

# Modelling stars with Gaussian Process Regression: Augmenting Stellar Model Grid

Tanda Li<sup>1</sup>\*, Guy R. Davies<sup>1</sup>†, Alex Lyttle<sup>1</sup>, Warrick Ball<sup>1</sup>, Lindsey Carboneau<sup>1</sup>, Rafael A. García<sup>2</sup>,

<sup>1</sup> School of Physics and Astronomy, University of Birmingham, Birmingham, B15 2TT, United Kingdom

<sup>2</sup> AIM, CEA, CNRS, Université Paris-Saclay, Université Paris Diderot, Sorbonne Paris Cité, F-91191 Gif-sur-Yvette, France

Accepted XXX. Received YYY; in original form ZZZ

## ABSTRACT

Grid-based modelling is widely used for estimating stellar parameters. However, stellar model grid is sparse because of the computational cost. This paper demonstrates an application of a machine-learning algorithm using the Gaussian Process (GP) Regression that turns a sparse model grid onto a continuous function. We train GP models to map five fundamental inputs (mass, equivalent evolutionary phase, initial metallicity, initial helium fraction, and the mixing-length parameter) to observable outputs (effective temperature, surface gravity, radius, surface metallicity, and stellar age). We test the GP predictions for the five outputs using off-grid stellar models and find no obvious systematic offsets, indicating good accuracy in predictions. As a further validation, we apply these GP models to characterise 1,000 fake stars. Inferred masses and ages determined with GP models well recover true values within one standard deviation. An important consequence of using GP-based interpolation is that stellar ages are more precise than those estimated with the original sparse grid because of the full sampling of fundamental inputs.

**Key words:** Star: Modelling – Machine Learning – Gaussian Process

## 1 INTRODUCTION

Theoretical stellar model has been developed for decades to simulate star structure and evolution. Star modelling is mostly grid-based (e.g. Choi et al. 2016) because computing many stellar models are time-consuming especially when a number of free input parameters are considered. Varying one of these adjusted parameters (mass, metallicity, helium fraction, mixing-length parameter, etc.) adds on an input demission and hence exponentially increases the computational cost.

A sparse grid is not ideal for the statistics analysis. Classical method like interpolation has been applied to overcome this disadvantage. For instance, Dotter (2016) developed a method to transform stellar evolution tracks onto a uniform basis and then interpolate to construct stellar isochrones. More recently, Rendle et al. (2019) uses Bayesian statistics and a Markov Chain Monte Carlo approach to find a representative set of interpolated models from a grid. The interpolation of both works achieve good accuracy for 3-demission girds (inputs are mass, age, and metallicity). However, this approach becomes less reliable for high-demission grid.

Machine learning is being applied to the field of stellar research in many ways to efficiently characterise stars. Verma et al. (2016)

applied artificial neural network, which is a series of algorithms that endeavours to recognise underlying relationships in a set of data, to determine the evolutionary parameters of the sun and sun-like stars based on spectroscopic and seismic measurements. Using a similar artificial neural network interference, Hendriks & Aerts (2019) developed a method to provide the optimal starting point of model competitions for more detailed forward asteroseismic modelling. Moreover, Mombarg et al. (2021) trained neural networks to predict theoretical pulsation periods of high-order gravity modes, as well as the luminosity, effective temperature, and surface gravity for a given mass, age, overshooting parameter, diffusive envelope mixing, metallicity, and near-core rotation frequency. Using different machine learning tools, Bellinger et al. (2016) trained a random forest regressors (Ho 1995), which is an ensemble learning method for regression and operates by constructing a multitude of decision trees, to rapidly estimate fundamental parameters of solar-like stars based on classical and asteroseismic observations. Hon et al. (2018) developed a convolutional neural network classifier that analyses visual features in asteroseismic frequency spectra to distinguish between red giant branch stars and helium-core burning stars. Wu et al. (2019) determined masses and ages for massive RGB stars from their spectra with a machine-learning method based on kernel principal component analysis, which is a nonlinear form of principal component analysis using integral operator kernel functions and can efficiently compute principal components in high dimen-

\* E-mail: t.li.2@bham.ac.uk

† E-mail: G.R.Davies@bham.ac.uk

<sup>47</sup> sional feature spaces related to input space by some nonlinear map  
<sup>48</sup> (Schölkopf et al. 1997). Hon et al. (2020) applied the mixture density network (Bishop 1994), which learns a transformation from a  
<sup>49</sup> set of input variables to a set of output variable, to determine stars’  
<sup>50</sup> fundamental parameters like mass and age based on observed mode  
<sup>51</sup> frequencies, spectroscopic, and global seismic parameters.  
<sup>52</sup>

In above studies, the discriminative machine-learning model is mostly used. The discriminative model treats observables as given facts to directly infer star fundamental parameters. The method is efficient and easy for computation, while the downside is not allowing any priors for star properties like mass. In an opposite direction, the generative machine-learning model uses the star fundamental parameters as given facts to predict observables. This approach offers flexibility to prior fundamental parameters in the sampling. For instance, Lytle et al. (2021) determined initial helium fraction and mixing-length parameters for a sample of *Kepler* dwarfs and subgiants with an artificial neural network to provide the generative model. This allowed them to prescribe prior distributions over the fundamental stellar parameters and, by extension, over population-level parameters such as a helium enrichment law. Priors encode our current knowledge and assumptions into inference from new data. This is especially important with noisy observations which span a large portion of parameter-space.

Constructing a comprehensive and fine model grid is computationally expensive. In this work, we aim to apply the machine learning tool to transform a sparse model grid onto a continuous function. We apply a machine learning algorithm that involves a Gaussian process (GP) that measures the similarity between data points (i.e., the kernel function) to predict values for unseen points from training data. We use the generative model and treat fundamental parameters as given facts to predict observables. This gives us flexibility to prior fundamental inputs when modelling stars. We organise the rest of the paper as follow. Section 2 contents descriptions about the computation of a representative stellar model grid. We then introduce the underline theory of GP and the setup of GP model in Section 3. We then demonstrate some preliminary studies for low-demission problems in Section 4. Section 5 demonstrates GP predictions and their systematic uncertainties. Subsequently, we augment the grid to have a set of continuously-sampled stellar models and model 100 fake stars for testing the accuracy of our method in Section 6. Lastly, we discuss advantages and limitations of this approach, highlight areas where improvements can be found in the near future, and summary conclusions in Section 7.

