luminosity as we probe fainter towards LSST-like depths (e.g., Temple et al. 2021). While an application of modern statistical techniques to color data could be used to select AGNs, we expect the addition of multi-parameter data to result in a purer and more complete AGN selection function.

As some AGNs and stars have similar colors, the fact that AGNs lack proper motions (unlike Galactic stars) has long been used as a discriminant (Sandage & Luyten 1967; Kron & Chiu 1981). LSST’s use of astrometric data will be no different in that regard. One way that LSST will be unique is in its ability to take advantage of differential chromatic refraction (DCR) of AGN (Kaczmarczik et al. 2009; Yu et al. 2020) which makes use of the astrometric offset of an emission-line object from that expected (in the astrometric solution) for a power-law source—to break degeneracies in photometric redshifts of luminous AGNs (henceforth quasars or QSOs).

Selection of AGNs via time-series due to their variability is another well known method (Bonoli et al. 1979; Trevese et al. 1989; Butler & Bloom 2011; De Cicco et al. 2019; Poulain et al. 2020). As quasars display higher fractional variability in their brightness than the average star and with different characteristics than the typical variable star, variability will be a cornerstone of AGN classification for LSST (Suberlak et al. 2021). For luminous quasars, it has been shown that variability combined with colors works better for selection than variability alone (Peters et al. 2015). Lower-luminosity AGNs are expected to have the most variable nuclei, however increased contamination from the host galaxy could compromise variability-selection methods if insufficient care is taken. The implementation of Difference Image Analysis (DIA) in the LSST reduction pipelines will completely revolutionize the detection of AGNs through variability by removing the contribution from the host galaxy (Zebrun et al. 2001; Bramich 2008;

Kozłowski et al. 2010a, 2016). Finally, the addition of high quality resolution images is expected to considerably increase the performance of the selection methods (Doorenbos et al. 2022).

Several data challenges (DCs) have been created in the past to facilitate the preparation of LSST, by other groups and for science use cases other than the study of AGN (e.g., Sánchez et al. 2020; Hložek et al. 2020). In 2021, the LSST AGN Science Collaboration (AGNSC<sup>1</sup>) organized a DC to get more people involved in the work needed for the AGN science with the upcoming LSST data. The main goal of the LSST AGN DC was to address the problem of AGN selection.

Unlike the previous data challenges, which relied on simulated datasets, this AGN DC utilizes real observational data. Major tasks also include establishing public training/test sets that will be used as a benchmark to test different machine learning (ML) algorithms. There were five proposed solutions submitted to the DC: one using a classical approach; and four applying ML-based AGN selection. We will present the solution using classical approach and three ML-based solutions, while the one remaining ML-based solution is addressed by Doorenbos et al. (2022).

The paper is organized as follows: in Section 2 we address the data retrieval and the construction of training and blinded datasets. We elaborate on applied ML methods in Section 3. The results obtained from the various methodologies employed in this DC are presented in Section 4. We discuss further issues relevant to our work and summarize our findings in Section 5.

## 2. DATASETS CONSTRUCTION

LSST will deliver three levels of data products and services: prompt data products that are computed and released within 24 hours of observation, data release data products that are computed during annual processing campaigns, and user generated data products. The Data Products Definitions Document (DPDD<sup>2</sup>) is the ultimate reference for descriptions of the planned LSST data products and pipelines (see also Ivezić et al. 2019).

The input data for the AGNSC-DC were modified such that the column names and units used for different measurements (e.g., flux) comply with the DPDD standards for data release catalogs (DPDD, Section 4.3). We refer to distinct astrophysical bodies that emit light detected as “Objects” and individual instances (detection) of those objects as “Sources”. Observations from a specific point in time will appear in the `Source` tables in

the data releases, while “co-added” (averaging/summing over time) information will appear in the `Object` tables. So-called “light curves” (brightness as a function of time) will appear in the `ForcedSource` tables (with summary statistics in the `Object` tables). “Simulated” training data will attempt to heel to this data structure as closely as possible.

The datasets released in this data challenge are pulled from different sources (public archives) and put together to mimic the architecture of future LSST data release catalogs as much as possible, but without taking into account the expected number of objects that will appear in certain regions of LSST sky. Details on how the tables are constructed can be found under the ‘docs’ folder in the main github repository (Yu & Richards 2021). Here, we provide a brief overview of the datasets.

The astronomical objects included in the release dataset are drawn from three main sources: spectroscopically identified objects in an extended Sloan Digital Sky Survey (SDSS<sup>3</sup>; York et al. 2000a) Stripe 82 area with the spectroscopy collected from the 16th data release of SDSS (DR16; Ahumada et al. 2020), X-ray detected and classified objects in the the XMM-LSS<sup>4</sup> (Pierre et al. 2007) region, and unidentified variable objects in the original SDSS Stripe 82 area<sup>5</sup> (Ivezić et al. 2007). Fig. 1 illustrates the source survey region footprint on the LSST sky for the baseline observing strategy. The XMM-LSS area is encompassed by one of the LSST deep drilling fields.

The total number of objects (both combined) in the `Object` table is  $\sim 440\,000$ , after removing  $\sim 5000$  duplicates found in more than one sources from above. The total number of epochs in the `ForcedSource` table is  $\sim 5$  million. The total number of features (parameters) in the object table is 374. Features are divided into main categories with number of features in each category indicated in the parentheses:

1. **Astrometry(5)**: ra, dec, proper motion and parallax.
2. **Photometry(48)**: point and extended source photometry, in both AB magnitudes and fluxes (nJy).
3. **Color(10)**: derived from the flux ratio between different photometric bands.

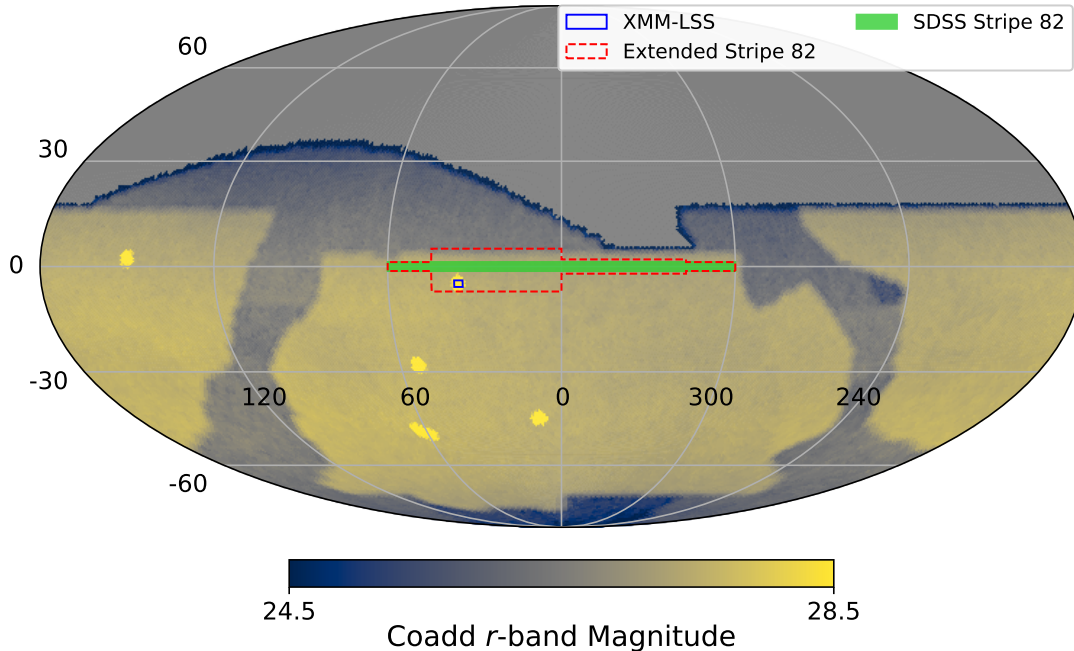
<sup>3</sup> <https://www.sdss.org/>

<sup>4</sup> <https://www.cosmos.esa.int/web/xmm-newton>

<sup>5</sup> <http://faculty.washington.edu/ivezic/sdss/catalogs/S82variables.html>

<sup>1</sup> <https://agn.science.lsst.org/>

<sup>2</sup> <https://docushare.lsst.org/docushare/dsweb/Get/LSE-163>



**Figure 1.** SDSS Stripe 82 (filled green), extended Stripe 82 (dashed red) and XMM-LSS (solid blue) areas projected on the LSST observable sky. Color map indicates the final depth of the coadds in the  $r$ -band. XMM-LSS region coincides completely with one of the LSST deep drilling fields (bright-yellow regions). The map was generated using the LSST simulation project OpSim (Delgado et al. 2014).

4. **Morphology(6)**: a real-value quantity between 0 and 1. Values closer to 1 for extended sources while values closer to 0 indicate point-like sources.
5. **Light Curve Features(302)**: extracted on the SDSS light curves if available.
6. **Redshift(2)**: both spectroscopic and photometric, wherever available.
7. **Class Labels(1)**: Star/Gal/Qso (Agn, high-ZQso), wherever available.

Distributions of a subset of 30 features are shown in Fig. 2. It is notable that stars, galaxies and quasars overlap in the feature space, but they may be separated by combining a selection of features containing the most information about interclass difference. For example, galaxies occupy different parts for morphology features when compared to stars and quasars, since the latter two are observed as point sources (Fig. 2, 4th row middle and right panels); and similarly for quasars when compared to stars and galaxies for variability features (Fig. 2, bottom row from left to right).

Astrometry measurements were obtained by matching the main catalogs (SDSS Stripe 82 and XMM-LSS)

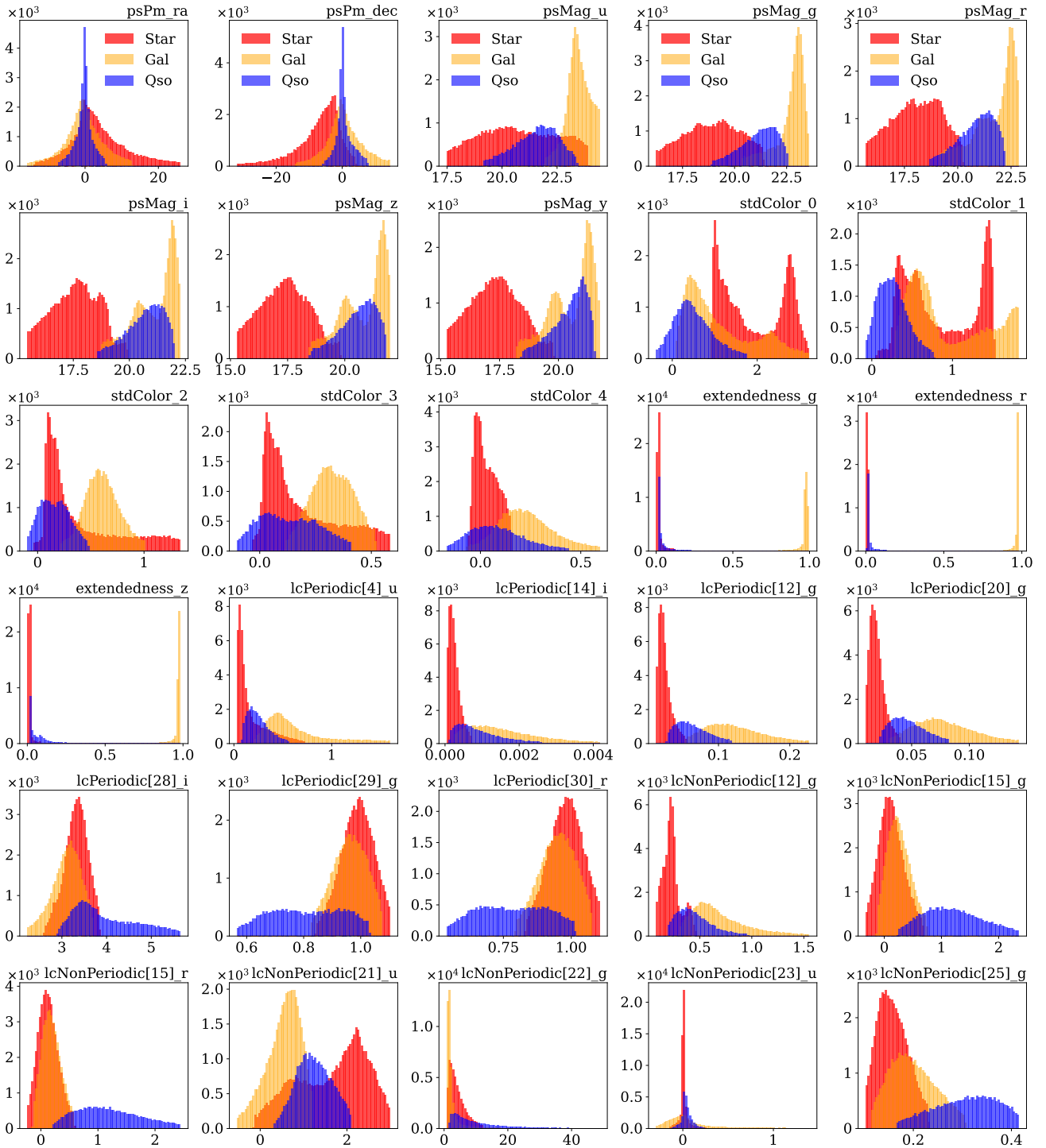
with Gaia EDR3<sup>6</sup> (Early Data Release 3 Gaia Collaboration et al. 2016, 2021) and the NOIRLab Source Catalog (NSC) data release 2 (DR2<sup>7</sup>; Nidever et al. 2021). Sources in NSC are extracted from reprocessed public images drawn from the NOIRLab Astro Data Archive<sup>8</sup>. The astrometry for NSC DR2 is calibrated using Gaia DR2 (Gaia Collaboration et al. 2018); its proper motion measurement achieves an accuracy of  $0.2 \text{ mas yr}^{-1}$  and a precision of  $2.5 \text{ mas yr}^{-1}$  relative to Gaia DR2 (Nidever et al. 2021). For objects with astrometry measurements found in both catalogs, we used the values from Gaia.

