

NeurIPS – Ariel Data Challenge 2025: 4th Place Winning Solution

Thomas Dueholm Hansen

October 7, 2025

1 Team description

- Competition name: NeurIPS – Ariel Data Challenge 2025
- Team Name: Thomas Dueholm Hansen
- Public leaderboard score: 0.587
- Public leaderboard place: 6th
- Private leaderboard score: 0.590
- Private leaderboard place: 4th

The team consisted of a single team member:

- Name: Thomas Dueholm Hansen
- Location: ...
- Email: ...
- LinkedIn: <https://www.linkedin.com/in/thomas-dueholm-hansen-333634178/>

1.1 Academic and professional background

I received a PhD in computer science in 2012 from Aarhus University, Denmark. I then did one-year postdocs at Tel Aviv University and Stanford University, respectively, after which I returned to Aarhus University for a tenure track assistant professor position. In 2018 I moved to The University of Copenhagen, Denmark, for a position as associate professor of algorithms. My most notable academic achievement is a best paper award from the 43rd Annual ACM Symposium on Theory of Computing (STOC 2011), one of the two flagship conferences in theoretical computer science.

In 2017 and 2018 I traded actively with equities, long and short, index futures, and options. The associate professor position at The University of Copenhagen was offered permanently full-time, with a half-time leave of absence to be reevaluated yearly. Half of the time was to be spent on trading. My discretionary approach to trading did not work out, however, and I instead joined Nordea Markets, Denmark, in 2019 as a quantitative analyst. For about the first three years, I worked with fixed income, and in November 2021, I moved to FX spot as a quantitative trader.

I decided to take an extended sabbatical in 2023 in order to learn more about machine learning and to attempt to develop a profitable, automatic trading algorithm for FX. Since then I have worked on a mix of this private project and Kaggle competitions. The first time I did GPU programming was after leaving Nordea.

1.2 Prior experience

I have worked with time-series data in the above mentioned trading strategy project, as well as the following Kaggle competitions:

- *Child Mind Institute – Detect Sleep States*
- *Optiver – Trading at the Close*
- *HMS – Harmful Brain Activity Classification*
- *Jane Street Real-Time Market Data Forecasting*

In some sense I am starting to specialize in extracting information from noisy time series data, which was also what appealed to me about the *NeurIPS – Ariel Data Challenge 2025*.

I also made use of experience from the competition *ICR – Identifying Age-Related Conditions*. Specifically, that competition taught me what *not* to do when working with a small dataset.

All of these experiences have contributed to proficiency with the automatic differentiation library JAX, understanding of how to limit overfitting, and various other machine learning skills.

From my time at Nordea, I have some experience with the Kalman smoother from the statsmodels python package. In this project I implemented a Kalman smoother myself with JAX. The automatic initialization of parameters for this Kalman smoother was something I invented for this project, i.e., it was not something I had used previously.

1.3 Time spent

About six weeks were spent on the competition.

2 Summary

TODO.

3 Overall model architecture

4 Preprocessing

I generally follow the approach suggested by the organizers (see <https://www.kaggle.com/code/gordonyip/calibrating-and-binning-ariel-data>). I perform no binning, so the time series for FGS1 have length 67500, and the time series for AIRS have length 5625.

I treat hot and dead pixels as missing values. I replace such missing values with linear interpolation of neighboring pixels, four neighbors for FGS1 and two neighbors with the same wavelength for AIRS. Note that clusters of missing values remain missing. If one of the four most significant pixels are missing for AIRS, I remove that entire wavelength. The later model is set up to handle such missing values.

