Laboratoire d'Étude du Rayonnement et de la Matière en Astrophysique

Paris, le 22 juillet 2024

Report on the PhD manuscript: *Identification and characterization of strong gravitational lenses and low surface brightmess galaxies using deep learning* by Hareesh Thuruthipilly

This manuscript presents the applications of a set of machine learning techniques to a series of important astrophysical problems. Chapter 1 introduces a basic overview of the current challenges in the current $\Lambda$CDM cosmological paradigm and two possible avenues to test and resolve the tension. Chapter 2 is a gentle and pedagogical introduction to deep learning and CNNs and present the transformers, a technique that will be applied in Chapters 3, 4 and 5 to actual datasets. Chapters 3 and 4 are published in refereed journals and Chapter 5 has been submitted. The manuscript concludes with a summary of the results and a brief perspective on follow-up projects.

The presentation is extremely good, from a clear introduction to the methods used and the astrophysical problems dealt with, even though the depth is somewhat irregular at times. For example, in Chapter 2, when discussing the positional encoding, the value of 12800 in Eqs. 2.7 and 2.8 should be justified and explained (p 31) and the explicit form of the softmax function provided as $\mathrm{softmax}(u_i) = \exp(u_i)/\sum_{j=1}^{N} \exp(u_j)$. These minor quibbles aside, Chapter two provides a very good overview which shows the deep knowledge of the subject and proposes the first-ever applications of the DL technique of transformers to astronomical images (although it has been also applied to astronomical time series by Allam, Peloton & McEwen in 2023).

Chapter 3 shows convincingly how the proposed technique using transformers ensembles (both detection transformers and vision transformers) can outperform the classic CNNs for strong gravitational lenses, a truly remarkable result which is well-explained, along with its limitations. A deeper discussion within a more general framework would have been welcome such as the actual use of these lenses to constrain cosmological

parameters as briefly touched upon in Chapter 1. The claim that $10^5$ spectra for strong gravitational lenses can be obtained should be discussed a bit further (p 18). This is actually the expected number of strong gravitational lenses from the images of the Euclid mission, not the spectra, and hence the question is whether and how the redshifts of both lenses and sources be actually measured in such a huge sample. Likewise the dispersion velocity of the lenses are difficult to observe and will limit the use of these samples (p 18). It would have been useful to frame the results in a Bayesian way, perhaps characterising the actual probabilities of detection given some prior distributions in the parameters, much in the way introduced by Sonnenfeld *et al.* (2023) *A&A*, 678, A4. Last but not least, given the critical role of the training sets, one could have envisioned using the Euclid-tailored Bologna sample to ground-based observations (*e.g.* with the LSST/Rubin) to further explore the role of seeing, PSF, etc beyond what has carried out with the KiDS sample.

Chapters 4 and 5 tackle an entirely different issue, namely the detection and characterisation of ultra-diffuse galaxies (UDGs) using the transformers techniques described in Chapter 2 and applied to strong gravitational lenses in Chapter 3. The method developped by the author, when applied to the sample extracted from the DES sample using the DECam instrument at CTIO, outperforms previous analyses and finds 17% more low surface brightness galaxies, a truly remarkable result given that the (biased) input sample is the same. Some of the artefacts defined by previous analyses are shown to be *bona-fide* UDGS which were misclassified. The newly-discovered galaxies are, mostly, red and within clusters of galaxies, and the (auto)correlation analysis confirms they are a population far more clustered than blue UDGs. A more sensitive probe of the differences in clustering could haven done using cross-correlations and would have measured their relative (linear) bias factors. It would also have been interesting to know what fraction of the newly-discovered UDGs are nucleated. A discussion on the selection biases introduced by the pre-selection in the imaging data would have been welcome beyond the statements in Section 4.2. For example, if the raw DES images would have been used, rather than the ones resulting from the standard pipeline (which was not optimised for LSB detection), presumably a much larger number of UDGs would have been detected simply because the volume explored in the $\mu_{eff}$ vs. $r_{eff}$ diagram would have expanded (this issue will come back in Chapter 5, where some fainter galaxies were missed, although the context is different). Alternatively, making further cuts to the training set could have been set up for perhaps mimicking part of the biases. All in all, this is a deeply original work carried out very carefully and keeping in mind the various selection effects.