The photometry were assembled following a mix-and-match approach. In the extended Stripe 82 region, we cross-matched our sources against the Dark Energy Survey (DES; Dark Energy Survey Collaboration et al. 2016) Data Release 2 (DR2; DES Collaboration et al. 2021) photometry catalog, the SDSS Stripe 82 coadded photometry catalog (Annis et al. 2014), and the SDSS DR16 single-epoch photometry catalog (Ahumada et al. 2020). DES provides photometry in  $grizY$  bands and

<sup>6</sup> <https://www.cosmos.esa.int/web/gaia/early-data-release-3>

<sup>7</sup> <https://datalab.noirlab.edu/nscdr2/index.php>

<sup>8</sup> <https://astroarchive.noirlab.edu>



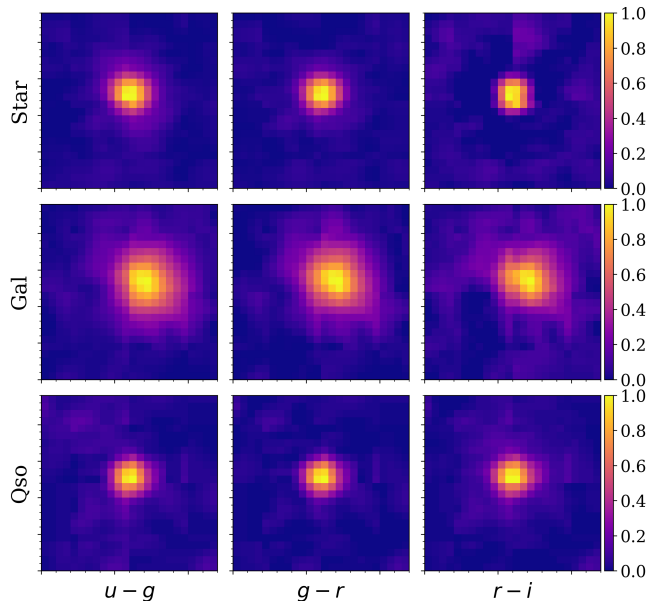
**Figure 2.** Distributions of 30 features drawn from astrometry (proper motion `psPm_ra` and `psPm_dec`); photometry (point-source magnitudes `psMag_u/g/r/i/z/y` and colors `stdColor_0/1/2/3/4`); morphology (extendedness `extendedness_g/r/z`); LC features (`lcPeriodic[12/20/29]_g`, `lcPeriodic[14/28]_i`, `lcPeriodic[30]_r`, `lcNonPeriodic[12/15/22/25]_g`, `lcNonPeriodic[4/21/23]_u`, `lcNonPeriodic[15]_r`). Color-coded per class: star (red), galaxy (orange), quasar (blue).



SDSS provides photometry in *ugriz* bands. When a source is matched with photometry in the same band from catalogs, we choose the photometry following the precedence of: DES DR2 > SDSS Stripe 82 coadd > SDSS DR16. In the XMM-LSS region, the *griz* photometry were collected from the HSC-VISTA joint catalog (see Section 2.1.2) and the *u*-band photometry comes from the Canada-France-Hawaii Telescope Legacy Survey (CFHTLS) catalog (Gwyn 2012).

The broad-band colors were first derived using the flux ratios of two adjacent bands and then converted into magnitudes. The error on the color was computed following the standard uncertainty propagation method. For the photometry taken from SDSS, the morphology (i.e., *extendedness* in LSST’s nomenclature) is defined as  $1 - \text{probPSF}$ . For DES photometry, we define *extendedness* as  $1 - \text{class\_star}$ . Both *probPSF* and *class\\_star* describe how close the photometry is to a true point source. For the photometry taken from the HSC-VISTA catalog, the LSST pipeline directly outputs the *extendedness* column.

All objects, except for  $\sim 130$  of them, in the *Object* table have pre-generated image thumbnails/cutouts of size 64x64 pixel from SDSS DR16. The objects without image cutouts were not removed from the main sample. An example of image cutouts for a randomly selected star, galaxy and quasar is shown in Fig. 3.



**Figure 3.** Image cutouts for a randomly selected star (top), galaxy (middle) and quasar (bottom). From left to right are colors  $u - g$ ,  $g - r$  and  $r - i$ . Images are normalized to unity for each color.

We list the percentage per class label and the assigning method:

- **Star**(24.5%): Spectroscopically confirmed stars (both variable and non-variable)
- **Gal**(55.9%): Spectroscopically confirmed galaxies
- **Qso**(19.0%): Spectroscopically confirmed AGNs or quasars
- **highZQso**(0.3%): A separate catalog of high redshift ( $z > 4.5$ ) quasars
- **AgN**(0.3%): X-ray classified AGNs in the XMM-LSS region and spectroscopically classified galaxies having emission properties consistent with being a Seyfert or LINER in the extended Stripe 82 region.

In addition to the described datasets that are publicly available, a distinct *blinded* dataset, that is  $\sim 10\%$  of the training dataset, was constructed and set aside in order to evaluate the performance of the ML and non-ML methods submitted by the participants of the DC. The total object count for a blinded dataset is  $\sim 45\,100$ , among which,  $\sim 1000$  come from the XMM-LSS region,  $\sim 44\,000$  come from the Stripe 82 region, and 100 come from the separate high- $z$  quasar catalog. Class labels and spectroscopic redshift were removed for those objects in the blinded dataset. Around 21 000 objects have pre-computed variability features.

We stress important caveats:

1. About  $\sim 13,000$  objects have no labels, and they come from the SDSS Stripe 82 unidentified variables source catalog (Ivezić et al. 2007). About  $\sim 2500$  objects from the XMM-LSS are classified using X-ray data and infrared photometry (Chen et al. 2018), which also have no spectroscopic redshift available.
2. Approximately  $\sim 1\%$  of the objects do not have optical counterpart i.e., they are bright in X-ray, but are too faint to be detected in the optical band.
3. CARMA(1,0) and CARMA(2,1) fits are presented as is. Potential bad fits (e.g., perhaps due to limited temporal sampling and/or poor S/N of the photometry) were not removed<sup>9</sup>

In the following two subsections, we describe how multi-wavelength datasets; high-redshift quasar catalog and how the light curve features are computed.

<sup>9</sup> A robust goodness-of-fit metric for CARMA models is not available. The definition for a bad fit can also change given the problem at hand.

Class label	Agn/Qso	highZQso	Gal	Star
Stripe 82	73 000	90	213 000	106 000
GALEX	36 000	10	45 000	38 000
UKIDSS	36 000	30	87 000	92 000
Spitzer	12 000	20	27 000	30 000
Herschel	2500	10	39 000	12 000
FIRST	2000	50	43 000	250 000

**Table 1.** The distribution of class labels per catalog. From top to bottom: Stripe 82; GALEX; UKIDSS; Spitzer; Herschel and FIRST.

Class label	Agn/Qso	Gal	Star
XMM-Newton	4000	N/A	N/A
GALEX	400	300	100
VISTA/VIDEO	3500	3500	370
Spitzer	3500	3500	350
Herschel	1500	1500	50

**Table 2.** The distribution of class labels per catalog for objects found in XMM-LSS dataset.

### 2.1. Multi-wavelength data

Multi-wavelength tables are provided along with the Object table. The multi-wavelength data are obtained by performing positional cross-match between our source positions and other catalogs. A summary of class label distribution per observing mission is listed in Table 1 and 2.

#### 2.1.1. XMM-LSS

In the X-ray band of the XMM-LSS field, a number of XMM-Newton surveys of different sensitivities have been collected (e.g. Fig. 3 of Brandt & Alexander 2015, also Table 2 of Chen et al. 2018). The X-ray source catalog from Chen et al. (2018) is adopted for our dataset, which makes use of XMM-Newton observations taken from 2000 to 2017 in the XMM-LSS<sup>10</sup> field, including the 1.3 Ms new endeavour from the XMM-SERVS survey (Chen et al. 2018; Ni et al. 2021). With the XMM-SERVS survey, a flux limit of  $6.5 \times 10^{-15}$  erg cm<sup>-2</sup> s<sup>-1</sup> over 90% of the XMM-LSS area is achieved in the 0.5–10 keV band. We only include X-ray sources that are classified as AGNs in the dataset (Chen et al. 2018, Section 6 of) for the source classification details). The X-ray AGNs are matched to other optical/IR catalogs with likelihood-ratio matching methods as described in Section 4 of Chen et al. (2018).

#### 2.1.2. HSC and VISTA joint catalogue

We have provided a jointly processed optical and near-infrared dataset from the HSC<sup>11</sup> (Hyper Suprime-Cam) Public Data Release 2 (Aihara et al. 2019) deep and ultradeep regions, and the VISTA<sup>12</sup> (Visible and Infrared Survey Telescope for Astronomy) VIDEO (Jarvis et al. 2012) surveys. The dataset was produced using the LSST Science Pipelines as described in Bosch et al. (2018). An object detected in any one of the ten bands across these two surveys is measured in every band ensuring that each object will have a measurement in each band. This dataset is a prototype developed in preparation for the upcoming LSST data.

#### 2.1.3. UKIDSS

In the near-infrared *YJHK* bands, we include 2 arcsec diameter aperture magnitudes (*AperMag3*) where available from the UKIDSS DR11plus (Hewett et al. 2006; Lawrence et al. 2007; Hambly et al. 2008; Hodgkin et al. 2009). SDSS Stripe 82 is partially covered by the UKIDSS Large Area Survey (LAS) in the *YJHK* bands to approximate depths of 20.2, 19.6, 18.8, and 18.2 respectively. The UKIDSS-LAS covers the original Stripe 82 footprint (shown in green in Fig. 1) but not the full extended area used in this challenge. We remove duplicate detections from overlapping tiles using *PriOrSec*, and remove noise and saturated detections using *mergedClass*.

#### 2.1.4. Herschel forced photometry

Far infrared measurements from the *Herschel*<sup>13</sup> Space Observatory come from HELP: The *Herschel* Extragalactic Legacy Project (Shirley et al. 2021). This dataset was produced by taking a prior list of Spitzer IRAC detections and providing full Bayesian probability posterior distributions on the object fluxes to account for blending in the low resolution far infrared maps at 250, 350, and 500 microns.

### 2.2. High-redshift quasars

The catalog of high-redshift known quasars is constructed by collecting all quasars at  $z \geq 4.5$  known before October 2020. These quasars are mainly selected using the optical/near-infrared colors, based on the wide-field optical and infrared photometric surveys, e.g., SDSS (York et al. 2000b), the Pan-STARRS1 survey (PS1, Chambers et al. 2016), the DESI Legacy Imaging Surveys (Dey et al. 2019), the Hyper Suprime-Cam

<sup>10</sup> <https://personal.psu.edu/wnb3/xmmservs/xmmservs.html>

<sup>11</sup> <https://www.naoj.org/Projects/HSC/>

<sup>12</sup> <https://www.eso.org/public/teles-instr/paranal-observatory/surveytelescopes/vista/>

<sup>13</sup> <https://www.herschel.caltech.edu/>

Subaru Strategic Program survey (Aihara et al. 2018), the UKIRT Hemisphere Survey (Dye et al. 2018), the UKIDSS-LAS (Lawrence et al. 2007), the VISTA Hemisphere Survey (McMahon et al. 2013), and the Wide-field Infrared Survey Explorer (*WISE*, Wright et al. 2010).