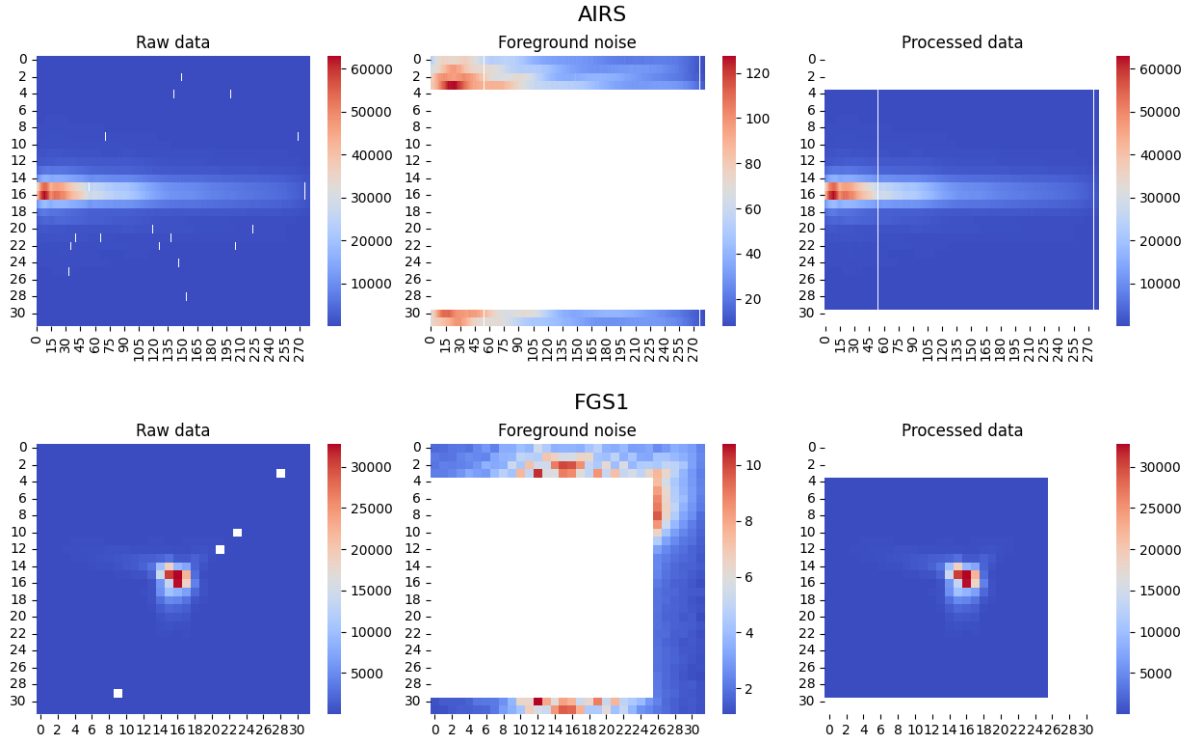


Figure 1: Foreground noise processing for planet 34983.

I use a rolling window of length 255 to discard outliers beyond five standard deviations from the centered rolling mean, again treating such values as missing.

Inspired by last year’s first place solution (see <https://www.kaggle.com/competitions/ariel-data-challenge-2024/writeups/c-number-daiwakun-1st-place-solution>), I use the rim of the image to estimate foreground noise, which I then subtract from the center of the image. Except, I do not explicitly use the rim as in last year’s first place solution. Instead I use the six rows with lowest intensity for AIRS, and the six columns and then six rows with lowest intensity for FGS1. Figure 1 illustrates the foreground noise processing step.

I was pressed on time and did not put effort into verifying that this form of preprocessing was beneficial. It seemed reasonable, so I spent my time elsewhere.

5 Initial depth and transit window estimation

My main model explicitly takes into account limb darkening, but it is non-trivial to optimize without ending up in an undesired local optimum. To improve the search, an initial solution is provided by a simpler model. This initial model was my baseline submission, and it is likely more complex than it needs to be when further refinement is performed later.

The idea for the initial model is to apply a Kalman Smoother to a Local Linear Trend model (see, e.g., https://www.statsmodels.org/dev/examples/notebooks/generated/statepace_local_linear_trend.html), to get a very smooth underlying level and trend from which the beginning and end of the transit can easily be extracted. The Kalman Smoother is designed to handle missing values, and the computation is performed twice: Once with all the data, and once with the data from the transit removed (set to NaN). The depth of the transit is extracted