## 90 2 REPRESENTATIVE MODEL GRID

### 91 2.1 Grid computation

We compute a stellar model grid as the training dataset. We aim to cover stars with approximate solar mass on the main-sequence and the subgiant phases. We consider four independent fundamental inputs which are stellar mass ( $M$ ), initial helium fraction ( $Y_{\text{init}}$ ), initial metallicity ( $[\text{Fe}/\text{H}]_{\text{init}}$ ), and the mixing-length parameter ( $\alpha_{\text{MLT}}$ ). We calculated three model dataset for different purposes. The primary dataset is a standard model grid with uniform mass step. This model grid is used for all preliminary tests and also for the final training. We also calculate an additional dataset to increase the grid resolution for  $M > 1.05M_{\odot}$ . Because we find that the blue hook feature (where global parameters sharply vary) is relatively hard to train. This dataset is only used in the final training. Details of parameter ranges and steps of the two grids are listed in Table 1. For

**Table 1.** Computation of Stellar model grid.

Primary dataset		
Input Parameter	Range	Increment
$M (M_{\odot})$	0.80 – 1.20	0.01
$[\text{Fe}/\text{H}] (\text{dex})$	-0.5 – 0.2/0.2 – 0.5	0.1/0.05
$Y_{\text{init}}$	0.24 – 0.32	0.02
$\alpha_{\text{MLT}}$	1.7 – 2.5	0.2
Additional dataset		
Input Parameter	Range	Increment
$M (M_{\odot})$	1.055 – 1.195	0.01
$[\text{Fe}/\text{H}] (\text{dex})$	0.25 – 0.45	0.1
$Y_{\text{init}}$	0.25 – 0.31	0.02
$\alpha_{\text{MLT}}$	1.8 – 2.4	0.2

<sup>105</sup> validating and testing GP predictions, we computed off-grid models  
<sup>106</sup> with randomly sampled fundamental inputs as a third dataset. The  
<sup>107</sup> computation of evolutionary tracks starts at the Hayashi line with  
<sup>108</sup> pre-main-sequence central temperature at 300,000K and terminates  
<sup>109</sup> at the base of red-giant branch (RGB) where  $\log g = 3.6\text{dex}$ . Note  
<sup>110</sup> that we only use models after the zero-age-main-sequence (ZAMS),  
<sup>111</sup> which is defined as the point where core-hydrogen burning con-  
<sup>112</sup>tributes over 99.9% of the total luminosity.

### 113 2.2 Input physics

We use the stellar code Modules for Experiments in Stellar Astrophysics (MESA, version 12115) to construct stellar grids. MESA is an open-source stellar evolution package which is undergoing active development. Descriptions of input physics and numerical methods can be found in Paxton et al. (2011, 2013, 2015). We adopted the solar chemical mixture ( $(Z/X)_{\odot} = 0.0181$ ) provided by Asplund et al. (2009). The initial helium fraction ( $Y_{\text{init}}$ ) and initial metallicity ( $[\text{Fe}/\text{H}]_{\text{init}}$ ) are both independent inputs.

We use the MESA  $\rho - T$  tables based on the 2005 update of OPAL EOS tables (Rogers & Nayfonov 2002) and OPAL opacity supplemented by low-temperature opacity (Ferguson et al. 2005). The grey Eddington  $T - \tau$  relation is used to determine boundary conditions for modelling the atmosphere. The mixing-length theory is implemented and the convection is adjusted by the mixing-length parameter ( $\alpha_{\text{MLT}}$ ). We also apply the MESA convective premixing scheme (Paxton et al. 2019), which an approach to handling mixing in convection zones that improves model structures at the convective boundary. Atomic diffusion of helium and heavy elements was also taken into account. MESA calculates particle diffusion and gravitational settling by solving Burger’s equations using the method and diffusion coefficients of Thoul et al. (1994) as well as radiative turbulence formula given by Morel & Thévenin (2002). We consider eight elements ( $^1\text{H}$ ,  $^3\text{He}$ ,  $^4\text{He}$ ,  $^{12}\text{C}$ ,  $^{14}\text{N}$ ,  $^{16}\text{O}$ ,  $^{20}\text{Ne}$ , and  $^{24}\text{Mg}$ ) for diffusion calculations, and have the charge calculated by the MESA ionization module, which estimates the typical ionic charge as a function of  $T$ ,  $\rho$ , and free electrons per nucleon from Paquette et al. (1986). We only compute diffusion during the main-sequence stage before the central hydrogen abundance drops below 0.05, because its effects can be neglected in post main-sequence stages. The MESA inlist used for the computation is available on [https://github.com/litanda/mesa\\_inlist/](https://github.com/litanda/mesa_inlist/).

145 **2.3 Equivalent Evolutionary Phase**

146 Apart from the four independent model inputs, i.e., mass, metallicity,  
 147 helium fraction, and the mixing-length parameter, stellar age is  
 148 the fifth fundamental input. However, the dynamical range of age  
 149 varies track by track. This makes GP models hard to map from age  
 150 to global parameters. We need a uniform input to replace the age.  
 151 The fractional age is an option but we find that global parameters  
 152 (e.g. effective temperature) sharply change with the fractional age  
 153 around the blue hook and the turn-off point (as shown in the left  
 154 panel in Figure 1). It requires a complex and spiky kernel function  
 155 to fit the curvatures in this area and hence difficult for GP to learn.  
 156 Dotter (2016) has introduced a quantity Equivalent Evolutionary  
 157 Phase (*EEP*), which numbers evolutionary stages and transform  
 158 stellar tracks onto a uniform basis. We follow this idea but define  
 159 *EEP* in a different way to make global parameters change relatively  
 160 smoothly. On each evolutionary track, we compute the displace-  
 161 ment between consecutive models on the  $\log T_{\text{eff}} - \log g$  diagram.  
 162 For instance, the displacement between model  $n$  and model  $n - 1$   
 163 can be calculated as

$$\delta d_n = ((\log T_{\text{eff},n} - \log T_{\text{eff},n-1})^2 + (\log g_n - \log g_{n-1})^2)^c, \quad (1)$$

164 where  $c$  is an adjusted parameter to scale the displacement. The  
 165 total displacement of model  $n$  from the ZAMS (model 0) can be  
 166 calculated with

$$d_n = \sum_{i=1}^{i=n} \delta d_i. \quad (2)$$

167 We then normalise  $d_n$  to the 0 – 1 range and define it as *EEP*.  
 168 On an evolutionary track, *EEP* equals to 0 at the ZAMS and 1 on  
 169 the RGB where  $\log = 3.6$  dex. The factor  $c$  in Eq. 1 is introduced  
 170 for modulating *EEP* because the track step on the  $\log T_{\text{eff}} - \log g$   
 171 diagram is not uniform. To avoid obvious data gaps, we test some  
 172 cases and find that  $c = 0.18$  gives the most uniform data distribution.  
 173 In Figure 1, we demonstrate how the effective temperature changes  
 174 with fractional age and *EEP*. It can be seen that *EEP* is a better  
 175 choice than the fractional age because global parameters change  
 176 smoother around the blue hook and the turn-off point.

177 **2.4 Sampling method**

178 There is a limitation of the data size in the GP framework, because  
 179 the computational and memory complexity exponentially increase  
 180 with the number of data points. In practice, the typical data size  
 181 is on an order of  $10^4$ . Given that the grid contents  $\sim 10,000,000$   
 182 stellar models, only a small subset can be used for training. The sam-  
 183 pling method is hence important. A flat sampling is not appropriate,  
 184 because the evolving step is not uniform at different evolutionary  
 185 stages due to the MESA step-control strategy. For instance, stellar  
 186 models are dense at the main-sequence and lower RGB but quite  
 187 sparse at the subgiant stage. We test a few methods and find that  
 188 using the displacement ( $\delta d_n$ ) defined in Eq. 1 as the weight to sam-  
 189 ple models on an evolutionary track gives a relatively uniform data  
 190 distribution at different evolutionary stages.