About half of the  $z < 6$  quasars are from SDSS quasar catalogs and the rest are from several major quasar surveys (e.g., McGreer et al. 2013; Bañados et al. 2016; Wang et al. 2016; Yang et al. 2019a). The  $z > 6$  quasars were mostly collected from quasar surveys like the SDSS high-redshift quasar survey (e.g., Fan et al. 2006; Jiang et al. 2016), the Canada–France High- $z$  Quasar Survey (e.g., Willott et al. 2010), the PS1 distant quasar survey (e.g., Bañados et al. 2016; Venemans et al. 2015; Mazzucchelli et al. 2017), the Subaru High- $z$  Exploration of Low-Luminosity Quasars project (e.g., Matsuoka et al. 2018), the DES quasar survey (e.g., Reed et al. 2017), and the reionization-era quasar survey (e.g., Wang et al. 2019; Yang et al. 2019b). Quasars included in this catalog are all identified through spectroscopic observations. Their redshifts are mainly from the quasar broad emission lines in the rest-frame UV (e.g., Ly $\alpha$ , Si IV, C IV, and Mg II), which result in a redshift uncertainty up to  $\sim 0.05$ . A small number of them have [C II]-based redshifts, with a typical uncertainty of  $< 0.001$ .

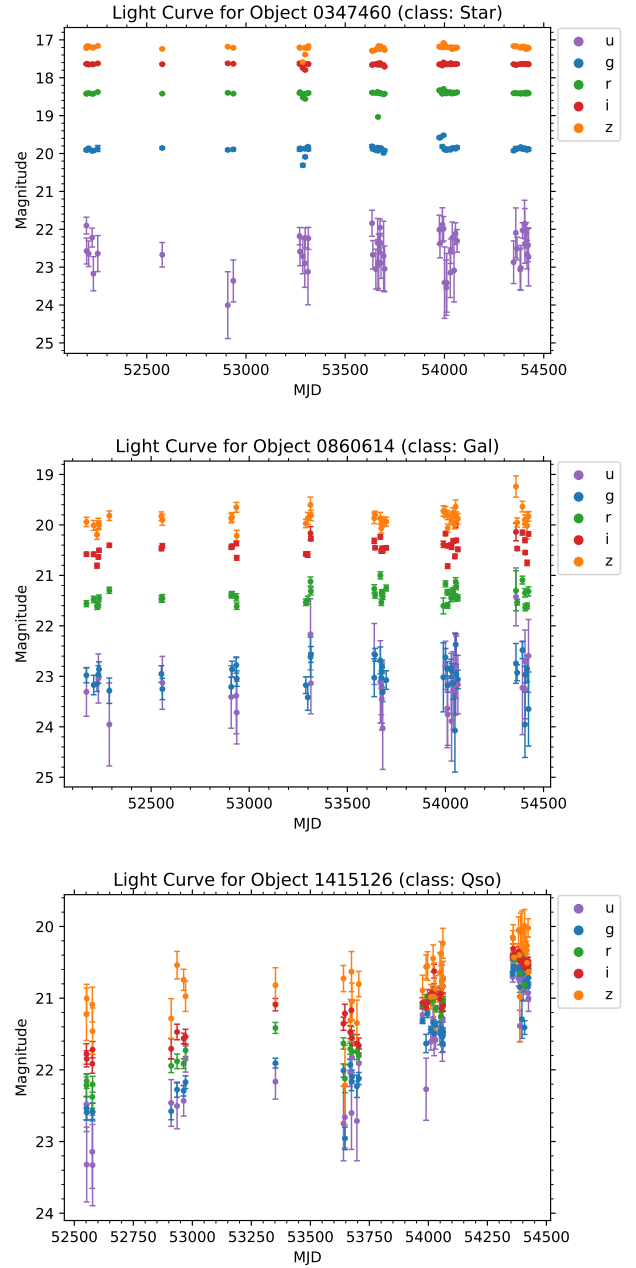
### 2.3. Light curve features

The light curve (LC) features are computed for the sources that have corresponding time-domain data from SDSS. These features populate the `lcPeriodic` and `lcNonPeriodic` columns. The majority of the features computed are described in Richards et al. (2011), hereafter R11. Some additional features computed include those introduced by the Feature Analysis for Time Series (FATS) project<sup>14</sup> and best-fit CARMA(1,0) (continuous-time auto-regressive moving average model, otherwise known as a damped random walk; DRW) and CARMA(2,1) (otherwise known as a damped harmonic oscillator; DHO) parameters (Kelly et al. 2009; Kasliwal et al. 2017; Moreno et al. 2019; Yu et al. 2022) obtained using EzTAO<sup>15</sup> (Yu & Richards 2022). Both the R11 and FATS features are computed using the CESIUM<sup>16</sup> software package. An example of LC for a random star, galaxy and QSO for the SDSS *ugriz*-bands are shown in Fig. 4.

<sup>14</sup> <https://isadoranun.github.io/tsfeat/FeaturesDocumentation.html>

<sup>15</sup> <https://github.com/ywx649999311/EzTao>

<sup>16</sup> <https://cesium-ml.org/>

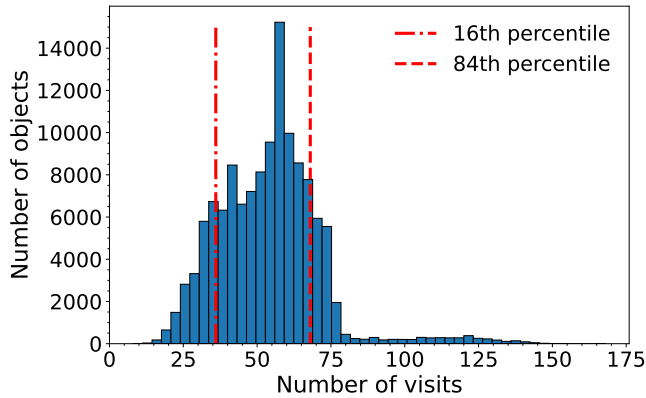


**Figure 4.** Example LCs for the SDSS *ugriz*-bands: *u*-band (purple), *g*-band (blue), *r*-band (green), *i*-band (red) and *z*-band (orange) for a star (top panel), galaxy (middle panel) and QSO (bottom panel). Time units are given in modified Julian date (MJD).

We note that a  $5\text{-}\sigma$  clipping (in magnitude) was applied before the variability metrics were computed, this decision is mostly driven by the ‘spurious dimming’ of SDSS light curves as discussed by Schmidt et al. (2010). The total count of time domain objects in Stripe 82 is  $\sim 210\,000$ . Each SDSS filter, for which LC features have been computed, has the same number of visits and the



same cadence. These numbers vary from object to object. The distribution of the number of visits per object is shown in Fig. 5. The bulk of objects have the number of visits between 30 and 70. For comparison, LSST sources will have a larger number of observations (Bianco et al. 2022; Kovacevic et al. 2022; Raiteri et al. 2022; Pozo Nuñez et al. 2023; Czerny et al. 2023).



**Figure 5.** Distribution of the number of visits in the `ForcedSource` table for objects with known class label. Vertical red lines mark the 16th and 84th percentile values.

Table 3 summarizes the LC features and notation. A full description of the LC features is also publicly available<sup>17</sup>.

### 3. MACHINE LEARNING METHODS

The classification of sources from wide-field surveys is one of the most fundamental problems in astronomy. Efficient classification will be difficult if using classical techniques for manual inspection and it is expected that ML applications will be helpful in automating the process. When trained on big astronomical data, ML methods tend to outperform traditional methods based on explicit programming (e.g., Banerji et al. 2010; Lochner et al. 2016; Baron 2019).

Broadly speaking, ML methods can be divided into supervised and unsupervised methods (Berry et al. 2019). Supervised Learning is a ML paradigm that relies on labeled data for acquiring the input-output relationship information for classification and regression problems. The supervised ML methods we used are: support vector machine (SVM, Cortes & Vapnik 1995), random forest (RF, Ho 1995), extreme gradient boosting (XGB, Chen & Guestrin 2016) and artificial neural networks (ANN, Cybenko 1989). These methods have been widely applied to numerous classification problems in astronomy

(e.g. Baron 2019) and for AGN classification (Carballo et al. 2008; Cavuoti et al. 2014; Doert & Errando 2014; Chen et al. 2021; De Cicco et al. 2021; Poliszczuk et al. 2021; Chang et al. 2021). For more details on each of supervised ML method we used, we refer to Berry et al. (2019).

Unsupervised methods discover hidden patterns in the data without the need for human intervention and are mostly used for clustering and dimensionality reduction. Commonly used are:  $k$ -means (Lloyd 1982), Gaussian mixture of models (Reynolds 2009) for clustering; and principal component analysis (PCA, Jolliffe 1986), autoencoders (Kramer 1991), t-distributed stochastic neighbor embedding (t-SNE, van der Maaten & Hinton 2008) for dimensionality reduction.

Deep learning (DL) is another broader family of ML methods that relies exclusively on the use of ANNs with a large number of hidden layers, hence *deep*. The abundance of imaging data in astronomy has naturally led to the application of deep convolutional neural networks (CNNs, Lecun et al. 2015). One advantage of DL is that it allows for simultaneous training of models with multiple inputs (e.g., catalog features + images) and for multiple outputs (e.g., regression and classification, Chollet et al. 2015). An important property of DL is the application of transfer learning (Tan et al. 2018) i.e., a model trained for one task is re-purposed on a second related task (Doorenbos et al. 2022). For the application of DL in astronomy, we refer to a review by Smith & Geach (2022).

Many ML methods are developed as a combination of the above mentioned methods. Therefore, it is common to train a dimensionality reduction model and then train supervised or unsupervised models using the latent features. We used one unsupervised model built in this manner: deep embedding clustering (DEC; Xie et al. 2015; Guo et al. 2017), which consists of an ANN based on autoencoder for dimensionality reduction and the preservation of the latent space, followed by clustering using latent features. Additionally, latent features are used to for visualizing complex multidimensional data space in 2 or 3 dimensions (Clarke et al. 2020; Jankov et al. 2021), which is often the very beginning of the ML experiment setup.

### 4. RESULTS

In this section, we summarize four out of the five “solutions” submitted to the DC (since one using CNNs and transfer learning has been already been presented by Doorenbos et al. 2022). One submission used a non-ML method extending the traditional approach based on color-color diagrams by adding magnitude standard

<sup>17</sup> [https://github.com/RichardsGroup/AGN\\_DataChallenge/blob/main/docs/04\\_LC\\_features.ipynb](https://github.com/RichardsGroup/AGN_DataChallenge/blob/main/docs/04_LC_features.ipynb)

LC feature	Description
lcPeriodic 0-3	Four best-fit CARMA(2,1)/DHO parameters, fitted in flux using a 3-point median filter with a 5- $\sigma$ clipping for outliers removing. Fitted data are for $g$ -, $r$ - and $i$ -band light curves only. Fitted light curves have more 30 epochs.
lcPeriodic 4-32	Generalized Lomb Scargle fit/parameters as described in the R11 paper. The exact matching from index to features is given below. Note that first relative phase ( <code>rel_phase_0</code> ) are not included since it is negligible.
lcNonPeriodic 0-20	Non-periodic features introduced in the R11 paper.
lcNonPeriodic 21	Variance divided by the median.
lcNonPeriodic 22	Reduced $\chi^2$ for a constant model with a given degrees of freedom: $\chi^2/\text{d.o.f.} = \frac{1}{N-1} \sum_{i=1}^N \left( \frac{m_i - \bar{m}}{\sigma_i} \right)^2,$ <p>where <math>\bar{m}</math> is the inverse variance weighted average.</p>
lcNonPeriodic 23	Excess variance defined by Bernstein et al. (2018): $\sigma_{\text{sys}}^2 \equiv \langle \Delta m_i^2 - \sigma_{\text{stat},i}^2 \rangle$ $\Delta m_i \equiv \frac{m_i - \bar{m}}{\sqrt{1 - w_i / \sum w_j}},$ $w_i \equiv \sigma_{\text{stat},i}^{-2},$ $\bar{m} \equiv \frac{\sum w_j m_j}{\sum w_j},$
lcNonPeriodic 24	where $\sigma_{\text{sys}}$ is the excess variance and $\sigma_{\text{stat},i}$ is the photometric uncertainty. Normalized excess variance (Allevato et al. 2013): $\sigma_{\text{sys, norm}}^2 \equiv \frac{\sigma_{\text{sys}}^2}{N \bar{m}^2}.$
lcNonPeriodic 25	Range of a cumulative sum (Kim et al. 2011): $R_{\text{CS}} = \max(S) - \min(S),$ $S_l = \frac{1}{N\sigma} \sum_{i=1}^l (m_i - \bar{m}).$
lcNonPeriodic 26	The von Neumann ratio: $\eta = \frac{1}{(N-1)\sigma^2} \sum_{i=1}^{N-1} (m_{i+1} - m_i)^2.$
lcNonPeriodic 27-28	The best fit CARMA(1,0)/DRW parameters. Light curves with less than 10 epochs were not fitted.
lcNonPeriodic 27	The driving amplitude $\sigma$ , or $\beta_0$ in the CARMA notation (Kelly et al. 2009).
lcNonPeriodic 28	The characteristic timescale described by the DRW model, or $1/\alpha_1$ in the CARMA notation.