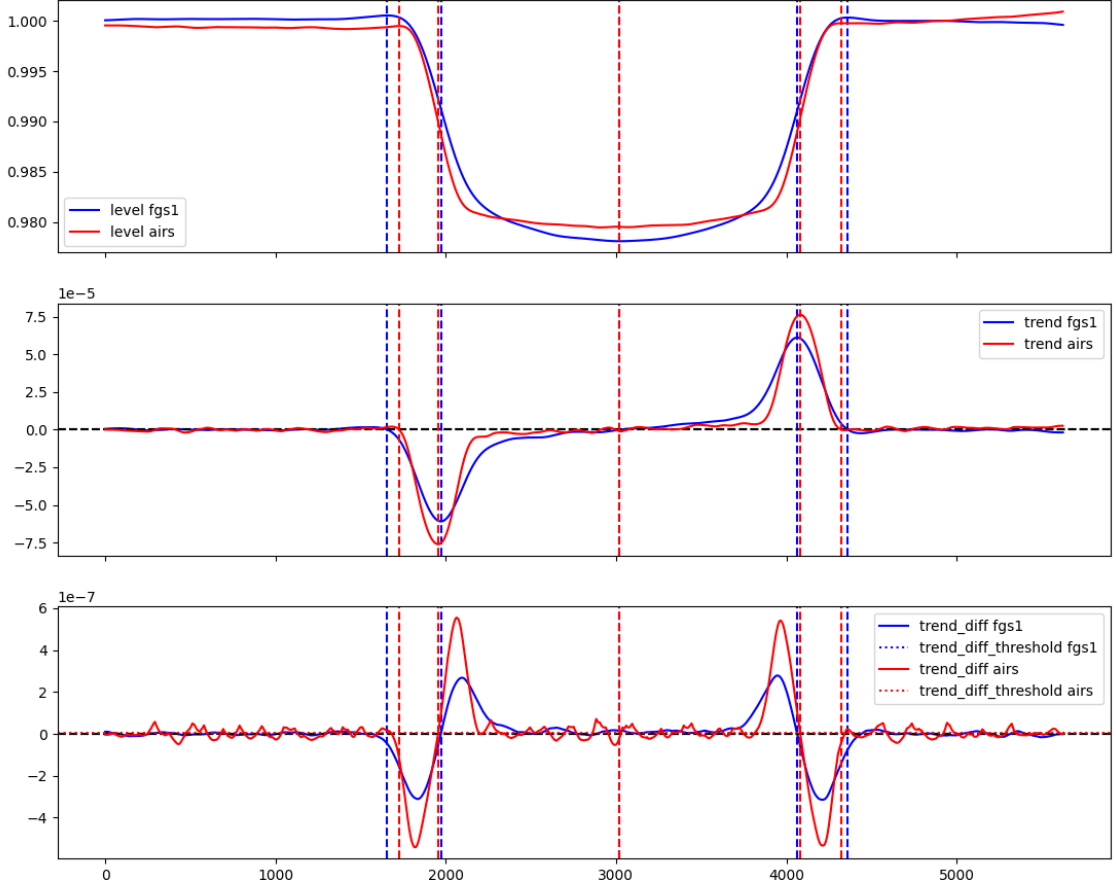


Figure 2: Transit identification for planet 34983.

from the difference of the levels resulting from these two computations.

The beginning, middle, and end of the transit are identified with the first Kalman Smoother. I identify as breakpoints the indices of the minimum and maximum values of the trend component. The middle of the transit is defined to be the average over these breakpoints, including both AIRS and FGS1. For AIRS, I use the sum across the entire spectrum to decrease noise. The beginning of the transit is the latest point prior to the first breakpoint where either the trend reaches zero, or the one-step difference of the trend is above that of the breakpoint. The end of the transit is defined similarly, but for the earliest point after the second breakpoint. The data from AIRS and FGS1 is combined by choosing the latest of the two beginnings and the earliest of the two ends. Transit identification is illustrated in Figure 2.