191 **3 GAUSSIAN PROCESS MODEL**

192 A GP can be applied as a probabilistic model to a regression prob-  
 193 lem. Here we use the GP model to generalise a stellar model grid to

194 a continuous and probabilistic function that maps inputs to observ-  
 195 able quantities. This allows us to predict observable quantities for  
 196 off-grid regions. We intend to train GP models that maps five funda-  
 197 mental inputs, i.e., mass ( $M$ ), initial metallicity ( $[Fe/H]_{\text{init}}$ ), initial  
 198 helium fraction ( $Y_{\text{init}}$ ), the mixing-length parameter ( $\alpha_{\text{MLT}}$ ), and  
 199 equivalent evolutionary phase (*EEP*), to five model outputs includ-  
 200 ing effective temperature ( $T_{\text{eff}}$ ), surface gravity ( $\log g$ ), radius( $R$ ),  
 201 surface metallicity ( $[Fe/H]$ ), and stellar age ( $\tau$ ). We use the GP  
 202 model as a non-parametric emulator, that is emulating the compara-  
 203 tively slow calls to models of stellar evolution. This emulator can  
 204 be described as a function approximation problem. In fact, the way  
 205 we have implemented the GP as function approximation means that  
 206 we have used one GP for each of the outputs so that they can be  
 207 described as

$$T_{\text{eff}} = f_{T_{\text{eff}}} (M, EEP, [Fe/H]_{\text{init}}, Y_{\text{init}}, \alpha_{\text{MLT}}), \quad (3)$$

$$\log g = f_{\log g} (M, EEP, [Fe/H]_{\text{init}}, Y_{\text{init}}, \alpha_{\text{MLT}}), \quad (4)$$

$$R = f_R (M, EEP, [Fe/H]_{\text{init}}, Y_{\text{init}}, \alpha_{\text{MLT}}), \quad (5)$$

$$[Fe/H] = f_{[Fe/H]} (M, EEP, [Fe/H]_{\text{init}}, Y_{\text{init}}, \alpha_{\text{MLT}}), \quad (6)$$

and

$$\tau = f_{\tau} (M, EEP, [Fe/H]_{\text{init}}, Y_{\text{init}}, \alpha_{\text{MLT}}). \quad (7)$$

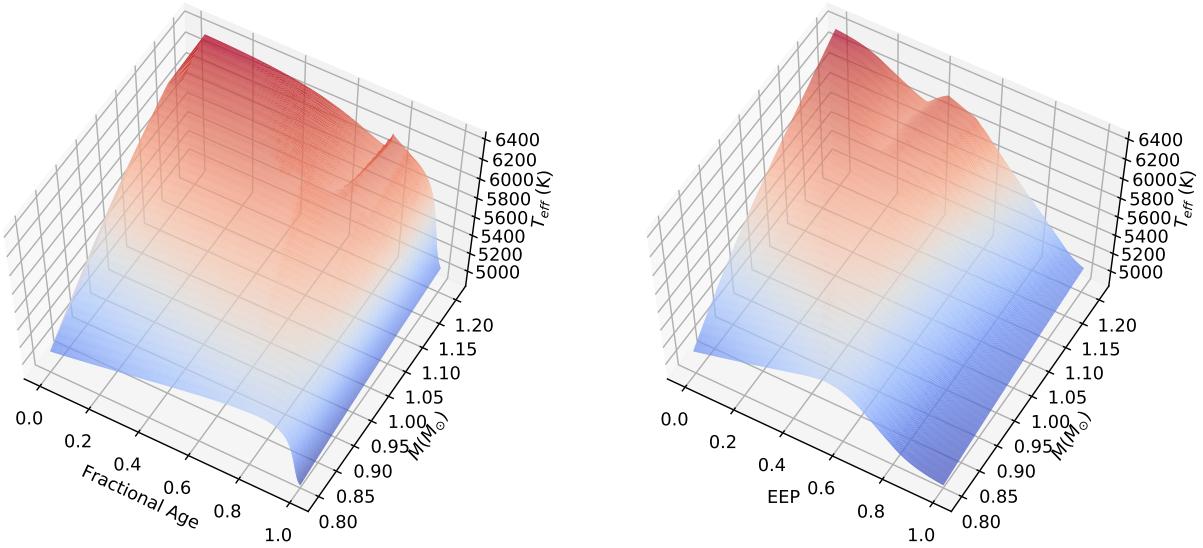
In the following, we introduce the underline theory of GP regression  
 208 and the setup of training GP models.

210 **3.1 Gaussian Process Application**

211 In our application to a stellar model grid, a GP has a number of  
 212 desirable properties. While a GP is a stochastic process, the dis-  
 213 tribution of a GP can be considered as a distribution of functions  
 214 with a continuous domain. In fact, the marginal likelihood consid-  
 215 ered in function space is equal to the likelihood of the data given  
 216 some function values, multiplied by the prior on those function val-  
 217 ues marginalised over all function values Williams & Rasmussen  
 218 (1996). That is to say that, the GP allows for the analytical evalua-  
 219 tion of a fit over many different functions (perhaps an infinite number)  
 220 weighted by some concept of a prior and the agreement with the  
 221 data. In addition, while the marginal likelihood will be assessed  
 222 on discrete data, predictions can be made using linear algebra for  
 223 new data in the continuous domain, but crucially again marginalised  
 224 over these many different functional forms. It is possible to see how  
 225 this might be useful for generalising (or emulating or augmenting)  
 226 a discrete grid of stellar models in order to obtain predictions in the  
 227 continuous domain.

228 In this section we will look at the required mathematics to  
 229 be able to implement a GP for our application to grids of stellar  
 230 models. We start with a series of definitions before dealing with the  
 231 marginal likelihood and the posterior predictive distributions.

232 We start with a grid of stellar models containing  $N$  models with  
 233 a label we want to learn, for example model effective temperature,  
 234 which we will denote with the general symbol  $\mathbf{y}$ , and a set of on-grid  
 235 inputs  $\mathbf{X}$  (e.g., mass, *EEP*, metallicity, ...). We can use a GP to make  
 236 predictions of the effective temperature (labelled  $y$ ) for additional  
 237 off-grid input values given by  $\mathbf{X}_{\star}$ . The vector  $\mathbf{y}$  is arranged  $\mathbf{y} =$   
 238  $(y_1, \dots, y_N)^T$  where the subscript label references the stellar model.



**Figure 1.** Surface plots of model effective temperature on the mass-fractional age (left) and mass-EEP (right) diagrams. Models in this figure are from the primary grid with fixed initial metallicity ( $[Fe/H]_{init} = 0.0$ ), helium fraction ( $Y_{init} = 0.28$ ) and mixing-length parameter ( $\alpha_{MLT} = 2.1$ ). It can be seen that the effective temperature changes much smoother on the mass-EEP diagram at the blue hook and turn-off points.

