

# Using Hybrid Machine Learning for Fundamental Stellar Parameter Inference

SUJAY SHANKAR,<sup>1</sup> MICHAEL GULLY-SANTIAGO,<sup>1</sup> AND CAROLINE V. MORLEY<sup>1</sup>

<sup>1</sup>*Department of Astronomy, The University of Texas at Austin, Austin, TX 78712, USA*

## ABSTRACT

Text

*Keywords:* Text

### 1. INTRODUCTION

Text

### 2. CLONING THE PHOENIX MODEL GRID

#### 2.1. The PHOENIX Subset

For the purposes of this study, we did not consider the full range of the PHOENIX synthetic spectral model grid. Instead, we focused on a subset of the grid, whose parameter ranges are given in Table 1.

PHOENIX Grid Subset		
Parameter	Symbol	Range
Alpha Element Abundance	$\alpha$	[0] dex
Iron Abundance	[Fe/H]	[-0.5, 0.5] dex
Effective Temperature	$T_{\text{eff}}$	[2300, 12000] K
Surface Gravity	$\log(g)$	[0, 5]
Wavelength	$\lambda$	[8050, 12850] Å

**Table 1.** The subset of the PHOENIX grid used in this study. These limits were imposed to reduce the computational cost of the algorithms. The wavelength limits in particular roughly line up with that of the Habitable Zone Planet Finder (HPF) spectrograph.

#### 2.2. Preprocessing with *gollum*

First, the PHOENIX subset was accessed directly from the PHOENIX website using *gollum*’s download option, which uses the website’s FTP server to download spectra without using any disk space. This came at the cost of longer load times due to internet connection being significantly more of a bottleneck than disk I/O, but was acceptable as this was a one-time operation.

The spectra were then put through a three-step preprocessing pipeline.

1. Blackbody Division: Since the  $T_{\text{eff}}$  of each spectrum is known, the according blackbody spectrum was divided out.
2. Percentile Normalization: The spectra were normalized by dividing them by their 99th percentile.
3. Continuum Normalization: The spectra were further normalized by dividing them by a 5<sup>th</sup> order polynomial continuum.

#### 2.3. Line Identification with *blase*

The next step was to convert the PHOENIX subset into an interpretable format. We wanted to represent the spectra as a list of spectral lines rather than a list of fluxes. This was done using *blase*, which detects spectral lines as Voigt profiles and tunes the profiles to mimic the original PHOENIX spectrum with back propagation. Four parameters were optimized: the line center  $\mu$ , the log-amplitude  $\ln(a)$ , the Gaussian width  $\sigma$ , and the Lorentzian width  $\gamma$ . The optimization used the Adam optimizer with a learning rate of 0.05 over 100 epochs. In addition, we limited two custom parameters: wing cut to 6000 and prominence to 0.005. Wing cut is a parameter that determines the extent of the Voigt profile to evaluate, saving computational resources by not evaluating small numbers. Prominence sets a lower limit for the amplitude of detected lines, which also saves resources by disregarding shallow lines. In short, larger values for wing cut and smaller values for prominence both increase the accuracy of using *blase*’s line detection and optimization at the expense of an increased computational cost.

Once optimization was complete, the list of identified lines, present in the state dictionary of *blase*’s model, was saved to a *.pt* file.

### 3. INTERPOLATING MANIFOLDS

#### 3.1. Cross-Model Line Identification

Earlier we mentioned that **blase** tuned the line centers of detected lines. This means that from one PHOENIX spectrum to the next, the same line could have a slightly different line center. Since the goal of this study is to interpolate the properties of each line, we needed to identify the presence of a particular line across the PHOENIX subset. We decided to do this by using the line centers of the detected lines pre-optimization. This introduced a new variable called the jitter  $\delta\mu$ , which measures the shift of the line center between pre-optimization and post-optimization. Now with each spectral line indexed by  $\mu$ , we had four parameters to interpolate:  $\delta\mu$ ,  $\ln(a)$ ,  $\sigma$ , and  $\gamma$ .

There was also a second issue, and that is that spectral lines were often only detected in some spectra from the PHOENIX subset. This means that in theory, interpolation would be inaccurate in regions where a line does not appear, and also breaks in practice because **scipy**'s regular grid interpolation relies on a rectilinear grid, which we would not have if regions of grid points were missing. To solve this, we artificially populated missing grid points with log-amplitudes of -1000, which retained interpolator stability by not being an infinity, but also essentially nullified the line when evaluated.