**Table 3.** LC features divided in two groups consisting of 33 periodic and 29 non-periodic features. Feature numeration starts with 0.

deviation and the coefficients of correlation between light curves in two wavebands. A second submission trained supervised models: XGB, RF, SVM, ANN and one unsupervised model (DEC) for the Star/Gal/QSO classification and clustering respectively by adding LC features. The third submission trained CNNs that use projected time series data onto 2D images for Star/Gal/QSO classification. A final submission trained a RF on a subsample consisting of stars and quasars only, when all data features, except for flags and redshift, are taken into account. A summary of the contributions, ML methods, train and test sample sizes and dimensionality  $\text{bf}(\text{number of features})$  and model performance is listed in Table 4, indicating which of the co-authors submitted the solution. For a measure of performance, we report accuracy and completeness.

In the following subsections, we elaborate in detail each of the submitted solution while following the main workflow: feature selection; preprocessing (if needed); model training and evaluating on test and blinded datasets. The uncertainties on the performance were estimated with  $k$ -fold cross-validation (Stone 1974). We split the training dataset into  $k = 10$  subsamples. One of the groups is used for testing while the rest are used for training. This process is repeated  $k$  times, with each group being used once for testing. The evaluation results are then averaged to give an overall train performance. During the cross-validation process, after every training iteration, we also evaluate the performance on the blinded dataset that is always kept aside.

#### 4.1. A preliminary classical approach – V.P. and M.P.

As already mentioned, variability is an intrinsic property of AGN and a promising selection tool, as both the time scales of the variations and the overall trends are different from other, mainly stellar, sources. Power spectra of AGN are characterized by a typical red noise behavior (Kelly et al. 2009; Kozłowski et al. 2010b; Zu et al. 2013), with most of the variation arising on long time scales (Uttley et al. 2002). This timescale behavior means that the longer the observed baseline, the better the selection through variability (see e.g. De Cicco et al. 2019). The LSST, with its dense and long coverage, will be the best survey to exploit this selection technique. Moreover, the implementation of DIA on the entire dataset will also enable the selection of low-luminosity AGN dominated by their host galaxy.

Unfortunately, the AGN DC does not contain difference images and the population of confirmed AGN in the region where optical light curves are available (SDSS Stripe82) is biased towards bright quasars. In spite of the limitation of the archival data, which prevented us

from testing the capabilities of DIA on AGN science, we examined some of the light curve features not included in the original dataset, before using ML with all the available LSST-like data products.

Intensive X-ray, UV and optical monitoring campaigns of AGN show that, whatever the intrinsic physical mechanism (thermal propagating fluctuations, hydrodynamical instabilities, reprocessing of high energy coronal photons by the accretion disk; McHardy et al. 2018) we can expect correlation between adjacent regions of the electromagnetic spectrum, such as the LSST bluer and redder bands. As SDSS images in the *ugriz*-bands were taken close in time (at 71.7 second intervals; Gunn et al. 1998), we calculated the Pearson correlation coefficient between pairs of light curves in different bands, along with the average magnitude and the standard deviation per each band. We restricted our analysis to the *gri*-bands as they have the highest signal-to-noise ratio. A preliminary analysis that include the *uz*-bands have shown to bring no improvement.

The distributions of the correlation coefficients for sources labeled as QSO, Star or Galaxy clearly show that they belong to three different populations (Fig. 6). We selected random samples of 10,000 objects and looked at how they distributed on a plot with *g*-band magnitude standard deviation vs. *gr*-bands correlation coefficient. The majority of quasars (more than 90%) tend to group in a well defined *wedge* of the space<sup>18</sup>, as shown in Fig. 7. Thus, we identified the region of interest for our sources and tested the *wedge* selection on the *blinded dataset* which was provided at the end of the AGN DC. We find that the light curve variance and correlation among bands alone, allow us to produce samples of AGN with a completeness<sup>19</sup> of 90.9%, albeit with low purity<sup>20</sup> 52.0%. The low purity is expected since these two features alone tend to identify intrinsic correlated variability above a certain variance threshold, but do not include a full characterization of AGN properties allowing one to disentangle them from, e.g., stars.

It is possible to remove most of the contaminants and reach a purity of 95% (and decreasing the completeness by less than 10%) by adding the extendedness and color information, which is particularly useful in the selection of candidate AGN (Richards et al. 2002). As our AGNs were mainly bright quasars, we made a cut requiring the LSST extendedness parameter to be greater than

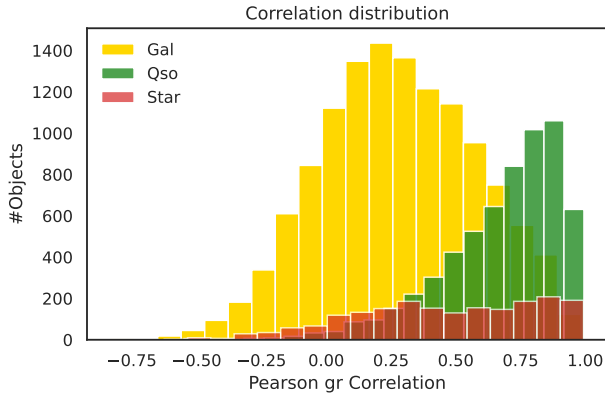
<sup>18</sup> Comprised in  $x > 0.25$ ,  $y < e^{1.3x-1.7}$ ,  $y > e^{1x-2}$ .

<sup>19</sup> The ratio of QSOs selected and total number of QSOs in the sample.

<sup>20</sup> The ratio of the selected QSOs and the total number of selected sources.

Contribution	ML method	Data type	Sample size			Performance			
			Train	Blinded	Dim.	Accuracy		Completeness	
						Train	Blinded	Train	Blinded
V.P. and M.P.	WEDGE+EXT+COL	tabular	10 000	6 000	10	$0.932 \pm 0.001$	$0.950 \pm 0.002$	$0.818 \pm 0.003$	$0.831 \pm 0.004$
	XGB	tabular	380 000	44 000	64	$0.980 \pm 0.001$	$0.970 \pm 0.001$	$0.894 \pm 0.002$	$0.834 \pm 0.001$
	XGB	tabular	128 000	15 000	64	$0.983 \pm 0.004$	$0.978 \pm 0.003$	$0.929 \pm 0.003$	$0.882 \pm 0.003$
	RF	tabular	128 000	15 000	64	$0.982 \pm 0.004$	$0.976 \pm 0.002$	$0.920 \pm 0.004$	$0.866 \pm 0.002$
Đ.S., I.J. and SER-SAG	SVM	tabular	128 000	15 000	64	$0.982 \pm 0.005$	$0.976 \pm 0.003$	$0.919 \pm 0.005$	$0.870 \pm 0.002$
	ANN	tabular	128 000	15 000	64	$0.982 \pm 0.004$	$0.975 \pm 0.004$	$0.914 \pm 0.005$	$0.858 \pm 0.006$
	ANN	tabular+im.	128 000	15 000	64	$0.982 \pm 0.005$	$0.974 \pm 0.004$	$0.913 \pm 0.005$	$0.852 \pm 0.006$
	DEC	tabular	128 000	15 000	64	$0.973 \pm 0.008$	$0.959 \pm 0.006$	$0.867 \pm 0.004$	$0.787 \pm 0.005$
W.Y.	CNN	tabular+im.	152 000	17 000	1070	$0.975 \pm 0.003$	$0.975 \pm 0.003$	$0.900 \pm 0.005$	$0.860 \pm 0.004$
G.T.R.	RF	tabular	61 000	3 000	380	$0.995 \pm 0.001$	$0.994 \pm 0.001$	$0.946 \pm 0.005$	$0.924 \pm 0.005$
L.D.	CNN	images	350 000	25 000	50 176	N/A	N/A	N/A	N/A

**Table 4.** Columns from left to right: contributions, listed by participants; ML methods used; data type used for model training (tabular, images or tabular+images); training and test sample sizes; dimensionality or the number of features used; and the performance (purity and completeness) for the train and blinded datasets. The sample sizes are rounded down to the nearest thousand. The performance of the CNN model by L.D. (bottom row) is described by [Doorenbos et al. \(2022\)](#), and the model is not applicable (N/A) to images in our datasets due to low resolution.

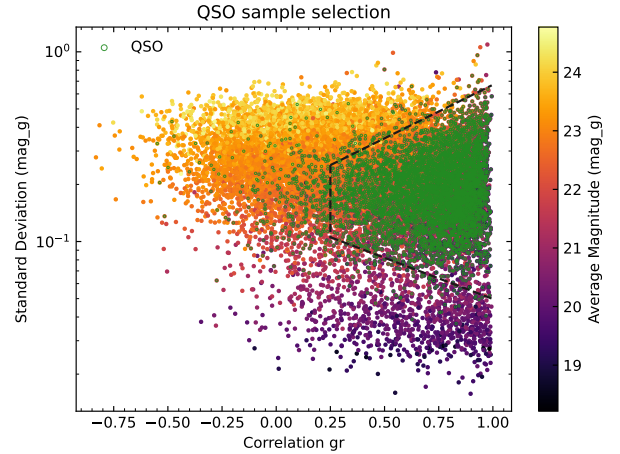


**Figure 6.** Distribution of the  $gr$ -bands Pearson correlation coefficients for galaxies (blue), QSOs (green) and stars (red).

0.95. Then, we used a  $r-i$  vs  $g-r$  color-color diagram to select candidate quasars among the sources deriving from the wedge+extendedness criteria by defining a box in which they tend to group<sup>21</sup> (see Fig. 8).

In order to evaluate the contribution of these three selection criteria, we tested them both singularly and combined together (see Table 5). The best result overall is obtained by the combination of all the selection criteria, with a purity of 95.0% (Tables 5, 4). All the methods alone, strongly suffer from contamination and the extendedness seems to show the best performance. However, this is mainly due to the bias towards point-like quasars in the AGN DC sample and does not reflect the true diversity of the AGN population. For this reason, it is also worth noting that the pair wedge+color

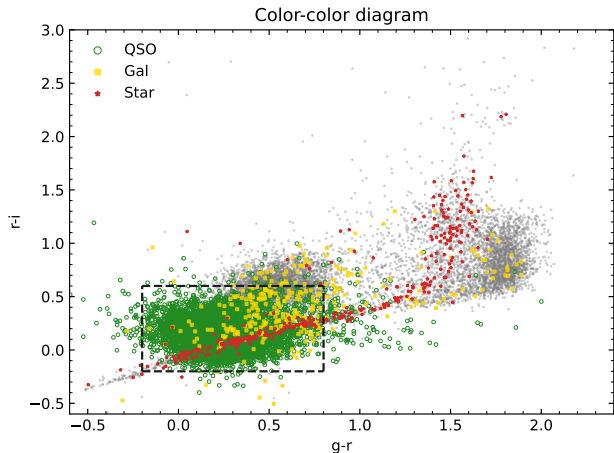
<sup>21</sup> Comprised in  $-0.2 < g-r < 0.8$  and  $-0.2 < r-i < 0.6$ .



**Figure 7.** Standard deviation of the light curves vs.  $g-r$  band correlation for a random sample of sources. Points are color-coded according to the  $g$ -band average magnitude. The black dashed lines define the *wedge* where QSOs (green points) tend to group.

returns completeness and purity of 83.0% and 85.7% respectively, without making any assumption on the morphology of the source (Table 5). This rate of success is extremely promising in the LSST perspective of applying the selection directly to sources detected on difference images, where sources will be point-like and contamination will only be due to transients, variable stars or bogus events.

We point out that the lower performance with respect to ML approaches presented in the following subsections has to be expected, since we did not use any advanced light curve feature. However, in spite of the lower performance, we demonstrate that correlation analysis among



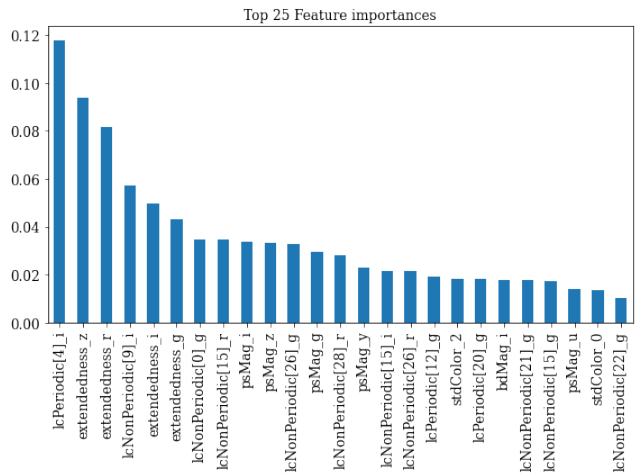
**Figure 8.** Color-color diagram showing the distribution of Galaxies, QSOs and Stars selected by the *wedge*, with a cut in LSST extendedness  $> 0.95$ . Gray dots represent the total sample, while the black dashed lines highlight the *box* where QSOs tend to group.

different bands enables a very fast and cheap first order selection of candidate QSOs (and possibly less luminous AGN). Furthermore, in the case of LSST where the low-luminosity AGN population will be detectable through DIA, correlation will help to probe the low S/N regime disentangling intrinsic variability from spurious uncorrelated noise.