The input labels are arranged into a  $N \times D$  matrix where  $D$  is the number of input dimensions (e.g.,  $D = 3$  for mass,  $EEP$ , and metallicity) so that  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$  where  $\mathbf{x}_i = (x_{1,i}, \dots, x_{D,i})^T$ . The matrix of additional inputs  $\mathbf{X}_\star$  has the same form as  $\mathbf{X}$  but size  $N_\star \times D$ .

Williams & Rasmussen (1996), from which our description below is based, define a GP as a collection of random variables, where any finite number of which have a joint Gaussian distribution. In general terms, a GP may be written so that our on grid labels are random variables drawn from our GP distribution,

$$\mathbf{y}(\mathbf{X}) \sim \mathcal{GP}(m(\mathbf{X}), \Sigma), \quad (8)$$

where  $m(\mathbf{X})$  is some mean function, and  $\Sigma$  is some covariance matrix. The mean function controls the deterministic part of the regression and the covariance function controls the stochastic part. The mean function defined here could be any deterministic function and we will label the additional parameters, or hyperparameters,  $\phi$ . Each element of the more familiar covariance matrix is defined by the covariance function or *kernel function*  $\mathbf{K}$  which has hyperparameters  $\theta$  and is given by,

$$\Sigma = \mathbf{K}(\mathbf{X}, \mathbf{X}, \theta), \quad (9)$$

or

$$\Sigma_{n,m} = k(\mathbf{X}_n, \mathbf{X}_m, \theta), \quad (10)$$

where the inputs  $\mathbf{X}_n$  and  $\mathbf{X}_m$  are  $D$ -dimensional vectors and the output is a scalar covariance. In addition to the covariance defined by the kernel function, we include additional white noise in the covariance matrix by adding an identity matrix  $\mathcal{I}$  multiplied by a scalar value  $\sigma_w^2$ , so that,

$$\Sigma = \mathbf{K}(\mathbf{X}, \mathbf{X}, \theta) + \sigma_w^2 \mathcal{I}, \quad (11)$$

where  $\sigma_w^2$  is another hyperparameter to be learnt during training.

### 3.1.1 The likelihood

Conceptually we value the GP because of its ability to marginalise over many functions  $\mathbf{f}$  and return a marginal likelihood,

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{f}, \mathbf{X}) p(\mathbf{f}|\mathbf{X}) d\mathbf{f}, \quad (12)$$

noting that this function space marginal likelihood is weighted by the probability of the data given the function and the probability of the function. This integral could be evaluated. However, by noting that a GP is a collection of random variables, where any finite number of which have a joint Gaussian distribution, the marginal probability of our data  $\mathbf{y}$  is also the joint likelihood of a multivariate normal distribution,

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(m(\mathbf{X}), \Sigma), \quad (13)$$

which can be straightforward to evaluate. Thus the marginal likelihood is,

$$p(\mathbf{y}|\mathbf{X}) = (2\pi)^{k/2} \det(\Sigma)^{-0.5} \exp\left(\frac{-1}{2} (\mathbf{X} - m(\mathbf{X}))^T \Sigma^{-1} (\mathbf{X} - m(\mathbf{X}))\right), \quad (14)$$

which can be evaluated without integrating over all possible function space. While this marginal likelihood expression is clearly more computationally feasible than the integral over functional space is not without its limitations. Because it is necessary to calculate the determinant and the inverse of the covariance matrix, typically applied algorithms, make this a  $O(N^3)$  or  $O(N^2 \log N)$  operation. This naturally limits the size of the data set for which the likelihood, and optimisations of the likelihood, can be applied.

## 283 3.1.2 Making predictions

284 If we want to obtain predictive distributions for the output  $\mathbf{y}_\star$  given  
 285 the inputs  $\mathbf{X}_\star$ , the joint probability distribution of  $\mathbf{y}$  and  $\mathbf{y}_\star$  is  
 286 Gaussian and given by

$$p\left(\begin{bmatrix} \mathbf{y} \\ \mathbf{y}_\star \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} m(\mathbf{X}) \\ m(\mathbf{X}_\star) \end{bmatrix}, \begin{bmatrix} \Sigma & \mathbf{K}_\star \\ \mathbf{K}_\star^T & \mathbf{K}_{\star\star} \end{bmatrix}\right), \quad (15)$$

287 where the covariance matrices  $\Sigma$  and  $\mathbf{K}$  are computed using the  
 288 kernel function so that,

$$\Sigma_{n,m} = k(\mathbf{X}_n, \mathbf{X}_m), \quad (16)$$

which is an  $N \times N$  matrix.

$$\mathbf{K}_{\star n,m} = k(\mathbf{X}_n, \mathbf{X}_{\star m}), \quad (17)$$

which is an  $N \times N_\star$  matrix, and finally

$$\mathbf{K}_{\star\star n,m} = k(\mathbf{X}_{\star n}, \mathbf{X}_{\star m}), \quad (18)$$

289 which is an  $N_\star \times N_\star$  matrix. The predictions of  $\mathbf{y}_\star$  are again a  
 290 Gaussian distribution so that,

$$\mathbf{y}_\star \sim \mathcal{N}(\hat{\mathbf{y}}_\star, \mathbf{C}), \quad (19)$$

where

$$\hat{\mathbf{y}}_\star = m(\mathbf{X}_\star) + \mathbf{K}_\star^T \Sigma^{-1} (\mathbf{y} - m(\mathbf{X})), \quad (20)$$

and

$$\mathbf{C} = \mathbf{K}_{\star\star} - \mathbf{K}_\star^T \Sigma^{-1} \mathbf{K}_\star. \quad (21)$$

291 At this point we can make predictions on model properties  
 292 given a grid of stellar models using equation 19. But these predic-  
 293 tions will likely be poor unless we select sensible values for the form  
 294 and hyperparameters of the mean function and covariance function.  
 295 In the following section we detail a number of kernel functions that  
 296 will be tested against the data. We will then discuss the method for  
 297 determining the values of the hyperparameters to be used.

## 298 3.2 Setup of GP Models

## 299 3.2.1 Tool package

300 We adopt a tool package named GPyTorch, which is a GP frame-  
 301 work developed by Gardner et al. (2018). It is a Gaussian process  
 302 library based on an open source machine-learning framework Py-  
 303 Torch<sup>1</sup>. The package provides significant GPU acceleration, state-  
 304 of-the-art implementations of the latest algorithmic advances for  
 305 scalability and flexibility, and easy integration with deep learning  
 306 frameworks.<sup>2</sup> We train GP models on a NVidia Tesla V100 graph-  
 307 ics processing unit (GPU) with 32GB GPU Memory. The GPU  
 308 captivity allows a training dataset with up to  $\sim 20,000$  data points.

<sup>1</sup> <https://pytorch.org>

<sup>2</sup> Source codes and detailed introductions are available on <https://gpytorch.ai>.

