

# Causal Machine Learning in Astrophysics and Beyond

Prof. Ofer Lahav (P&A), Dr. Nikos Nikolaou (P&A)

October 2023

Recent advancements in machine learning (ML) have enabled the efficient training of powerful statistical models from large amounts of high-dimensional data in various application domains, Astrophysics included [1-8]. Yet current learning systems are still almost exclusively operating on the level of statistical associations/correlations among the observed variables.

The next big step in the field should involve *causal modelling* [9, 10]; moving beyond simply capturing statistical associations to modelling cause-and-effect relationships among the underlying variables. The latest advancements in causal modelling are already finding practice in diverse fields such as healthcare & epidemiology, bioinformatics & pharmaceutical research, policymaking in social sciences, energy & climate, economics & finance -and more recently- in the physical sciences.

*Causal discovery* [11-13] aims to identify causal structure from data ('Which variables have a causal influence on variable A?') and *causal inference* to predict the results of *intervening* on variables ('What if I do X?') or -going a step further- of asking *counterfactual* questions ('What if I had done Y instead?') [14-16]. This project will be among the first to explore applications of causal ML algorithms in Astrophysics, particularly in the study of (i) galaxies and (ii) exoplanets.

- (i) In the exoplanetary literature, causal ML methods (in particular *half-sibling regression*) have so far been applied for decoupling observations from instrument systematics only in the context of exoplanet detection from transit light curves [17-19]. The project will apply and extend these methods to exoplanet characterization, i.e. inferring exoplanet atmospheric parameters from observed spectra [2-4].
- (ii) In extragalactic astronomy, causal ML methods have previously been applied by members of the project team Mucesh, Hartley, Lahav, in collaboration with Gilligan-Lee (Spotify) [20] to simulations (IllustrisTNG) for understanding the effect of environment on star formation in galaxies ('*nature vs. nurture*'). Key findings include that local density is found to be suppressing star formation at redshift  $z < 1$ , while the situation is reversed at higher redshift and that the mass of the halo is found to be a confounder. This project will explore application of these methodologies to real data from DES & DESI, Euclid & Rubin-LSST. We have access to all these data sets through Lahav's membership in these projects.

We propose to apply and extend the following subtasks for the CDT PhD project on causal inference (each to take ~1 year, but likely to be partially explored in parallel):

1. **Further to [17-19], to decouple observations for instrument systematics in exoplanet characterization<sup>1</sup>.**
2. **Further to [1], to explore and develop for real exoplanet data from JWST.**
3. **Further to [20], to apply to real galaxy data from DES & DESI<sup>2</sup>.**
4. **Further to [20], to apply to future galaxy data from Euclid & Rubin-LSST<sup>3</sup>.**

There is synergy between the two areas (exoplanets, galaxies). For example, when exploring subtasks (3 & 4), the methods developed in subtask (1) can be used to disentangle systematics from observations. This is an ambitious project with the potential to advance both the field of causal ML and the two subareas of Astrophysics.

Both supervisors have extensive experience in the intersection of ML & Astrophysics (e.g. [1-8]).

**Prof. Lahav** is Perren Professor of Astronomy and his research is in observational cosmology, using large galaxy surveys. He also co-directs CDT/DISI. He is currently supervising 4 CDT students on related projects (another student is about to submit his thesis; 2 others already defended their theses).

**Dr. Nikolaou** is a Lecturer at DISI and an expert in ML. He has worked on applications of ML in biomedicine & astronomy and his theoretical research interests include causal ML modelling [13]. He is currently supervising 2 CDT students (co-supervising another 2) & several MSc projects exploring ML applications in these and related areas.

---

<sup>1</sup> A good practice would be to first use a simulated dataset, in which the ground truth for both observations & systematics is known. This will allow us to validate the methods, before applying them to real data. A suggested initial dataset (already available) is the simulated one used for the Ariel Data Challenge: <https://zenodo.org/record/6770103#.Y2PmuuzP1qs>

<sup>2</sup> Data are already available and accessible to any student who joins.

<sup>3</sup> Data expected by the time the student reaches this stage.

## References:

- [1] Yip, K. H., et al. (2019) "Pushing the Limits of Exoplanet Discovery via Direct Imaging with Deep Learning". *ECML-PKDD 2019*
- [2] Yip, K. H., et al. "Peeking inside the Black Box: Interpreting Deep-learning Models for Exoplanet Atmospheric Retrievals". *The Astronomical Journal* 162.5 (2021): 195.
- [3] Morvan, M., et al. (2022). "Don't Pay Attention to the Noise: Learning Self-supervised Representations of Light Curves with a Denoising Time Series Transformer". In *Machine Learning for Astrophysics: Workshop at ICML 2022* (Vol. 162). *arXiv:2207.02777*
- [4] Nikolaou, N., et al. (2023). "Lessons Learned from the 1st ARIEL Machine Learning Challenge: Correcting Transiting Exoplanet Light Curves for Stellar Spots". *To Appear in RAS Techniques & Instruments*, *arXiv:2010.15996*
- [5] Morgan, R., et al. (2023). "DeepZipper. II. Searching for Lensed Supernovae in Dark Energy Survey Data with Deep Learning". *The Astrophysical Journal*, 943(1), 19.
- [6] Iess, A., et al. (2023). "LSTM and CNN application for core-collapse supernova search in gravitational wave real data". *arXiv preprint arXiv:2301.09387*.
- [7] Bhambra, P., et al. (2022). "Explaining deep learning of galaxy morphology with saliency mapping". *Monthly Notices of the Royal Astronomical Society*, 511(4), 5032-5041.
- [8] Henghes, B., et al. (2022). "Deep learning methods for obtaining photometric redshift estimations from images". *Monthly Notices of the Royal Astronomical Society*, 512(2), 1696-1709.
- [9] Pearl, J., et al. (2016). "Causal inference in statistics: A primer". John Wiley & Sons.
- [10] Peters, J., et al. (2017). "Elements of causal inference: foundations and learning algorithms" (p. 288). The MIT Press.
- [11] Heinze-Deml, C., et al. (2018). "Causal structure learning. *Annual Review of Statistics and Its Application*", 5, 371-391.
- [12] Glymour, C., et al. (2019). "Review of causal discovery methods based on graphical models". *Frontiers in genetics*, 10, 524.
- [13] Nikolaou, N., & Sechidis, K. (2020). "Inferring Causal Direction from Observational Data: A Complexity Approach". In *CoRR. PharML 2020*. *arXiv:2010.05635*
- [14] Spirtes, Peter. "Introduction to causal inference." *Journal of Machine Learning Research* 11, no. 5 (2010).
- [15] Prospero, M., et al. (2020). "Causal inference and counterfactual prediction in machine learning for actionable healthcare". *Nature Machine Intelligence*, 2(7), 369-375.
- [16] Yao, L., et al. (2021). "A survey on causal inference". *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15(5), 1-46.
- [17] Schölkopf, B., et al. (2015). "Removing systematic errors for exoplanet search via latent causes". In *International Conference on Machine Learning* (pp. 2218-2226). PMLR.
- [18] Wang, D. et al. (2016). "A causal, data-driven approach to modeling the Kepler data". *Publications of the Astronomical Society of the Pacific*, 128(967), 094503.
- [19] Schölkopf, B., et al. (2016). "Modelling confounding by half-sibling regression". *Proceedings of the National Academy of Sciences*, 113(27), 7391-7398.
- [20] Mucesh S., PhD thesis (2023) [in preparation], University College London [draft available upon request]