Chapter 5 deals with the application of this original ensemble transformers technique to UDGs in the cluster of galaxies Abell 194. The author tackles here the key question of how to transfer learning from a DL model trained on a set to another, with an eye towards the deeper data that LSST/Rubin and Euclid will gather. The identification and recovery results are impressive given the widely different datasets in terms of depth. Two limitations, however, are that the effect of the different PSFs have not (and cannot) be included, even though the instruments are widely different and the way to deal with subtraction are also different,, and that the pre-DL analysis uses Sextractor, which is known to fail in the régime of very low surface brightness levels (beyond the need to mask or flag as contaminants some features). That being said, the results of this chapter, which also makes very interesting astrophysical inferences on the populations of UDGs in the this cluster, are truly remarkable.
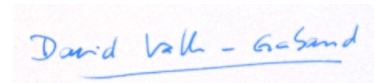
Minor points to be corrected:

1. Include table with the (long) list of acronyms used (for instance CNN is used on page *v* but not defined till much later).

2. Bing bang $\rightarrow$ Big Bang (p 3)

3. cephid $\rightarrow$ cepheid (p 4)

4. Figure 1.2: $n = 4$ (De Vaucouleurs' profile) would be more interesting to plot for its relevance, while values below $n = 0.5$ have never-observed unphysical depressions.

5. LSST/Rubin also observes in the $u$ band (p 10), as noted in Fig 1.5.

6. FIg $\rightarrow$ Fig (p 11)

7. It would be worth mentioning that dark energy could also change with direction, with anisotropic acceleration in the expansion (p 12)

8. area with 1 arcsec$^2$ region: this number is irrelevant for the argument (p 19)

9. In Eq. 2.6, $V$ should be included in the argument of the softmax function (p 31)

10. S´ersic $\rightarrow$ Sérsic (p 125)

11. keck $\rightarrow$ Keck (p 130)

12. Bibliographic references need to be completed, e.g.,

    - Allam & McEwen (2021) has been published in 2024 in *RAS Techniques and Instruments*.

    - Alzubaidi *et al.* (2021): page number (1) is missing

    - Chen *et al.* (2021): reference is incomplete

- de Block (2010): reference is incomplete

- Dietterich (2000): reference is incomplete

- E. Greene *et al.* (2022) → Greene, E. *et al.* (2022)

- Correct the spelling in the journal name in Einstein (1915) and Einstein (1917)

- Glorot & Bengio (2010a): reference is incomplete

- Ivezić *et al.* (2014): missing publisher's name (Cambridge Univ. Press)

- ditto for Mo *et al.* (2010): Cambridge Univ. Press

- ditto for Peebles (1980): Princeton Univ. Press

- Kigma & Ba (2015): *et al.* : reference is incomplete

- Popolo & Le Delliou (2017) should be Del Popolo & Le Delliou, and the review was published in Volume 5, id 17 of that journal.

- Sammut & Webb (2010): reference is incomplete

- Simonyan & Zisserman (2015): reference is incomplete

- Wang *et al.* (2022): reference is incomplete

- Wortsman *et al.* (2022): reference is incomplete

- Yosinski *et al.* (2014): reference is incomplete

The author should be commended to have placed in a public repository the codes created and used. This is essential for reproducible research and should be the norm.

It is recommended that the degree is awarded and the author congratulated for a deeply original research work and results.

Prof. David Valls–Gabaud
Directeur de recherche CNRS
Overseas Fellow, Churchill College, Univ. Cambridge
david.valls-gabaud@obspm.fr