Robustness of transit identification relies on the produced curves being very smooth, and the Kalman Smoother is parameterized to produce such smooth curves. The Local Linear Trend model is defined by the linear system:

$$\begin{aligned}
 y_t &= \mu_t + \epsilon_t & \epsilon_t &\sim N(0, \sigma_\epsilon^2) \\
 \mu_{t+1} &= \mu_t + \nu_t + \xi_t & \xi_t &\sim N(0, \sigma_\xi^2) \\
 \nu_{t+1} &= \nu_t + \zeta_t & \zeta_t &\sim N(0, \sigma_\zeta^2)
 \end{aligned}$$

y_t is the observed data, μ_t is the level, and ν_t is the trend. The smoothness is determined by

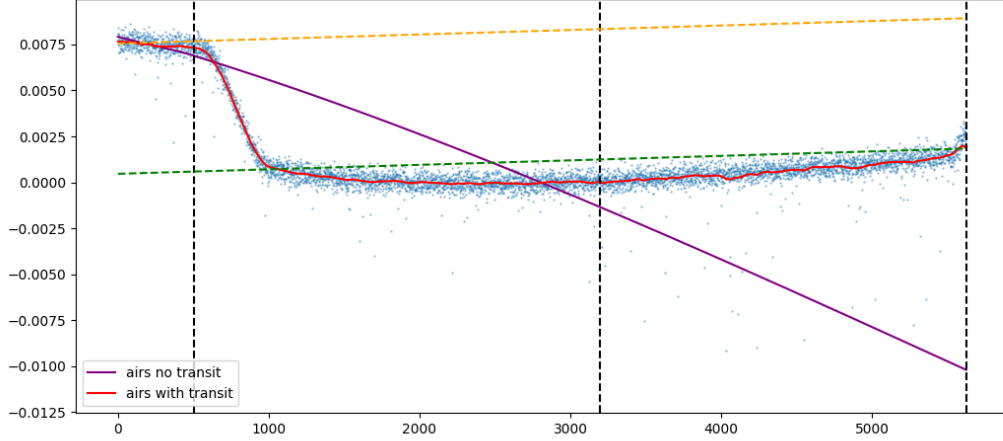


Figure 3: Example of incomplete transit for planet 1738121950.

the noise parameters σ_ϵ , σ_ξ , and σ_ζ . Instead of fitting such parameters to the data, I generate the parameters for a particular graph from three corresponding hyperparameters and a window length k . More precisely, I compute a moving average with the given window length k as a stand-in for the level, and a moving average of that as a stand-in for the trend. σ_ϵ , σ_ξ , and σ_ζ are products of the respective hyperparameters and the following quantities. For σ_ϵ , the variance around the stand-in level. For σ_ξ , the maximum absolute change in stand-in level. For σ_ζ , the maximum absolute change in stand-in trend.

Additionally, the Kalman Smoother requires μ_0 and ν_0 for initialization. Let y_1, y_2, \dots, y_N be the sequence of observed data. μ_0 and ν_0 are set to:

$$\mu_0 = \frac{1}{k} \sum_{t=1}^k y_t$$

$$\nu_0 = \left(\left(\frac{1}{k} \sum_{t=N-k+1}^N y_t \right) - \mu_0 \right) / (N - k)$$

Finally, the Kalman Smoother requires as input the covariance matrix Σ_0 for μ_0 and ν_0 . This is set to the steady state for σ_ϵ , σ_ξ , and σ_ζ , which is computed by solving Discrete Algebraic Riccati Equations (DARE). (See, e.g., section 2.2.1 of <https://www.robots.ox.ac.uk/~ian/Teaching/Estimation/LectureNotes2.pdf>.)

In some cases the beginning or the end of the transit falls outside the available data, and it becomes meaningless to apply a Kalman Smoother after removing the transit. In that case I replace the second Kalman Smoother by a line with slope extracted from data near the middle of the transit, and the highest intercept that allows the line to hit data either at the beginning or the end of the transit. An example of such an incomplete transit is shown in Figure 3. In this case, the second Kalman Smoother is replaced by the orange dashed line.