## 309 3.2.2 Training procedure

310 The training procedure of a GP model includes training, validating,  
 311 and testing. In the training process, we iteratively optimise hyper-  
 312 parameters of a GP model to learn the underline function which  
 313 maps inputs to outputs from on-grid evolutionary tracks (training  
 314 dataset). In each iteration, the GP model is validated by comparing  
 315 true and GP predicted values of some off-grid tracks (validating  
 316 dataset). Although the validating dataset is not directly involved in  
 317 training hyperparameters, it still constructs the GP model to some  
 318 extend because the optimal solution is the one that best fits the  
 319 validating dataset. For this reason, the validating dataset does not  
 320 give a completely independent validation for a GP model. We hence  
 321 have a testing process after the training. The testing dataset contains  
 322 some other off-grid tracks which are reserved from the training and  
 323 validating process. The testing dataset are also used to estimate the  
 324 systematic uncertainties of GP model.

325 Here we briefly summary the setup of GP model training. We  
 326 apply an Architecture Neural Network (ANN) including 6 hidden  
 327 layers and 128 nodes per layer as the mean function. Note that this is  
 328 not training an ANN to learn the data in detail. The ANN is quickly  
 329 trained at the beginning to interpret the complex mean function in  
 330 multiple-dimension space to accelerate the whole training process.  
 331 In the GP model training, the mean function is normal uninteresting  
 332 because all the inference effort is spent on estimating the correct  
 333 covariance function. In our tests, GP models with the linear or the  
 334 constant mean function could achieve similar results, but it takes  
 335 more time for models to converge. The GPyTorch standard likeli-  
 336 hood for regression, which assumes a standard homoskedastic noise  
 337 model, is applied as the likelihood function. We use the negative  
 338 logarithm of the likelihood as the loss function. The optimiser for  
 339 training is called ‘Adam’, which is a combination of the advantages  
 340 of two other extensions of stochastic gradient descent, specifically,  
 341 Adaptive Gradient Algorithm and Root Mean Square Propagation  
 342 (Kingma & Ba 2014). More detailed discussions about these choices  
 343 can be seen in the Appendix A.

344 We set up the so-called ‘Early Stopping’ procedure to decide  
 345 when to terminate the training. The procedure evaluates GP models  
 346 on a holdout validation dataset after each iteration. If the perfor-  
 347 mance of the GP model on the validation dataset starts to degrade  
 348 or stops upgrading after many iterations, then the training process is  
 349 terminated (see discussions in Anzai 2012; Goodfellow et al. 2016).  
 350 The ‘Early Stopping’ procedure can reduce overfitting and improve  
 351 the generalisation of GP models. We use the aforementioned vali-  
 352 dating EI to monitor the training and terminate it when there is no  
 353 improvement for 300 iterations.

354 To save the best learned GP models, we check the validating  
 355 EI after every iteration. The current model will be saved to replace  
 356 the last saving if it has the so-far lowest validation errors. This is to  
 357 say, the final saved model is the one with the best performance in  
 358 the training process.

## 359 3.2.3 Kernel Function

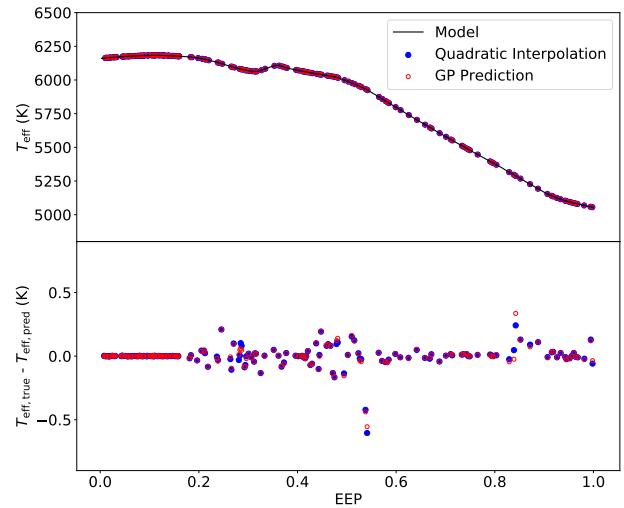
360 To select the proper kernel function for training GP models, we test  
 361 four basic kernels and a number of combined kernels. The four basic  
 362 kernels are listed as follow.

- 363 • **RBF**: Radial Basis Function kernel (also known as squared  
 364 exponential kernel)
- 365 • **RQ**: Rational Quadratic Kernel (equivalent to adding together  
 366 many RBF kernels with different lengthscales)

- 367 • **Mat12**: Matern 1/2 kernel (equivalent to the Exponential Kernel)  
 368  
 369 • **Mat32**: Matern 3/2 kernel

370 These four kernels are all universal, and we can integrate each of  
 371 them against most functions. Every function in its prior has infinitely  
 372 many derivatives (Williams & Rasmussen 1996). The differences  
 373 between these kernels, in a simply way, can be understood as their  
 374 smoothness/flexibility levels. The RBF kernel is very smooth function  
 375 and can be expressed as a product of a polynomial. It hence  
 376 suits for the case when the data follow a slowly varying function.  
 377 The RQ kernel, as a combination of many RBF kernels, is more  
 378 complex and is able to fit to data with a number of smooth underline  
 379 functions (e.g. when the output depends on multiple inputs). On the  
 380 opposite, the Mat12 gives the absolute exponential kernel, which  
 381 is hence very spiky. It can fit to any sharp variations in the data.  
 382 The Mat32 kernel has a smoothness somewhere between the RBF  
 383 and Mat12 kernels, because it is a combination of an exponential  
 384 and a polynomial. It is a smooth function but has significantly more  
 385 extrema than the RBF kernel. For each kernel, there are hyper para-  
 386 meters (e.g. the lengthscale) to modulate the smoothness as well.  
 387 The hyper parameters can be either prioritised or purely determined  
 388 in training process. A less smooth kernel like Mat32 with a large  
 389 lengthscale could behave similarly to the RBF kernel. However,  
 390 RBF kernel with small lengthscale can not reproduce those extrema  
 391 in Mat32 kernel. These differences are not chance coincidence, and  
 392 the origin of these differences are crucial for interpreting the results.  
 393 Choosing a good kernel for a particular application is necessary for  
 394 good predictions. If the kernel function is too spiky for the data,  
 395 the learnt function could over explain some random variations. On  
 396 the contrary, if the kernel function is too smooth, it may not fully  
 397 capture the variability of the underlying function and lead to an in-  
 398 crease in bias. Combining basic kernels can increase the flexibility  
 399 of a kernel. For instance, the combination of a RBF kernel and a  
 400 Mat12 kernel is able to fit to both smooth and spiky features in the  
 401 data. However, overfitting is a risk when using a combined kernel  
 402 because the kernel could be over-flexible.

