Ariel Machine Learning Data Challenges

September 2022





Horizon 2020



Ingo Waldmann





Science & Technology Facilities Council



European Research Council Established by the European Commission



State of ML in Exoplanets

- Still very new to exoplanets
- Mainly by three sub-fields: Transit detection, direct imaging, atmospheric modelling
- Number of serious ML/AI papers increased dramatically since 2018/19
- Many concurrent issues with other fields but little cross-talk yet

Search term

(full:"extrasolar planet" or full:"exoplanet") and (full:"machine learning" or full:"artificial intelligence" or full:"neural network")

stability multiplanet planet system global

> identify based signals transit deep

atmospheric transmission atmosphere science application

tess transit star young tic

imaging

plane

microlensing

components

high

Generate using Astrophysics Data System



A natural synergy with the AI community

<u>Astronomy</u>

- large data sets

- Complicated noise
- III posed problems

Machine learning

- Data detrending
- Pattern recognition
- Noise models

Three Machine Learning Data Challenges

ECML 2019 Würzburg







ECML 2021 Online

NeurIPS 2022 **New Orleans**





Accidentally entered pop culture in Portugal...



Can we model exoplanet atmospheres with AI?

- One AI that works for varying instruments, noise levels and systematics?





• What if we can train an AI to quickly and reliably classify and measure planet atmospheres?









Current ML retrieval progress in the field

- Radom Forrests (e.g. Marquez-Neila et al. 2018)
- Feed forward models (e.g. Waldmann 2016)
- Convolutional models (Ardevol Martinez et al. 2022)
- Generative models (Zingales & Waldmann 2018)
- Ensemble forward networks (e.g. Cobb et al 2019)
- Variational Inference (Yip et al. 2022
- Surrogate models (e.g. Himes et al. 2021)
- Transverse Vector Decomposition (Matchev et al. 2022)



Why a challenge on retrievals?

in contemporary ML

important to an adoption of these methods in science

parameter spaces -> we're asking for 6D

Generating approximate posterior distributions using machine learning over a wide parameter range is one of the hardest tasks

- **Understanding decision making** and biases is fundamentally
- Last year's NeurIPS Bayesian inference challenge featured 2D



Largest simulation set to date: 100,000 forward models - 26,000 retrievals



Frequency

Putting retrievals to the ML community



SciML AI Benchmark Dataset and continued development platform

- The dataset will be made into a standard AI benchmarking set
- and for future data workshops





Open Source Software for Scientific Machine Learning

We hope to maintain this data set / platform as development platform



Science & Technology Facilities Council Rutherford Appleton Laboratory



Ariel Competition Tracks







Ariel Competition Tracks

Light Track



- Provide summary statistics
- 500 planets for leader board
- 3000 planets for final solution
- 1st Prize: \$1000
- 2nd Prize: \$500



- Provide full posteriors
- 200 planets for leader board
- 3000 planets for final solution
- 1st Prize: \$2000
- 2nd Prize: \$500



Participating and Sponsoring partners

DiRAC Cesa SPACEFLUX





European Research Council Established by the European Commission



Open Source Software for Scientific Machine Learning



Science & Technology Facilities Council Rutherford Appleton Laboratory



CENTRE FOR SPACE EXOPLANET DATA



The Alan Turing Institute

FLATIRON INSTITUTE



Important beyond planets

- Bayesian Sampling of high parameter spaces is a general problem - Applicable in all complex inverse processes in physical science - Big problem in medical physics for example

PneumoniaMNIST



DermaMNIST

RetinaMNIST

Yang et al. 2021

Ariel ML Data Challenge

Start: 30th June **End: 17th October** NeurIPS 2022, 26 November - New Orleans

https://www.ariel-datachallenge.space

About The Science Leader Board Personal Score Download Information - Logout