### 3.2. Continuously Evaluable Manifolds

For each line, the inputs to the interpolator were the three input parameters  $T_{\text{eff}}$ ,  $\log(g)$ , and  $[Fe/H]$ , and the output was a list of four parameters,  $\delta\mu$ ,  $\ln(a)$ ,  $\sigma$ , and  $\gamma$ . For each line, one of these interpolator objects was created using linear interpolation, and these interpolators were aggregated into a single list, which was when pickled and written to a **.pkl** file. These interpolators represent multiple manifolds mapping fundamental stellar parameters to spectral line properties. These interpolators could now be evaluated at any point lying within the domain of the PHOENIX subset, turning a discretely sampled PHOENIX subset into a continuous one. With the given size of the PHOENIX subset, the interpolator list takes up 19 GB of disk space.

## 4. BAYESIAN OPTIMIZATION

Text

## 5. INFERENCE TESTING

### 5.1. Test Data

Text

### 5.2. Performance Results

Text

## 6. DISCUSSION

It should be noted that this study represents a proof of concept, and that there are numerous design considerations that could be improved upon with future work. These considerations include but are not limited to the following:

- *Limited PHOENIX Subset*: The PHOENIX subset used in this study was just that, a subset. The full PHOENIX grid not only expands the  $[Fe/H]$  range to  $[-4.0, 1.0]$  dex, but also includes the alpha element abundance, which we elected to fix at 0 for this study. In addition to the actual fundamental stellar parameters, we also took a subset of the PHOENIX wavelength range, with the full  $[500, 55000]$  Å wavelength range also being left to future work.
- *Strict Wavelength Range*: Currently, the generator only supports inference on spectra whose wavelength limits are either equal to it, or encompass that of the generator and have been truncated to match. However, when the spectrum in question has a smaller wavelength range than the generator, currently there is no functionality to truncate the generator. This would require externally indexing the generator's individual interpolators by line center position and selectively evaluating those to eliminate wasteful computation.
- *Single Model Grid*: The PHOENIX grid is not the only model grid of synthetic spectra available, and it does not apply to all types of stars. Future work would extend the reach of this study's algorithm to encompass other model grids such as the Sonora series for substellar models. **blase** should be able to have an option for the user to input which model grid they would like to base the inference on, and to get even more advanced, perhaps even have the ability to intelligently determine which model grid to use automatically.
- *Memorization vs. Generalization*: The current design of the algorithm constructs manifolds using interpolation. This means that performance is good at points close to PHOENIX subset grid points, but is highly dependent on the type of interpolation used. As interpolators require memorization of the data, advanced interpolation becomes extremely expensive in terms of disk utilization. Future work would involve constructing manifolds using regression, which would allow for much better generalization and lower disk utilization at the expense of some accuracy.

- *Extrinsic Absence:* The current design of our algorithm does not account for extrinsic parameters that modify the appearance of spectra such as rotational broadening and doppler shifting. Future work would need to develop ways to tune these extrinsic parameters alongside the fundamental stellar parameters.
- *Framework Overhead:* As this algorithm is currently more proof of concept than practical, it uses convenience functions from various libraries, which naturally introduces some level of overhead and leaves performance on the table. Future work would involve writing custom functions expressly designed for **blase**, most likely a complete rewrite of the library from the ground up.
- *Pseudo-Interpretability:* Our algorithm boasts interpretability by considering spectral lines as the objects of interest as opposed to the rather uninterpretable flux values of other approaches. However, this is only a step in the direction of interpretability. True interpretability would decompose a spectrum not into a set of spectral lines, but into a set of species component spectra, which requires a much more advanced understanding of different species and their behavior, as well as direct access to a radiative transfer code as opposed to an off-the-shelf model grid. This approach would also extend the inference from just fundamental stellar parameters defined by a grid to any set of parameters accounted for in the radiative transfer model, down to specific species abundances.
- *The Continuum Black Box:* Continuum normalization is a process that is not yet completely understood, and is currently done as a preprocessing step with a fairly simple algorithm. Future work would dive deeper into the science of continuums and develop more advanced methods that can discern continuums with greater accuracy and less modeling restrictions.
- *One Voigt Fits All:* The current assumption of **blase** is that every spectral line is a Voigt profile. This assumption is largely true, but there are situations where that is simply not enough. Future studies need to account for more advanced spectral line profiles and procedures to deal with phenomena such as ro-vibrational bands.

In short, it is quite clear that **blase** is a long way from being a powerhouse in spectral inference, however it represents a step down a road not yet traveled, and the

potential for future growth is immense. At the end of the day, we aim to develop a tool that can analyze spectrum based on how nature behaves rather than how data behaves.

## ACKNOWLEDGEMENTS

Text