6 Depth estimation with explicit limb darkening

Outliers

6.1 Transit simulation

My model that explicitly takes into account limb darkening simulates a planet moving in front of a star, and estimates the amount of light observed at every time step. Limb darkening means that the intensity of the star is highest at the center and is reduced toward the edge (the limb). The formula that I use for limb darkening is based on <https://www.astro.uvic.ca/~tatum/stellatm/atm6.pdf>. The exact formula used is:

$$intensity(x, u, v) = 1 - u \cdot ((1 - v) \cdot (1 - \sqrt{1 - x^2}) + v \cdot x^2),$$

where x is the distance to the center of the star, assuming that the star has radius 1, and u and v are parameters that depend on the star. Unfortunately, I made a sign error in my code and changed the last plus into a minus, which meant that v behaved unexpectedly, so I gave up on directly optimizing v and instead treated it as a hyperparameter.

The model for constructing the transit uses six parameters:

- r : The radius of the planet, scaled such that the radius of the star is 1.
- u : The limb darkening parameter u described above.
- v : The limb darkening parameter v described above.
- d_{min} : The closest distance from the center of the planet to the center of the star during the transit.
- t_{min} : The time at which the planet is closest to the center of the star. The timeframe goes from 0 to 1.
- t : The transit starts when the planet first overlaps the star, and ends when the planet no longer overlaps the star. The transit length t is the time between these two points. Note that the transit length can be greater than 1 if the timeframe does not cover the entire transit.

The parameters are assumed to satisfy:

$$\begin{aligned} 0 &< r \\ 0 &\leq u \leq 1 \\ -1 &\leq v \leq 1 \\ 0 &\leq d_{min} < 1 + r \end{aligned}$$

Importantly, the transit is differentiable with respect to these six parameters, which means that the parameters can be optimized to produce the transit with the best fit for some training data.

Figure 4 shows examples of the transits generated by my model for the same parameters, except without and with limb darkening. Note that without limb darkening, the valley of the transit is very flat, because the observed light mainly depends on whether the planet is fully inside the boundary of the star or not. With limb darkening, the observed light continues being reduced as the planet gets closer to the center of the star, because it then covers an area with higher intensity.

6.2 Parameter optimization

initialization

Stages:

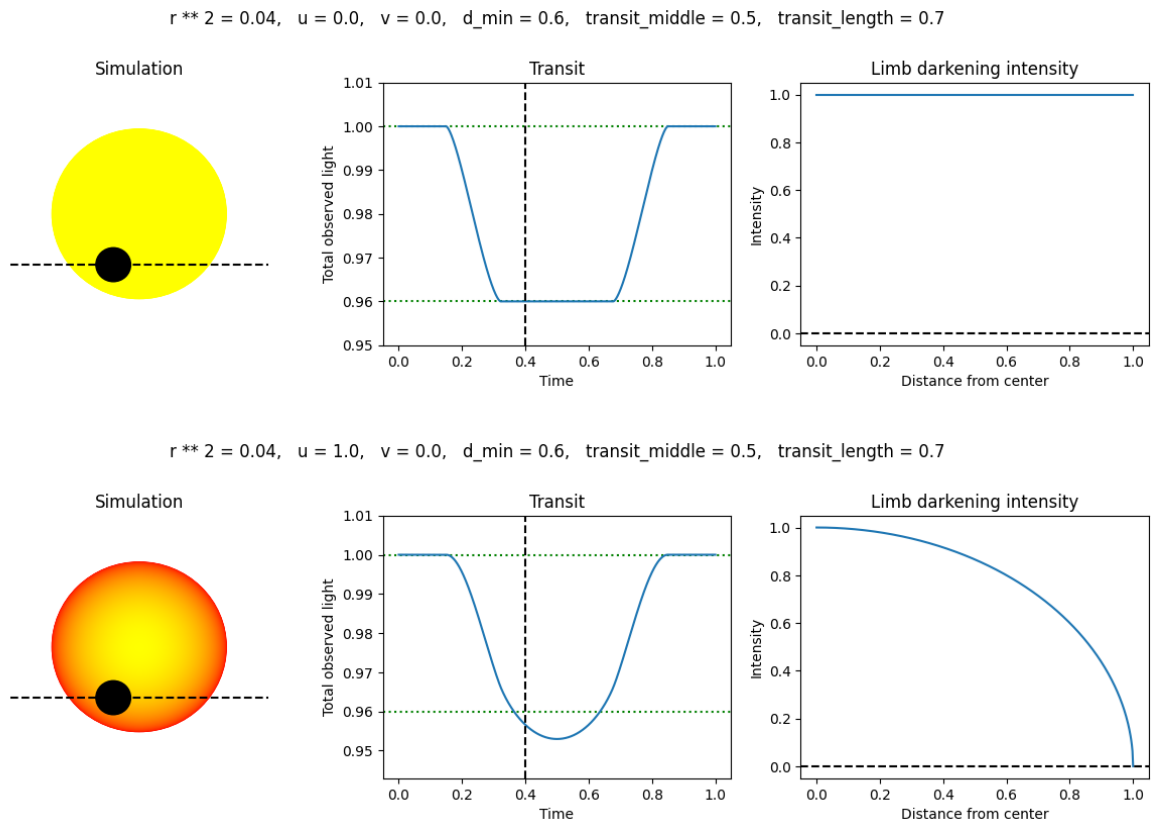


Figure 4: Example of a simulated transit without and with limb darkening.

7 Computation of uncertainty

The standard deviation component of the solution is computed with two independent Neural networks, one for AIRS and one for FGS1. I use 11 features for AIRS and 7 features for FGS1. The features are described below.

The neural networks are essentially multilayer perceptrons (MLPs), but with only a single hidden layer. Instead of using a single activation function, the hidden layer computes two vectors in parallel, applies a sigmoid activation function to one and relu to the other, and takes the product of the two vectors before applying dropout. I also use batch normalization for both vectors. The output is transformed with one plus an elu activation function and scaled by a parameter. The code is simple enough that it is worth showing in its entirety in Figure 5.

The inputs for the MLPs are normalized with a StandardScaler, i.e., I subtract the mean and divide by the standard deviation. The inputs are clipped by the MLP itself to simplify hyperparameter optimization. Hyperparameter optimization was performed with manual experimentation. I initially planned to use Optuna, but I found satisfactory parameters before I got around to it. I used five fold cross validation with three repetitions. I selected features in the same way.

For AIRS, the features are clipped at plus-minus 7 standard deviations, the hidden dimension and dropout are 256 and 0.4, and the log_scale parameter is -9. For FGS1, the features are clipped at plus-minus 5 standard deviations, the hidden dimension and dropout are 32 and 0.1, and the log_scale parameter is -8. In both cases I use an Adam optimizer with weight decay. For AIRS I trained for 5 epochs, and for FGS1 I trained for 150 epochs. Note that there was much more training data for AIRS due to working with a spectrum, which explains why it worked better with fewer epochs.

I use 11 features for AIRS, the most important feature being the limb darkening coefficient u . The features are based on the output of the limb darkening model described in Section 6. Instead of directly using the output of the depth estimation model, I use five fold cross validation with three repetitions to generate three times as much training data with out-of-sample properties. The features are as follows:

1. The predicted value of r^2 for the planet. If there are multiple instances of the planet, then the mean of such predictions.
2. The absolute value of the difference between feature 1 and the predicted r^2 for the particular instance, i.e., if there is only one instance for a planet, then this feature is zero.
3. The predicted value of the limb darkening coefficient u .
4. The predicted value of the transit length t .
5. The predicted minimum distance d_{min} .
6. The multiplicity of the planet, i.e., how many instances of data that are available.
7. The standard deviation of the predictions for r^2 for the spectrum.
8. The index of the data point in the spectrum.
9. For a particular index in the spectrum, the standard deviation of the predictions for r^2 for that index, but ignoring outliers.
10. A measure for the non-linearity of the polynomials for the gain drift: The sum over the outer product of the absolute values of coefficients from f with degree at least 2 and the

```

import flax.linen as nn
import jax.numpy as jnp
from typing import Tuple

class MLP(nn.Module):

    input_clipping_threshold: float
    n_hidden: Tuple[int]
    dropout: Tuple[float]
    log_scale: float

    @nn.compact
    def __call__(self, x, train: bool):
        x = x.astype(jnp.float64)
        x = jnp.clip(x, -self.input_clipping_threshold, self.input_clipping_threshold)

        for n_hidden, dropout in zip(self.n_hidden, self.dropout):
            x_sigmoid = nn.BatchNorm(use_running_average=(not train))(x)
            x_sigmoid = nn.Dense(features=n_hidden, use_bias=True)(x_sigmoid)
            x_sigmoid = nn.sigmoid(x_sigmoid)

            x = nn.Dense(features=n_hidden, use_bias=True)(x)
            x = nn.BatchNorm(use_running_average=(not train))(x)
            x = nn.relu(x)

            x = x * x_sigmoid

        x = nn.Dropout(rate=dropout, deterministic=(not train))(x)

        x = nn.Dense(features=1, use_bias=True)(x)[..., 0]
        x = 1 + nn.elu(x)

        log_scale = self.param('log_scale no_weight_decay',
                                nn.initializers.constant(self.log_scale), (), jnp.float32)

        x = jnp.exp(log_scale) * x

    return x