403 According to the stellar theory, changing fundamental input  
 404 parameters smoothly changes the dependent outputs. Hence the  
 405 kernel function for our application needs to be smooth. Moreover,  
 406 we can notice in Figure 1 that the observable quantities fast vary  
 407 at some particular regions (e.g., around the blue hook). This means  
 408 that a slowly varying function like the RBF kernel may under fit in  
 409 these areas. We do a number of preliminary studies of training GP  
 410 models with different kernel functions to choose the proper kernel  
 411 function. Details of training will be mentioned in Section 4.2. Here  
 412 we only summary the results. We apply four basic kernels and a  
 413 number of their combinations (RBF + Mat21, RQ + Mat21, Mat32  
 414 + Mat21, RBF + Mat32, RQ + Mat32) to train GP models. The  
 415 combined kernel RBF+Mat21 gives the best fit to the training data,  
 416 however, its testing errors are large. This indicates that the kernel  
 417 is too flexible and hence overfits to the data. The kernel having  
 418 the best performance is the Mat32. The GP model with Mat32 fits  
 419 training data reasonably well and gives the best predictions for the  
 420 testing models. The results match our expectations. What we need is  
 421 a smooth function but not too smooth to fit to the quick variations at  
 422 some particular evolutionary phases. The Mat32 kernel is apparently  
 423 suitable for our application.



**Figure 2.** GP application on 1D problem. Models on this track are split into training and testing data by 70-to-30. Top: the evolution of effective temperature for a  $1.1M_{\odot}$  track. The grey line is the evolutionary track computed with MESA; blue and red circles indicate predictions for the testing data from the quadratic interpolator and the GP model. Bottom: residuals of predictions in the top graph.

#### 4 PRELIMINARY STUDIES

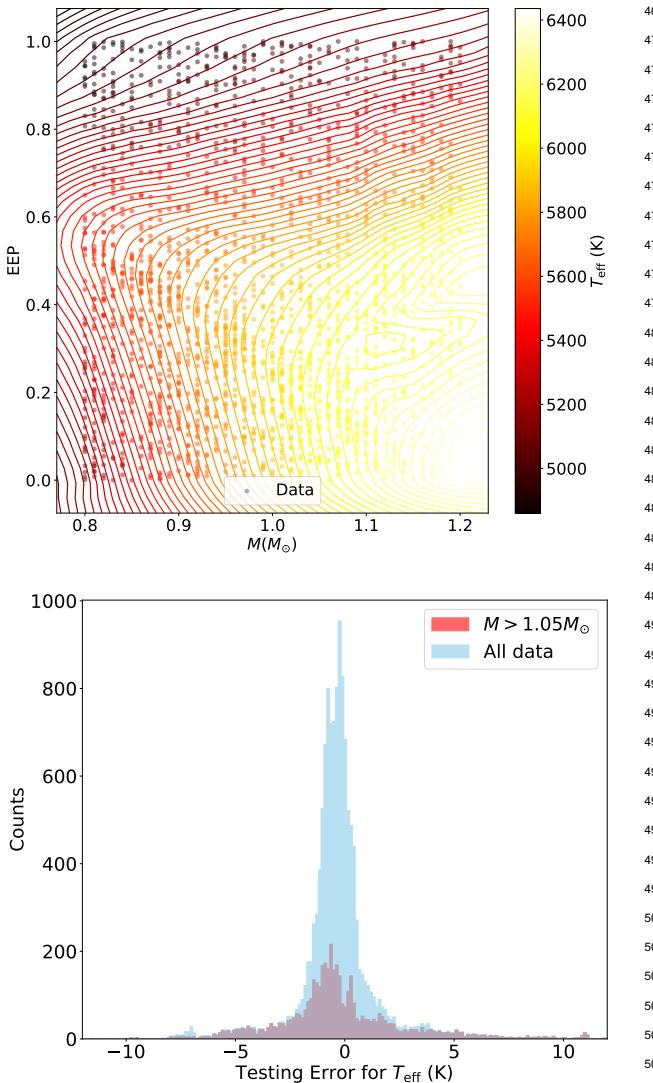
424 Before training the whole model grid (with five input demissions),  
 425 we start with a number of preliminary studies on low-demission  
 426 dataset. These preliminary studies are for several purposes. In the  
 427 1D problem (training data on a single evolutionary track), we compare  
 428 GP predictions with the classical interpolator. In the 2D problem,  
 429 we train GP models on a mass – *EEP* platform to test the  
 430 performances of using different kernel functions. We also discuss  
 431 about introducing a new error index for validating and testing GP  
 432 models instead of using a global error quantity such like Root Mean  
 433 Square Error. In the 3D problem, where GP maps three fundamental  
 434 inputs (mass, *EEP*, and metallicity) onto observables, we solve the  
 435 training strategy for the large dataset whose data size excesses the  
 436 practical limitation.

##### 4.1 1D Problem

437 We first demonstrate an example of GP application on an 1D prob-  
 438 lem. We train a GP model using the Mat32 kernel to learn the  
 439 evolution of effective temperature for a  $1.1M_{\odot}$  track. We split the  
 440 model data points on this track into training and testing data by  
 441 70-to-30. We train a GP model which maps *EEP* to effective tem-  
 442 peratures and then test GP-predicted effective temperatures with  
 443 truths. As it can be seen in Figure 2, the GP model gives very good  
 444 predictions with residuals less than  $\pm 0.5$  K. As a comparison, we  
 445 fit the training data with the quadratic function and use the fitted  
 446 function to do the same prediction. We find very similar results  
 447 from the two methods. It suggests that GP can be an alternative of  
 448 classical interpolators on the 1D problem.

##### 4.2 2D Problem

451 As a further step, we train GP models on a 2D problem where GP  
 452 models map mass and *EEP* to the five observable outputs (Outputs



**Figure 3.** Top: The 2D GP model for  $T_{\text{eff}}$ . Bottom: probability distributions of validating errors of the GP model.

=  $f(M, EEP)$ ). Training data are selected from the primary grid with fixed  $[Fe/H]_{\text{init}}(0.0)$ ,  $Y_{\text{init}}(0.28)$ , and  $\alpha_{\text{MLT}}(2.1)$ . There are 41 evolutionary tracks which content 24,257 models, and we sample 20,000 of them as training data. To validate and test GP models, we compute 44 evolutionary tracks with the same  $[Fe/H]_{\text{init}}$ ,  $Y_{\text{init}}$ , and  $\alpha_{\text{MLT}}$  but randomly sampled  $M$ . We split off-grid tracks half-to-half as validating and testing datasets. The script developed based on the SIMPLE GP REGRESSION example<sup>3</sup>. We change the mean function and optimiser in the example and add an early stopping and a model saving modules. We follow the aforementioned training procedure to train, validate, and test GP models. We illustrate the learned GP model for effective temperature on the mass– $EEP$  diagram in Figure 3. As it shown that, GP transforms the sparse data onto a continuous function and hence is able to predict values for unseen points in the grid.