#### 4.2. AGN, galaxy, star classification – *D.S., I.J. and SER-SAG*

Initial feature selection was done through trial and error. Bearing in mind that XGB supports missing (null) values by default, we first train a few XGB classifiers on a total set of  $\sim 380\,000$  containing all objects with known labels using a smaller subset of features. We further examined feature histograms with good visual separation between at least two classes (especially star-QSO and galaxy-QSO). After numerous tests, we find optimal 64 features for which we report accuracy of  $98.0 \pm 0.1\%$  and  $97.0 \pm 0.1\%$  on a test and blinded datasets. A large fraction of the dataset contains missing values for many of the features. We’ve performed a few methods of data imputation: median, hot deck and KNN imputation, however, the preliminary results are poor. Therefore, we keep the same 64 features and filter out objects with missing values for which we obtain a subset of  $\sim 128\,000$  objects divided into  $\sim 55\,000$  stars,  $\sim 45\,000$  galaxies and  $\sim 28\,000$  quasars. With such setup, we proceed with ML steps before training separate XGB, RF, SVM and NN models. The most important 20 features in the dataset from the initial XGB analysis are given in Fig. 9. More

than a half of the top 25 features are LC features, both periodic and non-periodic.



**Figure 9.** Top 25 features ranked by importance for XGB classifier run on a full dataset.

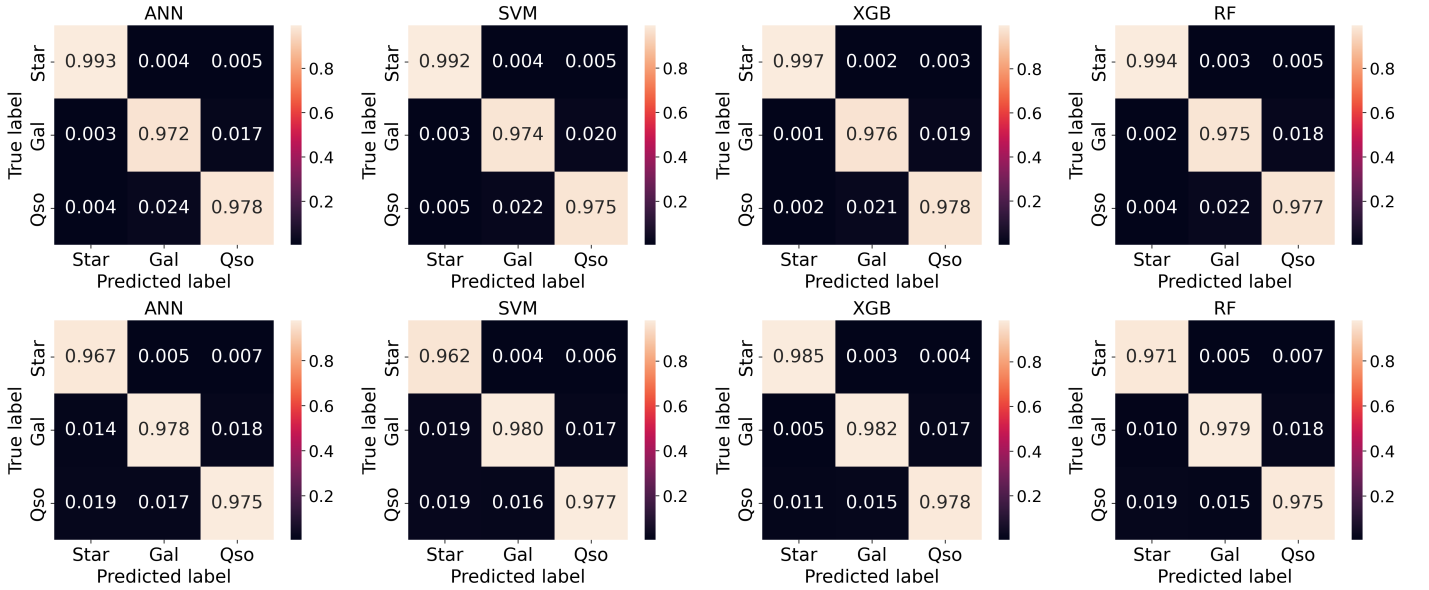
We divide the subset dataset(s) into training, validation and test sets, which account for 70, 20 and 10% of the subset object count respectively. Before applying any of the ML methods, data pre-processing is required. We standardize the training data to zero mean and standard deviation of unity. Using mean and standard deviations obtained for each train feature, we normalize validation and test sets. We used both supervised and unsupervised ML methods with the goal of comparing the performance of each method in order to establish a solid foundation towards building more advanced models. For supervised learning, we were able to achieve a high accuracy for each method (e.g., XGB accuracy of  $98.3 \pm 0.4\%$  and  $97.8 \pm 0.3\%$  for the training and blinded datasets respectively) when the light curve features are taken into account. However, the dataset is dominated by bright quasars. The XGB and RF perform the best overall on this subset of data, which is very often the case for ML applications on tabular data. The performance rating measured by classification accuracy is XGB  $>$  RF  $>$  SVM  $>$  NN (Table 4). Confusion matrices are illustrated in Fig. 10. Feature importance for the latter four methods follow the similar trend as for the initial XGB analysis. Applying the same methods on the same set when the LC features are ignored, reduces the classification accuracy to 95%.

In addition to tabular data, we investigated whether the addition of pixel information of image cutouts could contribute to achieving accuracy  $> 99\%$ . We investigated deep ANNs with joint image and tabular data as inputs, as well as using an autoencoder with a simple bottleneck architecture for image dimensionality reduction to



	Galaxy	QSO	Star	Total	Completeness (QSO)	Purity (QSO)
Sample	13066	6177	1955	21238		
Wedge	4767	5614	404	10789	90.9 %	52.0 %
Extendedness	1385	6029	1796	9211	97.6 %	65.4 %
Color	2763	5538	956	9257	90.3 %	59.8 %
Wedge+Ext	325	5492	331	6148	89.5 %	89.3 %
Col+Ext	612	5492	950	7054	89.5 %	77.8 %
Wedge+Col	716	5093	130	5989	83.0 %	85.7 %
Wedge+Ext+Col	138	5055	125	5318	82.4 %	95.0 %

**Table 5.** Results of the selection of QSOs on the blinded dataset, using a classical approach and combinations of different selection criteria.



**Figure 10.** Confusion matrices normalized by purity for tabular data. From left to right: ANN, SVM, XGB, RF. Upper panels were computed on a test set, while lower panels were computed for a blinded set. True labels are placed on vertical axis, while the predicted labels are on the horizontal axis.

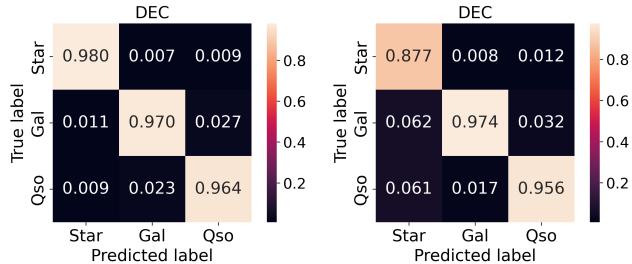
$\sim 10$  latent features that are concatenated to the tabular data. No improvement has been observed in comparison with the models trained using tabular data only. The pixel resolution is simply too low for any morphology traits to be learned by the model.

Unsupervised methods are usually statistically weaker and, are in general outperformed by supervised methods on similar problems. Essentially, clustering is grouping the instances based on their similarity without the help of class labels. In this ML paradigm, the model learns the similarities (i.e., the distance between instances in multidimensional parameter space), not the mapping function between input features and class label as is the case with supervised methods. The class naturally emerges in the representation of the data clusters. In this way, we can expect that latent features of

stars will differ at measurable level from the latent features of galaxies and QSOs. However, the number of clusters is in principle unknown. We set the number of clusters to be equal to 3, which is the number of ground-truth categories (star/galaxy/QSO). When the objects are passed to the model, it outputs soft assignments (the probabilities that given objects belong to given clusters). Each object is assigned a label of the most probable cluster. The DEC model performance is evaluated by unsupervised clustering accuracy (Xie et al. 2015, Eq. 10) computed in the following way. The optimal one-to-one matching between the set of cluster labels and the set of class labels is found, so as to maximise the agreement between the mapped cluster labels and true class labels on the given set of objects. Such mapped cluster labels constitute the predictions of the clustering model. Then

the accuracy of the predictions given the true classes is computed.

We obtained  $97.3 \pm 0.8\%$  and  $95.9 \pm 0.6\%$  (Table 4) clustering accuracy on the same test and blinded datasets respectively as for the supervised methods. When applying  $k$ -means or Gaussian mixture of models on the latent space, a clustering accuracy of  $\sim 94\%$  is obtained, however, these numbers are heavily influenced by the initial weights of the ML model. Confusion matrices for DEC are shown in Fig. 11.



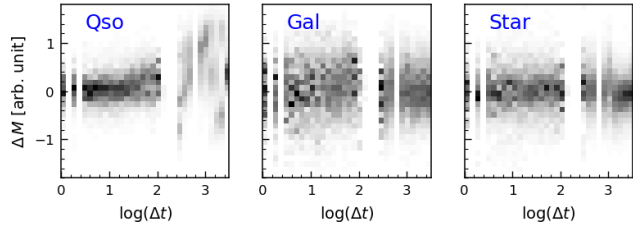
**Figure 11.** Same as Figure 10, for DEC model. The performance is shown for the test set (left) and the blinded set (right).

For the datasets used in our analysis, we expect a similar ranking of methods with respect to measured performance within a possible framework developed by Shy et al. (2022), that attempts to incorporate measurement error in astronomical classification problems. We expect a drop in performance if such experiments are repeated.

#### 4.3. Density maps – *W.Y.*

We also explored the approach of first projecting time-series data onto 2D images (i.e., density maps) and then performing classification using the constructed density maps (e.g., Mahabal et al. 2017). Density maps are essentially a 2D distribution of variability power in the timescale and magnitude space; example density maps and the distinguishing power brought by this method are shown in Figure 12. The actual algorithm used to generate those density maps is a modified version of the one used in Mahabal et al. (2017). Projecting time-series data onto fixed size 2D images allows us to exploit the power of CNNs (Lecun et al. 2015)—the gold standard for ML-based image classification.

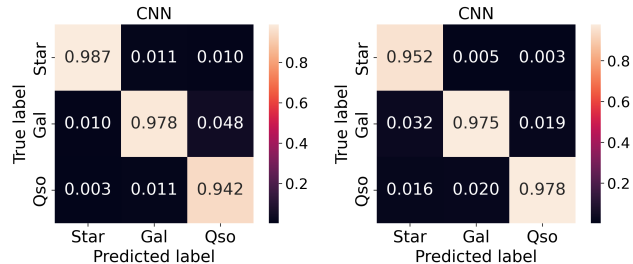
The total size of the density map sub-sample is  $\sim 152\,000$ , divided into  $\sim 32\,000$  stars,  $\sim 52\,400$  galaxies and  $\sim 67\,600$  QSOs. We removed objects with less than 5 epochs in their  $r$ -band light curves and dropped objects having missing values (e.g., NaN) for the features utilized to train the model (see next paragraph). Objects with missing values typically fall into one of two categories: 1) too faint to have been detected in all "ugrizY"



**Figure 12.** Computed density map in the  $r$ -band for one randomly selected quasar (left), galaxy (middle) and star (right).

bands; 2) too extended (and/or faint) to have a reliably determined proper motion from Gaia EDR3 (Gaia Collaboration et al. 2021) or NSC DR2 (Nidever et al. 2021). Two versions of the above dataset are created to test our ML model: one only containing objects with low temporal sampling (number of epochs is less than 30 in the  $r$ -band) and one containing 20% of all objects randomly selected from the parent sample (without any additional cut regarding temporal sampling).

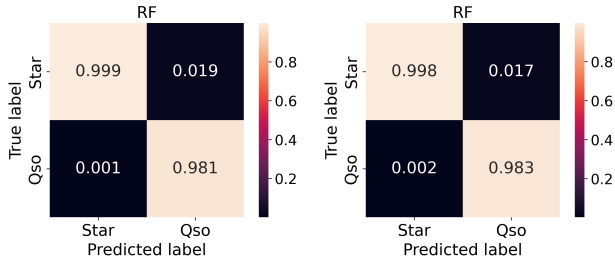
Our final classification model starts with extracting latent features from the constructed density maps using a convolutional neural network and then concatenate the latent features with a selection of features/columns from the Object table to form the final input for ANN that is used for classification. The columns picked from the Object table are the five optical colors (i.e., stdColor[0-4]), proper motion, and eight time-series features (i.e., lcNonPeriodic[9, 13, 15, 19, 20, 22, 25, 26]). The time-series features were cherry-picked from the feature importance rank of a random-forest classifier trained/tested on the full catalog of 374 features. A very high accuracy of 97.5% (Table 4) were achieved for both low-cadence and high-cadence sub-samples, which demonstrates the robustness of the classification trained on density maps which are computationally inexpensive. Confusion matrices are shown in Figure 13.



**Figure 13.** Confusion matrices for a CNN model with density maps included. The performance is shown for the test set (left) and the blinded set (right).