```

Figure 5: Modified multilayer perceptron.

absolute values of coefficients from g with degree at least 1, normalized by the mean of the raw signal.

11. The predicted spectrum: The difference between the predicted value of r^2 for individual indices and the mean over the entire spectrum.

For FGS1 I use the same features as for AIRS, computed with data from FGS1, but I only use the 7 first features. Note that the 7th feature describes the standard deviation of the predicted spectrum, which is computed independently of FGS1, but it was still helpful to include it as a feature for FGS1.

8 Principal component analysis (PCA)

9 Ensembling

Both with and without PCA I train five models and average the result, so in total I train ten models with different seeds. This number could have been much larger while staying within the time limit, but I did not experiment with this aspect.

10 Computation time

I processed data and trained models locally with an RTX 5090. In total, the running time of a Kaggle notebook with a P100 was about twice as long as with my RTX 5090. I expected a higher speedup with a significantly faster GPU, but the code has not been properly optimized (the electricity consumption was very low), and calculations are generally performed with 64 bit precision. The time spent on the individual steps with an RTX 5090 was:

- Preprocessing (Section 4): 17 minutes.
- Initial depth estimation (Section 5): 1 minute.
- Depth estimation with limb darkening (Section 6): 3 hours and 35 minutes.
- Model training: 11 minutes.

11 Reproducibility

During the competition, I noticed that the standard deviations produced by Kaggle notebooks differed from those produced locally. I thought this was due to using different batch sizes, but even with the same batch sizes there was a difference. Since the end of the competition, there was an automatic update of Ubuntu that broke my Nvidia drivers, and I had to reinstall the drivers. Now there is a slight difference between the standard deviations produced locally by the same code before and after the driver update, so the issue with Kaggle notebooks may also have been driver related.

I was careful to keep track of code and seed changes, since I had trouble with reproducibility in the *Jane Street Real-Time Market Data Forecasting* competition, but I did not keep track of drivers. The new driver is nvidia-driver-580-open. I believe my previous driver was nvidia-driver-575-open, which is no longer available via ubuntu-drivers. I also did a full purge of Nvidia software on my system in general. All in all, I decided that it was not worth trying to replicate my previous setup.

I removed some unused code for the final version, but I verified that this did not change the produced standard deviations.

12 Conclusion

...

13 Are we alone in the universe?

The Fermi Paradox asks "Why are there no signs of intelligence elsewhere in the universe?" The opportunities for life to appear are literally astronomical. I would argue, however, that astronomical numbers are quite small in this context, since they are only polynomial in nature. The probability of intelligent life appearing may very well be exponential in nature, for example if an initial critical step is for a sequence of genes to randomly come together in a specific order. Even with slower growing functions than exponentials, the inverse of the probability of intelligent life may dwarf any number in the astronomical range. Thus, it may be plausible that the answer to the question "Are we alone in an infinite universe?" is no, while the answer to "Are we alone in the observable universe?" is yes with very high probability. This all assumes that life appears spontaneously. In a distant future it seems plausible that humanity will seed other planets with life. Perhaps life on Earth itself originates from such seeds planted by other civilizations or higher entities. Just something to keep in mind: If life was seeded by ancient civilizations, how should that affect the way that we search for it?