<sup>3</sup> [https://docs.gpytorch.ai/en/stable/examples/01\\_Exact\\_GPs/Simple\\_GP\\_Regression.html](https://docs.gpytorch.ai/en/stable/examples/01_Exact_GPs/Simple_GP_Regression.html)

It can be seen in Figure 3 that kernels in the area of  $M \geq 1.05M_{\odot}$  and  $EEP \leq 0.7$  is more complex than those for other regions. There are two regions where global parameters vary relatively fast. The first is around the blue hook and main-sequence-turn-off point ( $EEP \sim 0.4$ ) where high-mass tracks sharply turns on the HR diagram. The second is at early subgiant phase ( $EEP \sim 0.6$ ) where stars fast restructure. Features in these particular areas are relatively difficult to learn and hence poorly predicted by the GP model. When there is a subregion in which the GP model performs worse than other areas, the error distribution would not follow a Normal function. As shown at the bottom of Figure 3, the density distribution of testing errors form long tails which contents about 10% data. The cases for other two global parameters surface gravity and radius are similar to the effective temperature.

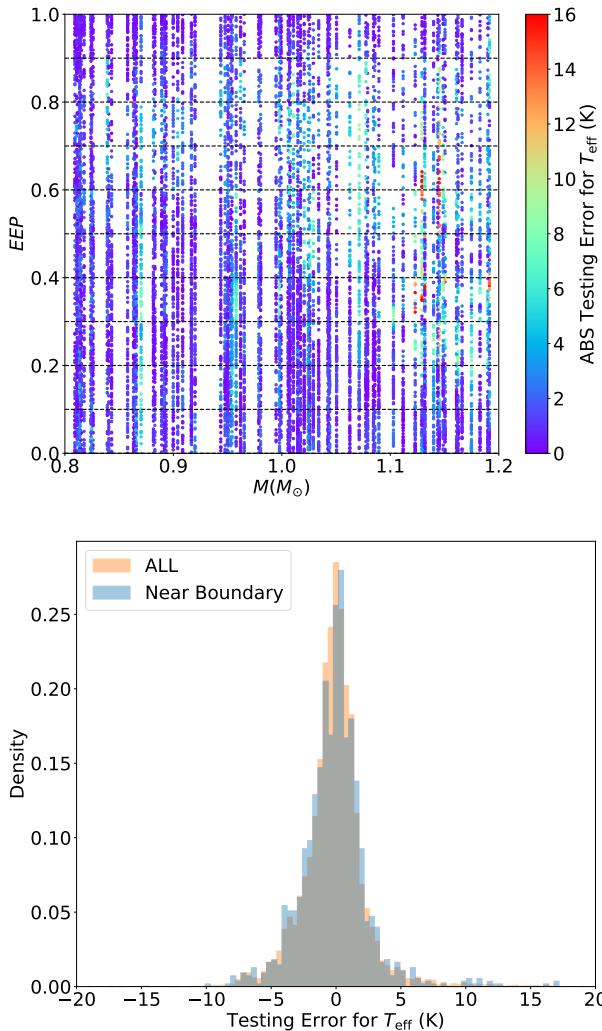
We also find substructures when inspecting testing errors for metallicity and age. The region where the surface metallicity quickly changes is at the early subgiant phase for relatively high-mass tracks. This is because high-mass models maintain shallow convective envelope and hence have strong diffusion effect during the main-sequence stage. At the early subgiant phase, the quick expansion of the surface convective envelope mixes up the settled heavy elements, leading to a fast raise of the surface metallicity. The accuracy of age prediction drops down for very old low-mass stellar models. This is because age values vary in a relatively big dynamic range (15 – 50 Gyr) in a small fraction of data points. Poor GP predictions are caused by the low age resolution.

The error distribution causes an issue in validating and testing GP models. What we normally use are some global errors, such as Root Mean Square Error (RMSE), to represent the validating or testing results. For our case, a global error is not able to point out how GP performs in regions where an observable quantity quickly varies. We want to have an error index that can reflect the GP model performance in general as well as in those sub-areas. By inspecting the error distributions of all five outputs, we find the data points in the tails (outside the 3 times of full width at half maximum) are around 10% (8 – 12% for different outputs). For the majority (90%) of data points, which from a Gaussian-like profile, the 68% confidential interval ( $1-\sigma$  uncertainty) can be used to reflect the global accuracy. For the worst 10% of the data, we could use the 95% and 99.7% confidential intervals (2- and 3- $\sigma$  uncertainties) to describe the median and the length of the tail. Thus, we define an Error Index (EI), which is the sum of 68%, 95%, and 99.7% cumulative values of the absolute errors. For the case in Figure 3, cumulative values at 68%, 95%, and 99.7% are 1.1, 4.9, and 11.1K, which give a testing EI equals to 17.1K. We apply this EI in all following training processes to validate and test GP models.

We train GP models using different kernel functions to investigate which is the best for our application. We do this with the 2D data because the training is fast to be able to test many different options. As mentioned in Section 3.2.3, we find Mat32 is the most suitable kernel for mapping the stellar model grid.

### 4.3 3D Problem: Strategy for Large Data Sample

We apply GP to a 3D problem where GP maps three fundamental inputs, i.e.,  $M$ ,  $EEP$ , and  $[Fe/H]_{\text{init}}$  to observables. The main purpose of this preliminary study is investigate the strategy for training large data sample that exceeds the data size limitation of 20,000. We select training data from the primary grid with  $Y_{\text{init}} = 0.28$  and  $\alpha_{\text{MLT}} = 2.1$ . The training dataset contents ~300,000 data points which is 15 times the limit of training data size (20,000). For validating and



**Figure 4.** Top: Testing errors of 3D GP model for  $T_{\text{eff}}$  on the  $M - \text{EEP}$  diagram. Dashes indicate section boundaries. Bottom: examination of the edge effects of the section scenario. Probability distributions of testing errors of all testing data and those near the boundary ( $\pm 0.01 \text{ EEP}$ ) in the upper graph are compared. As it can be seen, testing errors do not raise around the boundary.

problems. Because we need much more data give the fact that the parameter space exponentially increases with the demission. A simple way to overcome this issue is breaking the grid into many sections and train GP models for each section separately. We divide the training dataset into 10 equal sections by  $\text{EEP}$  and sample 20,000 training data in each. A set of GP models are then trained for each  $\text{EEP}$  section. Using this section scenario, we improve the testing EI for the five output parameters by around 10%. For instance, the testing EI for  $T_{\text{eff}}$  decreases from 23.5 to 21.6K. (1.7, 5.0, 14.9K at at 68th, 95th, and 99.7th). EI values for the five outputs before and after sectioning can be seen in Table 3. The section scenario outperforms the SVGP and SKI GP methods. We hence apply it as our training strategy.

The section scenario improves the performance of the GP model, but there is a major concern about the edge effect at the boundary between segments. If a GP model works significantly poorly at these boundaries, it will be difficult to map the systematic errors across the whole parameter space. We hence examine potential edge effects as illustrated in Figure 4. We inspect absolute testing errors for each output on the  $M - \text{EEP}$  diagram. No obvious edge effect is found. We also do a statistical comparison between all errors and those around section boundaries ( $\pm 0.01 \text{ EEP}$ ). As shown in the bottom graph, the density distributions of the two samples are very similar to each other.