#### 4.4. Separating quasars from stars – *G.T.R*

In this experiment, we use all data features except for flags and the redshift, which results in a subsample  $\sim 10\,000$  quasars and  $\sim 50\,000$  stars after excluding all objects with any missing values. The excluded objects either fall into the same two categories of objects containing missing values listed in Section 4.3 or have less than 30 epochs in their *r*-band light curves. The strict cut on temporal sampling is given by the fact that only light curves with more than 30 epochs are fitted with CARMA(2,1), whose parameters are `lcPeriodic[0-3]`. The main difference between this experiment and the one in Section 4.2 is that the training dataset contains more data features while at the same time fewer objects due to filtering out those with missing and/or NaN values. We used RF for separating quasars from stars with an incredibly high accuracy of 99.6% (Table 4) and a similar performance on the blinded dataset. The white dwarfs and M-type stars can be seen to not be confused for quasars, which is a common source of error. Confusion matrices are shown in Fig. 14.



**Figure 14.** Confusion matrices for a star/quasar RF classification model. The performance is shown for the test set (left) and the blinded set (right).

## 5. SUMMARY AND DISCUSSION

The LSST AGNSC-DC was designed to fulfill the growing need for efficient ML-based AGN selection methods. While we focus on simulated LSST data, these tools should benefit a much wider AGN community. We addressed the problem of star, galaxy, and quasar classification from both a classical and a ML perspective while mimicking the future LSST catalog data. We used both supervised and unsupervised ML approaches. We followed the standard procedure of dividing the dataset into training, validation, and test datasets. The blinded dataset was revealed much later, after the submission of proposed solutions was finished. The performance of each method, sample size and the dimensionality is listed in the Table 4. We obtain high performance for supervised models and slightly lower for the unsupervised models for train and test datasets. We obtain slightly

lower performance on the blinded dataset. The addition of LC features engineered from the *gri*-bands significantly improves the classification accuracy by a few percents, however, the dataset is skewed toward bright quasars. Fainter AGNs where the host galaxy contamination is strong, as well as sources in the regions where colors overlap will strongly benefit from variability (De Cicco et al. 2019). Further improvements of the ML methods that rely on LC features will be the inclusion of LC features computed using the missing *uzy*-bands. We are limited by the number of visits per object that is between 30 and 70 (Fig. 5), and uneven temporal sampling (Fig. 4), which are the two main obstacles for direct application of DL, but at the same time, the main reason for adopting LC feature engineering. The number of visits, higher by at least an order of magnitude, will allow us to train DL models directly on the light curves.

In addition, we find that 64x64 pixel image cutouts are not sufficient for extracting pixel-level information. However, in the parallel effort done on low-redshift/low-luminosity AGNs (Doorenbos et al. 2022), it was found that images of resolution 224x224 pixel can be used to directly extract meaningful information that can help in disentangling AGN vs non-AGN hosts. Moreover, the quality of images in all 6 bands that will be provided by LSST will allow us to develop DL models for separating AGNs from galaxies in a highly efficient manner at the cost of additional computing power.

In addition to LSST data, we can also make use of the data from other surveys – not just in the optical, but also in the X-ray, ultraviolet, infrared and radio wavebands. Such data allow us to make use of differences in object SEDs across a more extended wavelength range. For example, Myers et al. (2015) have shown that adding just a single infrared data point to SDSS *ugriz* optical bands significantly improves the classification probability for AGNs vs stars. However, these other surveys do not cover the same area of sky and adding multi-wavelength data will have the effect of breaking our largely monolithic LSST survey area into considerably smaller regions. Surveys at different wavelengths most often have different spatial resolutions, leading to situations where a single source in one survey is associated with multiple sources in another waveband. Forced photometry or probabilistic cross-matching are techniques to mitigate this problem (Budavári & Szalay 2008; Lang et al. 2016; Nyland et al. 2017; Buchner et al. 2021; Salvato et al. 2022). Simple recognition of the problem may enable it to be avoided in many cases; however, the depth of LSST means that there will be few sources that are truly isolated.

We conclude that LSST-like data enable the development of highly accurate star/galaxy/quasar classifiers, mainly using only *gri*-bands and without spectral information. Promising results have been achieved in spite of limitations in resolution, survey area, cadence, and baseline of the SDSS data (compared to what LSST will be capable of). Thus the new survey will allow us not only to improve the performances of the algorithms successfully tested within the DC, but also to develop more sophisticated techniques e.g., exploitation of pixel-level information on source cutouts, analysis of the population of weak AGNs dominated by their host galaxy through DIA and identifying new types of AGNs. The AGN DC collection of multi-wavelength data is an important legacy for future research beyond AGNs science that will be included in the follow up work.

#### DATA STATEMENT

The DC was hosted on the SciServer<sup>22</sup> platform. Each participant of the DC had a verified account on the platform. The datasets released in the DC are publicly available on Zenodo<sup>23</sup> (Yu et al. 2022) under a Creative Commons Attribution 4.0 International Public License. The notebooks which reproduce the results of this work are designed as end-to-end and follow the flowchart of this work. The notebooks are publicly available on GitHub<sup>24</sup>. The notebooks have been successfully executed on different servers with different operating systems which can also support GPU acceleration, including SciServer.

#### AUTHOR CONTRIBUTION STATEMENT

Đ.S. and I.J. submitted the winning solution to the data challenge and drafted the manuscript; W.Y. and G.T.R. curated the data challenge; M.J.T., Q.N. and R.S. supplied the multi-wavelength data and associated paper text; Đ.S., I.J., W.Y., V.P., A.K., M.N., D.I., L.P., M.P., A.C. and G.T.R. contributed to solutions to the data challenge. All authors contributed to the writing and editing of the draft.

#### ACKNOWLEDGEMENTS

Prizes for participating in data challenge were funded by the LSST Corporation’s Enabling Science Program. Đ.S. acknowledges the support by the F.R.S. FNRS

under grant PDR T.0116.21. Đ.S. and L.Č.P. acknowledge support by the Astronomical Observatory (the contract №451-03-68/2022-14/200002), through the grants by the Ministry of Education, Science, and Technological Development of the Republic of Serbia. Đ.S. acknowledges support by the Science Fund of the Republic of Serbia, PROMIS №6060916, BOWIE. D.I., A.B.K. and L.Č.P. acknowledge funding provided by the University of Belgrade - Faculty of Mathematics (the contract №451-03-68/2022-14/200104) through the grants by the Ministry of Education, Science, and Technological Development of the Republic of Serbia. D.I. acknowledges the support of the Alexander von Humboldt Foundation. A.B.K. and L.Č.P. thank the support by Chinese Academy of Sciences President’s International Fellowship Initiative (PIFI) for visiting scientist. M.J.T. acknowledges support from ANID (Fondecyt Proyecto 3220516). S.P. acknowledges financial support from the Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) Fellowship (№164753/2020-6) and the Polish Funding Agency National Science Centre, project №2017/26/A/ST9/00756 (MAESTRO 9). A.Ć. acknowledges support from the Fermi Research Alliance, LLC under Contract №DE-AC02-07CH11359 with the U.S. Department of Energy (DOE), Office of Science, Office of High Energy Physics.

The authors thank Feige Wang and Jinyi Yang for constructing and providing the `highZQso` catalog.

This research makes use of the SciServer science platform ([www.sciserver.org](http://www.sciserver.org)). SciServer is a collaborative research environment for large-scale data-driven science. It is being developed at, and administered by, the Institute for Data Intensive Engineering and Science at Johns Hopkins University. SciServer is funded by the National Science Foundation through the Data Infrastructure Building Blocks (DIBBs) program and others, as well as by the Alfred P. Sloan Foundation and the Gordon and Betty Moore Foundation.

*Software:* PYTHON (Van Rossum & Drake Jr 1995), JUPYTER (Kluyver et al. 2016).

*ML packages:* NUMPY and SCIPY (van der Walt et al. 2011), PANDAS (McKinney et al. 2010), SCIKIT-LEARN (Pedregosa et al. 2011), KERAS (Chollet et al. 2015), TENSORFLOW (Abadi et al. 2016).

*Data visualization:* MATPLOTLIB (Hunter 2007), SEABORN (Waskom et al. 2017).

#### REFERENCES

<sup>22</sup> <https://www.sciserver.org/>

<sup>23</sup> <https://doi.org/10.5281/zenodo.6878414>

<sup>24</sup> [https://github.com/RichardsGroup/AGN\\_DataChallenge](https://github.com/RichardsGroup/AGN_DataChallenge)



- Abadi, M., Barham, P., Chen, J., et al. 2016, arXiv e-prints, arXiv:1605.08695. <https://arxiv.org/abs/1605.08695>
- Ahumada, R., Prieto, C. A., Almeida, A., et al. 2020, ApJS, 249, 3, doi: [10.3847/1538-4365/ab929e](https://doi.org/10.3847/1538-4365/ab929e)
- Ahumada, R., Prieto, C. A., Almeida, A., et al. 2020, The Astrophysical Journal Supplement Series, 249, 3, doi: [10.3847/1538-4365/ab929e](https://doi.org/10.3847/1538-4365/ab929e)
- Aihara, H., Arimoto, N., Armstrong, R., et al. 2018, PASJ, 70, S4, doi: [10.1093/pasj/psx066](https://doi.org/10.1093/pasj/psx066)
- Aihara, H., ALSayyad, Y., Ando, M., et al. 2019, arXiv:1905.12221 [astro-ph], doi: [10.1093/pasj/psz103](https://doi.org/10.1093/pasj/psz103)
- Allevato, V., Paolillo, M., Papadakis, I., & Pinto, C. 2013, ApJ, 771, 9, doi: [10.1088/0004-637X/771/1/9](https://doi.org/10.1088/0004-637X/771/1/9)
- Annis, J., Soares-Santos, M., Strauss, M. A., et al. 2014, ApJ, 794, 120, doi: [10.1088/0004-637X/794/2/120](https://doi.org/10.1088/0004-637X/794/2/120)
- Antonucci, R. 1993, ARA&A, 31, 473, doi: [10.1146/annurev.aa.31.090193.002353](https://doi.org/10.1146/annurev.aa.31.090193.002353)
- Bañados, E., Venemans, B. P., Decarli, R., et al. 2016, ApJS, 227, 11, doi: [10.3847/0067-0049/227/1/11](https://doi.org/10.3847/0067-0049/227/1/11)
- Banerji, M., Lahav, O., Lintott, C. J., et al. 2010, MNRAS, 406, 342, doi: [10.1111/j.1365-2966.2010.16713.x](https://doi.org/10.1111/j.1365-2966.2010.16713.x)
- Baron, D. 2019, arXiv e-prints, arXiv:1904.07248. <https://arxiv.org/abs/1904.07248>
- Bellm, E. 2014, in The Third Hot-wiring the Transient Universe Workshop, ed. P. R. Wozniak, M. J. Graham, A. A. Mahabal, & R. Seaman, 27–33. <https://arxiv.org/abs/1410.8185>
- Bernstein, G. M., Abbott, T. M. C., Armstrong, R., et al. 2018, PASP, 130, 054501, doi: [10.1088/1538-3873/aaa753](https://doi.org/10.1088/1538-3873/aaa753)
- Berry, M. W., Mohamed, A., & Yap, B. W. 2019, Supervised and Unsupervised Learning for Data Science, 1st edn. (Springer Publishing Company, Incorporated)
- Bianco, F. B., Ivezić, Ž., Jones, R. L., et al. 2022, ApJS, 258, 1, doi: [10.3847/1538-4365/ac3e72](https://doi.org/10.3847/1538-4365/ac3e72)
- Bonoli, F., Braccisi, A., Federici, L., Zitelli, V., & Formiggini, L. 1979, A&AS, 35, 391
- Bosch, J., Armstrong, R., Bickerton, S., et al. 2018, PASJ, 70, S5, doi: [10.1093/pasj/psx080](https://doi.org/10.1093/pasj/psx080)
- Bramich, D. M. 2008, MNRAS, 386, L77, doi: [10.1111/j.1745-3933.2008.00464.x](https://doi.org/10.1111/j.1745-3933.2008.00464.x)
- Brandt, W. N., & Alexander, D. M. 2015, A&A Rv, 23, 1, doi: [10.1007/s00159-014-0081-z](https://doi.org/10.1007/s00159-014-0081-z)
- Buchner, J., Salvato, M., Budavári, T., & Fotopoulou, S. 2021, nway: Bayesian cross-matching of astronomical catalogs, Astrophysics Source Code Library, record ascl:2102.014. <http://ascl.net/2102.014>
- Budavári, T., & Szalay, A. S. 2008, ApJ, 679, 301, doi: [10.1086/587156](https://doi.org/10.1086/587156)
- Butler, N. R., & Bloom, J. S. 2011, AJ, 141, 93, doi: [10.1088/0004-6256/141/3/93](https://doi.org/10.1088/0004-6256/141/3/93)
- Carballo, R., González-Serrano, J. I., Benn, C. R., & Jiménez-Luján, F. 2008, MNRAS, 391, 369, doi: [10.1111/j.1365-2966.2008.13896.x](https://doi.org/10.1111/j.1365-2966.2008.13896.x)
- Cavuoti, S., Brescia, M., D’Abrusco, R., Longo, G., & Paolillo, M. 2014, MNRAS, 437, 968, doi: [10.1093/mnras/stt1961](https://doi.org/10.1093/mnras/stt1961)
- Chambers, K. C., Magnier, E. A., Metcalfe, N., et al. 2016, arXiv e-prints, arXiv:1612.05560. <https://arxiv.org/abs/1612.05560>
- Chang, Y.-Y., Hsieh, B.-C., Wang, W.-H., et al. 2021, ApJ, 920, 68, doi: [10.3847/1538-4357/ac167c](https://doi.org/10.3847/1538-4357/ac167c)
- Chen, B. H., Goto, T., Kim, S. J., et al. 2021, MNRAS, 501, 3951, doi: [10.1093/mnras/staa3865](https://doi.org/10.1093/mnras/staa3865)
- Chen, C. T. J., Brandt, W. N., Luo, B., et al. 2018, MNRAS, 478, 2132, doi: [10.1093/mnras/sty1036](https://doi.org/10.1093/mnras/sty1036)
- Chen, T., & Guestrin, C. 2016, arXiv e-prints, arXiv:1603.02754. <https://arxiv.org/abs/1603.02754>
- Chollet, F., et al. 2015, Keras, GitHub. <https://github.com/fchollet/keras>
- Clarke, A. O., Scaife, A. M. M., Greenhalgh, R., & Griguta, V. 2020, A&A, 639, A84, doi: [10.1051/0004-6361/201936770](https://doi.org/10.1051/0004-6361/201936770)
- Cortes, C., & Vapnik, V. 1995, Machine learning, 20, 273
- Cybenko, G. V. 1989, Mathematics of Control, Signals and Systems, 2, 303
- Czerny, B., Cao, S., Jaiswal, V. K., et al. 2022, arXiv e-prints, arXiv:2209.06563. <https://arxiv.org/abs/2209.06563>
- Czerny, B., Panda, S., Prince, R., et al. 2023, arXiv e-prints, arXiv:2301.08975, doi: [10.48550/arXiv.2301.08975](https://doi.org/10.48550/arXiv.2301.08975)
- Dark Energy Survey Collaboration, Abbott, T., Abdalla, F. B., et al. 2016, Mon Not R Astron Soc, 460, 1270, doi: [10.1093/mnras/stw641](https://doi.org/10.1093/mnras/stw641)
- De Cicco, D., Paolillo, M., Falocco, S., et al. 2019, A&A, 627, A33, doi: [10.1051/0004-6361/201935659](https://doi.org/10.1051/0004-6361/201935659)
- De Cicco, D., Bauer, F. E., Paolillo, M., et al. 2021, A&A, 645, A103, doi: [10.1051/0004-6361/202039193](https://doi.org/10.1051/0004-6361/202039193)
- Delgado, F., Saha, A., Chandrasekharan, S., et al. 2014, in Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, Vol. 9150, Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, 15, doi: [10.1117/12.2056898](https://doi.org/10.1117/12.2056898)
- DES Collaboration, Abbott, T. M. C., Adamow, M., et al. 2021, arXiv:2101.05765 [astro-ph]. <http://ascl.net/2101.05765>
- Dey, A., Schlegel, D. J., Lang, D., et al. 2019, AJ, 157, 168, doi: [10.3847/1538-3881/ab089d](https://doi.org/10.3847/1538-3881/ab089d)



- Doert, M., & Errando, M. 2014, *ApJ*, 782, 41, doi: [10.1088/0004-637X/782/1/41](https://doi.org/10.1088/0004-637X/782/1/41)
- Doorenbos, L., Torbaniuk, O., Cavuoti, S., et al. 2022, *A&A*, 666, A171, doi: [10.1051/0004-6361/202243900](https://doi.org/10.1051/0004-6361/202243900)
- Dye, S., Lawrence, A., Read, M. A., et al. 2018, *MNRAS*, 473, 5113, doi: [10.1093/mnras/stx2622](https://doi.org/10.1093/mnras/stx2622)
- Eckert, D., Gaspari, M., Gastaldello, F., Le Brun, A. M. C., & O'Sullivan, E. 2021, *Universe*, 7, 142, doi: [10.3390/universe7050142](https://doi.org/10.3390/universe7050142)
- Fabian, A. C. 2012, *ARA&A*, 50, 455, doi: [10.1146/annurev-astro-081811-125521](https://doi.org/10.1146/annurev-astro-081811-125521)
- Fan, X., Strauss, M. A., Richards, G. T., et al. 2006, *AJ*, 131, 1203, doi: [10.1086/500296](https://doi.org/10.1086/500296)
- Ferrarese, L., & Merritt, D. 2000, *ApJL*, 539, L9, doi: [10.1086/312838](https://doi.org/10.1086/312838)
- Gaia Collaboration, Prusti, T., de Bruijne, J. H. J., et al. 2016, *A&A*, 595, A1, doi: [10.1051/0004-6361/201629272](https://doi.org/10.1051/0004-6361/201629272)
- Gaia Collaboration, Brown, A. G. A., Vallenari, A., et al. 2018, *A&A*, 616, A1, doi: [10.1051/0004-6361/201833051](https://doi.org/10.1051/0004-6361/201833051)
- . 2021, *A&A*, 649, A1, doi: [10.1051/0004-6361/202039657](https://doi.org/10.1051/0004-6361/202039657)
- Gebhardt, K., Kormendy, J., Ho, L. C., et al. 2000, *ApJL*, 543, L5, doi: [10.1086/318174](https://doi.org/10.1086/318174)
- Gunn, J. E., Carr, M., Rockosi, C., et al. 1998, *AJ*, 116, 3040, doi: [10.1086/300645](https://doi.org/10.1086/300645)
- Guo, X., Gao, L., Liu, X., & Yin, J. 2017, in *Ijcai*, 1753–1759
- Gwyn, S. D. J. 2012, *The Astronomical Journal*, 143, 38, doi: [10.1088/0004-6256/143/2/38](https://doi.org/10.1088/0004-6256/143/2/38)
- Hambly, N. C., Collins, R. S., Cross, N. J. G., et al. 2008, *MNRAS*, 384, 637, doi: [10.1111/j.1365-2966.2007.12700.x](https://doi.org/10.1111/j.1365-2966.2007.12700.x)
- Hewett, P. C., Warren, S. J., Leggett, S. K., & Hodgkin, S. T. 2006, *MNRAS*, 367, 454, doi: [10.1111/j.1365-2966.2005.09969.x](https://doi.org/10.1111/j.1365-2966.2005.09969.x)
- Hložek, R., Ponder, K. A., Malz, A. I., et al. 2020, *arXiv e-prints*, arXiv:2012.12392, <https://arxiv.org/abs/2012.12392>
- Ho, T. K. 1995, in *Proceedings of 3rd international conference on document analysis and recognition*, Vol. 1, IEEE, 278–282
- Hodgkin, S. T., Irwin, M. J., Hewett, P. C., & Warren, S. J. 2009, *MNRAS*, 394, 675, doi: [10.1111/j.1365-2966.2008.14387.x](https://doi.org/10.1111/j.1365-2966.2008.14387.x)
- Hunter, J. D. 2007, *Computing in Science Engineering*, 9, 90, doi: [10.1109/MCSE.2007.55](https://doi.org/10.1109/MCSE.2007.55)
- Ivezić, Ž., Smith, J. A., Miknaitis, G., et al. 2007, *AJ*, 134, 973, doi: [10.1086/519976](https://doi.org/10.1086/519976)
- Ivezić, Ž., Kahn, S. M., Tyson, J. A., et al. 2019, *ApJ*, 873, 111, doi: [10.3847/1538-4357/ab042c](https://doi.org/10.3847/1538-4357/ab042c)
- Jankov, I., Ilić, D., & Kovačević, A. 2021, in *XIX Serbian Astronomical Conference*, Vol. 100, 241–246
- Jarvis, M. J., Bonfield, D. G., Bruce, V. A., et al. 2012, *Monthly Notices of the Royal Astronomical Society*, 428, 1281, doi: [10.1093/mnras/sts118](https://doi.org/10.1093/mnras/sts118)
- Jiang, L., McGreer, I. D., Fan, X., et al. 2016, *ApJ*, 833, 222, doi: [10.3847/1538-4357/833/2/222](https://doi.org/10.3847/1538-4357/833/2/222)
- Jolliffe, I. T. 1986, *Principal component analysis*
- Kaczmarszyc, M. C., Richards, G. T., Mehta, S. S., & Schlegel, D. J. 2009, *AJ*, 138, 19, doi: [10.1088/0004-6256/138/1/19](https://doi.org/10.1088/0004-6256/138/1/19)
- Kasliwal, V. P., Vogeley, M. S., & Richards, G. T. 2017, *Mon Not R Astron Soc*, 470, 3027, doi: [10.1093/mnras/stx1420](https://doi.org/10.1093/mnras/stx1420)
- Kelly, B. C., Bechtold, J., & Siemiginowska, A. 2009, *ApJ*, 698, 895, doi: [10.1088/0004-637X/698/1/895](https://doi.org/10.1088/0004-637X/698/1/895)
- Kim, D.-W., Protopapas, P., Byun, Y.-I., et al. 2011, *ApJ*, 735, 68, doi: [10.1088/0004-637X/735/2/68](https://doi.org/10.1088/0004-637X/735/2/68)
- Kluyver, T., Ragan-Kelley, B., Pérez, F., et al. 2016, in *Positioning and Power in Academic Publishing: Players, Agents and Agendas*, ed. F. Loizides & B. Schmidt (IOS Press), 87–90. <https://eprints.soton.ac.uk/403913/>
- Koo, D. C., & Kron, R. G. 1982, *A&A*, 105, 107
- Koo, D. C., Kron, R. G., & Cudworth, K. M. 1986, *PASP*, 98, 285, doi: [10.1086/131756](https://doi.org/10.1086/131756)
- Kormendy, J., & Ho, L. C. 2013, *ARA&A*, 51, 511, doi: [10.1146/annurev-astro-082708-101811](https://doi.org/10.1146/annurev-astro-082708-101811)
- Kovacevic, A. B., Radovic, V., Ilic, D., et al. 2022, *arXiv e-prints*, arXiv:2208.06203, <https://arxiv.org/abs/2208.06203>
- Kozłowski, S., Kochanek, C. S., Ashby, M. L. N., et al. 2016, *ApJ*, 817, 119, doi: [10.3847/0004-637X/817/2/119](https://doi.org/10.3847/0004-637X/817/2/119)
- Kozłowski, S., Kochanek, C. S., Stern, D., et al. 2010a, *ApJ*, 716, 530, doi: [10.1088/0004-637X/716/1/530](https://doi.org/10.1088/0004-637X/716/1/530)
- Kozłowski, S., Kochanek, C. S., Udalski, A., et al. 2010b, *ApJ*, 708, 927, doi: [10.1088/0004-637X/708/2/927](https://doi.org/10.1088/0004-637X/708/2/927)
- Kramer, M. A. 1991, *Aiche Journal*, 37, 233
- Kron, R. G., & Chiu, L. T. G. 1981, *PASP*, 93, 397, doi: [10.1086/130845](https://doi.org/10.1086/130845)
- Lang, D., Hogg, D. W., & Mykytyn, D. 2016, *The Tractor: Probabilistic astronomical source detection and measurement*, *Astrophysics Source Code Library*, record ascl:1604.008. <http://ascl.net/1604.008>
- Lawrence, A., Warren, S. J., Almaini, O., et al. 2007, *MNRAS*, 379, 1599, doi: [10.1111/j.1365-2966.2007.12040.x](https://doi.org/10.1111/j.1365-2966.2007.12040.x)
- Lecun, Y., Bengio, Y., & Hinton, G. 2015, *Nature*, 521, 436, doi: [10.1038/nature14539](https://doi.org/10.1038/nature14539)
- Lloyd, S. 1982, *IEEE Transactions on Information Theory*, 28, 129, doi: [10.1109/TIT.1982.1056489](https://doi.org/10.1109/TIT.1982.1056489)

- Lochner, M., McEwen, J. D., Peiris, H. V., Lahav, O., & Winter, M. K. 2016, *ApJS*, 225, 31, doi: [10.3847/0067-0049/225/2/31](https://doi.org/10.3847/0067-0049/225/2/31)
- LSST Science Collaboration, Marshall, P., Anguita, T., et al. 2017, arXiv e-prints, arXiv:1708.04058. <https://arxiv.org/abs/1708.04058>
- Luo, B., Brandt, W. N., Xue, Y. Q., et al. 2017, *ApJS*, 228, 2, doi: [10.3847/1538-4365/228/1/2](https://doi.org/10.3847/1538-4365/228/1/2)
- Macuga, M., Martini, P., Miller, E. D., et al. 2019, *ApJ*, 874, 54, doi: [10.3847/1538-4357/ab0746](https://doi.org/10.3847/1538-4357/ab0746)
- Mahabal, A., Sheth, K., Gieseke, F., et al. 2017, arXiv e-prints, arXiv:1709.06257. <https://arxiv.org/abs/1709.06257>
- Matsuoka, Y., Iwasawa, K., Onoue, M., et al. 2018, *ApJS*, 237, 5, doi: [10.3847/1538-4365/aac724](https://doi.org/10.3847/1538-4365/aac724)
- Mazzucchelli, C., Bañados, E., Venemans, B. P., et al. 2017, *ApJ*, 849, 91, doi: [10.3847/1538-4357/aa9185](https://doi.org/10.3847/1538-4357/aa9185)
- McGreer, I. D., Jiang, L., Fan, X., et al. 2013, *ApJ*, 768, 105, doi: [10.1088/0004-637X/768/2/105](https://doi.org/10.1088/0004-637X/768/2/105)
- McHardy, I. M., Connolly, S. D., Horne, K., et al. 2018, *MNRAS*, 480, 2881, doi: [10.1093/mnras/sty1983](https://doi.org/10.1093/mnras/sty1983)
- McKinney, W., et al. 2010, in *Proceedings of the 9th Python in Science Conference*, Vol. 445, Austin, TX, 51–56
- McMahon, R. G., Banerji, M., Gonzalez, E., et al. 2013, *The Messenger*, 154, 35
- Moreno, J., Vogeley, M. S., Richards, G. T., & Yu, W. 2019, *PASP*, 131, 063001, doi: [10.1088/1538-3873/ab1597](https://doi.org/10.1088/1538-3873/ab1597)
- Myers, A. D., Palanque-Delabrouille, N., Prakash, A., et al. 2015, *ApJS*, 221, 27, doi: [10.1088/0067-0049/221/2/27](https://doi.org/10.1088/0067-0049/221/2/27)
- Netzer, H. 2015, *ARA&A*, 53, 365, doi: [10.1146/annurev-astro-082214-122302](https://doi.org/10.1146/annurev-astro-082214-122302)
- Ni, Q., Brandt, W. N., Chen, C.-T., et al. 2021, *ApJS*, 256, 21, doi: [10.3847/1538-4365/ac0dc6](https://doi.org/10.3847/1538-4365/ac0dc6)
- Nidever, D. L., Dey, A., Fasbender, K., et al. 2021, *AJ*, 161, 192, doi: [10.3847/1538-3881/abd6e1](https://doi.org/10.3847/1538-3881/abd6e1)
- Nyland, K., Lacy, M., Sajina, A., et al. 2017, *ApJS*, 230, 9, doi: [10.3847/1538-4365/aa6fed](https://doi.org/10.3847/1538-4365/aa6fed)
- Padovani, P., Alexander, D. M., Assef, R. J., et al. 2017, *A&A Rv*, 25, 2, doi: [10.1007/s00159-017-0102-9](https://doi.org/10.1007/s00159-017-0102-9)
- Panda, S., Martínez-Aldama, M. L., & Zajaček, M. 2019, *Frontiers in Astronomy and Space Sciences*, 6, 75, doi: [10.3389/fspas.2019.00075](https://doi.org/10.3389/fspas.2019.00075)
- Pedregosa, F., Varoquaux, G., Gramfort, A., et al. 2011, *Journal of machine learning research*, 12, 2825
- Peters, C. M., Richards, G. T., Myers, A. D., et al. 2015, *ApJ*, 811, 95, doi: [10.1088/0004-637X/811/2/95](https://doi.org/10.1088/0004-637X/811/2/95)
- Pierre, M., Chiappetti, L., Pacaud, F., et al. 2007, *MNRAS*, 382, 279, doi: [10.1111/j.1365-2966.2007.12354.x](https://doi.org/10.1111/j.1365-2966.2007.12354.x)
- Poliszczuk, A., Pollo, A., Małek, K., et al. 2021, *A&A*, 651, A108, doi: [10.1051/0004-6361/202040219](https://doi.org/10.1051/0004-6361/202040219)
- Poulain, M., Paolillo, M., De Cicco, D., et al. 2020, *A&A*, 634, A50, doi: [10.1051/0004-6361/201937108](https://doi.org/10.1051/0004-6361/201937108)
- Pozo Nuñez, F., Bruckmann, C., Deesamutara, S., et al. 2023, *MNRAS*, 522, 2002, doi: [10.1093/mnras/stad286](https://doi.org/10.1093/mnras/stad286)
- Raiteri, C. M., Carnerero, M. I., Balmaverde, B., et al. 2022, *ApJS*, 258, 3, doi: [10.3847/1538-4365/ac3bb0](https://doi.org/10.3847/1538-4365/ac3bb0)
- Reed, S. L., McMahon, R. G., Martini, P., et al. 2017, *MNRAS*, 468, 4702, doi: [10.1093/mnras/stx728](https://doi.org/10.1093/mnras/stx728)
- Reynolds, D. 2009, *Gaussian Mixture Models* (Boston, MA: Springer US), 659–663, doi: [10.1007/978-0-387-73003-5\\_196](https://doi.org/10.1007/978-0-387-73003-5_196)
- Richards, G. T., Fan, X., Newberg, H. J., et al. 2002, *AJ*, 123, 2945, doi: [10.1086/340187](https://doi.org/10.1086/340187)
- Richards, J. W., Starr, D. L., Butler, N. R., et al. 2011, *ApJ*, 733, 10, doi: [10.1088/0004-637X/733/1/10](https://doi.org/10.1088/0004-637X/733/1/10)
- Risaliti, G., & Lusso, E. 2019, *Nature Astronomy*, 3, 272, doi: [10.1038/s41550-018-0657-z](https://doi.org/10.1038/s41550-018-0657-z)
- Salpeter, E. E. 1964, *ApJ*, 140, 796, doi: [10.1086/147973](https://doi.org/10.1086/147973)
- Salvato, M., Wolf, J., Dwelly, T., et al. 2022, *A&A*, 661, A3, doi: [10.1051/0004-6361/202141631](https://doi.org/10.1051/0004-6361/202141631)
- Sánchez, J., Walter, C. W., Awan, H., et al. 2020, *MNRAS*, 497, 210, doi: [10.1093/mnras/staa1957](https://doi.org/10.1093/mnras/staa1957)
- Sandage, A., & Luyten, W. J. 1967, *ApJ*, 148, 767, doi: [10.1086/149200](https://doi.org/10.1086/149200)
- Schmidt, K. B., Marshall, P. J., Rix, H.-W., et al. 2010, *ApJ*, 714, 1194, doi: [10.1088/0004-637X/714/2/1194](https://doi.org/10.1088/0004-637X/714/2/1194)
- Shirley, R., Duncan, K., Campos Varillas, M. C., et al. 2021, *Monthly Notices of the Royal Astronomical Society*, 507, 129, doi: [10.1093/mnras/stab1526](https://doi.org/10.1093/mnras/stab1526)
- Shy, S., Tak, H., Feigelson, E. D., Timlin, J. D., & Babu, G. J. 2022, *AJ*, 164, 6, doi: [10.3847/1538-3881/ac6e64](https://doi.org/10.3847/1538-3881/ac6e64)
- Smith, M. J., & Geach, J. E. 2022, arXiv e-prints, arXiv:2211.03796. <https://arxiv.org/abs/2211.03796>
- Stone, M. 1974, *Journal of the Royal Statistical Society: Series B (Methodological)*, 36, 111, doi: <https://doi.org/10.1111/j.2517-6161.1974.tb00994.x>
- Suberlak, K. L., Ivezić, Ž., & MacLeod, C. 2021, *ApJ*, 907, 96, doi: [10.3847/1538-4357/abc698](https://doi.org/10.3847/1538-4357/abc698)
- Tan, C., Sun, F., Kong, T., et al. 2018, arXiv e-prints, arXiv:1808.01974. <https://arxiv.org/abs/1808.01974>
- Temple, M. J., Hewett, P. C., & Banerji, M. 2021, *MNRAS*, 508, 737, doi: [10.1093/mnras/stab2586](https://doi.org/10.1093/mnras/stab2586)
- Trevese, D., Pittella, G., Kron, R. G., Koo, D. C., & Bershadsky, M. 1989, *AJ*, 98, 108, doi: [10.1086/115129](https://doi.org/10.1086/115129)
- Uttley, P., McHardy, I. M., & Papadakis, I. E. 2002, *MNRAS*, 332, 231, doi: [10.1046/j.1365-8711.2002.05298.x](https://doi.org/10.1046/j.1365-8711.2002.05298.x)

- van der Maaten, L., & Hinton, G. 2008, *Journal of Machine Learning Research*, 9, 2579.  
<http://www.jmlr.org/papers/v9/vandermaaten08a.html>
- van der Walt, S., Colbert, S. C., & Varoquaux, G. 2011, *Computing in Science and Engineering*, 13, 22, doi: [10.1109/MCSE.2011.37](https://doi.org/10.1109/MCSE.2011.37)
- Van Rossum, G., & Drake Jr, F. L. 1995, *Python reference manual (Centrum voor Wiskunde en Informatica Amsterdam)*
- Venemans, B. P., Bañados, E., Decarli, R., et al. 2015, *ApJL*, 801, L11, doi: [10.1088/2041-8205/801/1/L11](https://doi.org/10.1088/2041-8205/801/1/L11)
- Wang, F., Wu, X.-B., Fan, X., et al. 2016, *ApJ*, 819, 24, doi: [10.3847/0004-637X/819/1/24](https://doi.org/10.3847/0004-637X/819/1/24)
- Wang, F., Yang, J., Fan, X., et al. 2019, *ApJ*, 884, 30, doi: [10.3847/1538-4357/ab2be5](https://doi.org/10.3847/1538-4357/ab2be5)
- Warren, S. J., Hewett, P. C., Irwin, M. J., & Osmer, P. S. 1991, *ApJS*, 76, 1, doi: [10.1086/191563](https://doi.org/10.1086/191563)
- Waskom, M., Botvinnik, O., O’Kane, D., et al. 2017, *mwaskom/seaborn: v0.8.1 (September 2017)*, v0.8.1, Zenodo, doi: [10.5281/zenodo.883859](https://doi.org/10.5281/zenodo.883859)
- Willott, C. J., Delorme, P., Reylé, C., et al. 2010, *AJ*, 139, 906, doi: [10.1088/0004-6256/139/3/906](https://doi.org/10.1088/0004-6256/139/3/906)
- Wright, E. L., Eisenhardt, P. R. M., Mainzer, A. K., et al. 2010, *AJ*, 140, 1868, doi: [10.1088/0004-6256/140/6/1868](https://doi.org/10.1088/0004-6256/140/6/1868)
- Xie, J., Girshick, R., & Farhadi, A. 2015, arXiv e-prints, arXiv:1511.06335. <https://arxiv.org/abs/1511.06335>
- Yang, J., Wang, F., Fan, X., et al. 2019a, *ApJ*, 871, 199, doi: [10.3847/1538-4357/aaf858](https://doi.org/10.3847/1538-4357/aaf858)
- . 2019b, *AJ*, 157, 236, doi: [10.3847/1538-3881/ab1be1](https://doi.org/10.3847/1538-3881/ab1be1)
- York, D. G., Adelman, J., Anderson, John E., J., et al. 2000a, *AJ*, 120, 1579, doi: [10.1086/301513](https://doi.org/10.1086/301513)
- . 2000b, *AJ*, 120, 1579, doi: [10.1086/301513](https://doi.org/10.1086/301513)
- Yu, W., & Richards, G. T. 2021, *LSSTC AGN Data Challenge 2021*, Tech. rep., GitHub, doi: [10.17918/AGN\\_DataChallenge](https://doi.org/10.17918/AGN_DataChallenge)
- Yu, W., & Richards, G. T. 2022, *EzTao: Easier CARMA Modeling*. <http://ascl.net/2201.001>
- Yu, W., Richards, G. T., Vogeley, M. S., Moreno, J., & Graham, M. J. 2022, *ApJ*, 936, 132, doi: [10.3847/1538-4357/ac8351](https://doi.org/10.3847/1538-4357/ac8351)
- Yu, W., Richards, G. T., Yoachim, P., & Peters, C. 2020, *Research Notes of the American Astronomical Society*, 4, 252, doi: [10.3847/2515-5172/abd6e2](https://doi.org/10.3847/2515-5172/abd6e2)
- Yu, W., Richards, G., Buat, V., et al. 2022, *LSSTC AGN Data Challenge 2021*, 1.1, Zenodo, doi: [10.5281/zenodo.6878414](https://doi.org/10.5281/zenodo.6878414)
- Zebrun, K., Soszynski, I., Wozniak, P. R., et al. 2001, *AcA*, 51, 317. <https://arxiv.org/abs/astro-ph/0110623>
- Zel’dovich, Y. B., & Novikov, I. D. 1964, *Soviet Physics Doklady*, 9, 246
- Zu, Y., Kochanek, C. S., Kozłowski, S., & Udalski, A. 2013, *ApJ*, 765, 106, doi: [10.1088/0004-637X/765/2/106](https://doi.org/10.1088/0004-637X/765/2/106)