## 5 AUGMENTING THE STELLAR GRID

Based on what we find from preliminary studies, we now apply GP to mapping the whole 5D model grid. The setup of GP model is summarised in Table 2. The training data are sampled from both primary and additional girds (as described in Table 1). The additional grid increases the grid resolution for relatively high-mass models, and this gives more information about the blue hook for GP to learn. We also computed 4,880 off-grid tracks. These off-grid tracks are split by 50-to-50 for validating (in the training progress) and testing (after the training progress) GP models.

The section scenario is applied. For each section, we train a set of GP models for each output parameter with 20,000 training and 20,000 validating data. The number of sections need to be tested to obtain the best efficiency. To do this, we gradually increase the number of sections from 1 to 100 and track down the changes in testing EI. We find significant improvement from 1 to 10 sections but no further improvements for more than 10 sections. We list the testing EI with different numbers of sections in Table 3. It turns out that dividing the grid into 10 sections (corresponding to a 2% sampling rate) is the most efficient.

We use GP models for the 10-sections case as our final result. All following analysis and discussion are based on it. When testing GP models, we do not section the dataset because the data size limitation for testing is not strict. We sample 100,000 off-grid stellar models as the testing dataset. Note that we do not use models with  $\tau \geq 20.0 \text{ Gyr}$ ,  $[\text{Fe}/\text{H}]_{\text{surf}} \leq -0.6 \text{ dex}$ , or  $T_{\text{eff}} \geq 7000 \text{ K}$  for testing because we find strong edge effects in those ranges.

### 5.1 Overview of Results

A overview of testing errors (Truths - GP predictions) can be seen in Figure 5, where we plot rolling medians and rolling standard deviations for all outputs' errors against fundamental inputs. Median values are approximate along zero in most plots, indicating good

528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544

testing purposes, we compute another 174 evolutionary tracks with the same input  $Y_{\text{init}}$  and  $\alpha_{\text{MLT}}$  but random input  $M$  and  $[\text{Fe}/\text{H}]_{\text{init}}$ .

We start with sampling 20,000 training data and train a set of

GP models. We then obtain testing EI for each output. For instance,

the testing EI for  $T_{\text{eff}}$  is 23.5K (2.0, 5.8, and 15.7 K at 68th, 95th, and 99.7th).

We then apply two state-of-the-art approaches designed for

large dataset to train the model. These two approaches are named as

Stochastic Variational GP (SVGP) and Structured Kernel Interpolation (SKI GP). We only find some minor improvements in the GP

predictions. Comparing the testing EI for  $T_{\text{eff}}$ , SVGP gives a results

of  $\text{EI} = 24.1 \text{ K}$  (2.2, 6.8, and 15.1K at at 68th, 95th, and 99.7th)

and SKI GP ends up with  $\text{EI} = 22.9 \text{ K}$  (2.0, 6.1, and 14.8K at 68th,

95th, and 99.7th). Details about these two implements and some

discussions about the results can be seen in appendix B.

We seek for better strategy for training large data sample. The

GPU memory captivity limits the actual number of data that induce

the kernel. This limitation becomes crucial for the high-demision

**Table 2.** Setup of GP Models

GP model inputs		
Parameter	Notation	Range
Mass	$M$	0.8–1.2 $M_{\odot}$
Equivalent evolutionary phase	$EEP$	0 – 1
Initial metallicity	$[Fe/H]_{init}$	-0.5 – 0.5 dex
Initial helium fraction	$Y_{init}$	0.24 – 0.32
Mixing-length parameter	$\alpha_{MLT}$	1.7 – 2.5
GP model outputs		
Parameter	Notation	Trust-worth range <sup>a</sup>
Effective temperature	$T_{eff}$	$\leq 7000$ K
Surface gravity	$\log g$	-
Radius	$R$	-
Surface metallicity	$[Fe/H]$	$\geq -0.6$ dex
Stellar age	$\tau$	$\leq 20$ Gyr
Setup of training		
Item	Adopted	
Kernel	Mat32	
Mean Function	6 layers x 128 notes Neural Network	
Likelihood Function	Gaussian Likelihood Function	
Loss Function	Exact marginal likelihood	
Optimiser	Adam including AMSGRAD variant	
Termination	Early Stopping (monitoring the validating $EI$ )	

<sup>a</sup> The ranges without strong edge effects.

agreement between GP predictions and true values. The 68% confidential intervals are generally small and their dynamical ranges do not significantly vary across input ranges. However, the 95% confidential intervals have more significant changes and are not well scaled to the 68% confidential intervals. This corresponds to the tail feature as seen in Figure 3. GP predictions are relatively poor in some particular regions. For instance, predictions for the effective temperature are more scattered in high-mass because of the appearances of the blue hook. From these results, it can be seen that the model systematic uncertainty is not uniform across the parameter space. Proper estimates of model uncertainty are hence necessary.

## 5.2 Mapping Systematic Uncertainties

A learned GP model predicts output quantities with uncertainties based on its noise model. In our preliminary studies, uncertainties are properly determined for the 1D and 2D problems, but we find obviously underestimated uncertainties in the 3D problem. In the 5D problem, GP models also predict significantly small uncertainties: they are mostly one order of magnitude smaller than testing errors. This is to say, the learned GP models are over-confident for high-dimension cases. The reason could be the equally spaced training data, from which GP model learn few variations at the scale smaller than the grid step and hence turns to fit with large lengthscale values. Because GP models do not give reliable uncertainties, we intend to use testing errors to estimate the systematic uncertainty in GP prediction. As shown in Figure 5, systematical uncertainties relate to  $M$ ,  $EEP$ , and  $[Fe/H]_{init}$  but not to  $Y_{init}$  or  $\alpha_{MLT}$ . We can treat this as a 3D problem and train another GP model, in which GP model systematic uncertainty is a function of  $M$ ,  $EEP$ , and  $[Fe/H]_{init}$ .

We inspect the testing errors in the  $M$ - $EEP$ - $[Fe/H]_{init}$  space and find that their local medians vary smoothly. We hence apply

the constant mean function and the RBF kernel. The testing dataset contents 100,000 which exceeds the data size limitation. We use the SVGP approach but not the section scenario for this training, because the SVGP can well handle large data following smooth function (see Appendix B for detailed discussions about SVGP). We split the testing error data by 75-to-25 for training and validating. The variational evidence lower bound (ELBO) is adopted as the loss function because it is designed for when there is too much data for the exact inference. We set up Early Stopping by tracking the RMSE value and terminate training when the RMSE value stops decreasing for 100 iterations. The outputs of GP models are the local medians of testing errors. We use them to infer the systematic uncertainties for five observable quantities (referred as  $\sigma_{T_{eff}}$ ,  $\sigma_{\log g}$ ,  $\sigma_R$ ,  $\sigma_{[Fe/H]_{surf}}$ , and  $\sigma_{\tau}$ ). To differentiate these GP models, we refer to them as GP-SYS models. In Figure 6, we compare the actual local systematic uncertainties for  $T_{eff}$  at  $[Fe/H]_{init} \simeq 0.0$  with those given by the GP-SYS models. It shows that the GP-SYS model well reproduces the  $\sigma_{T_{eff}}$  distributions.

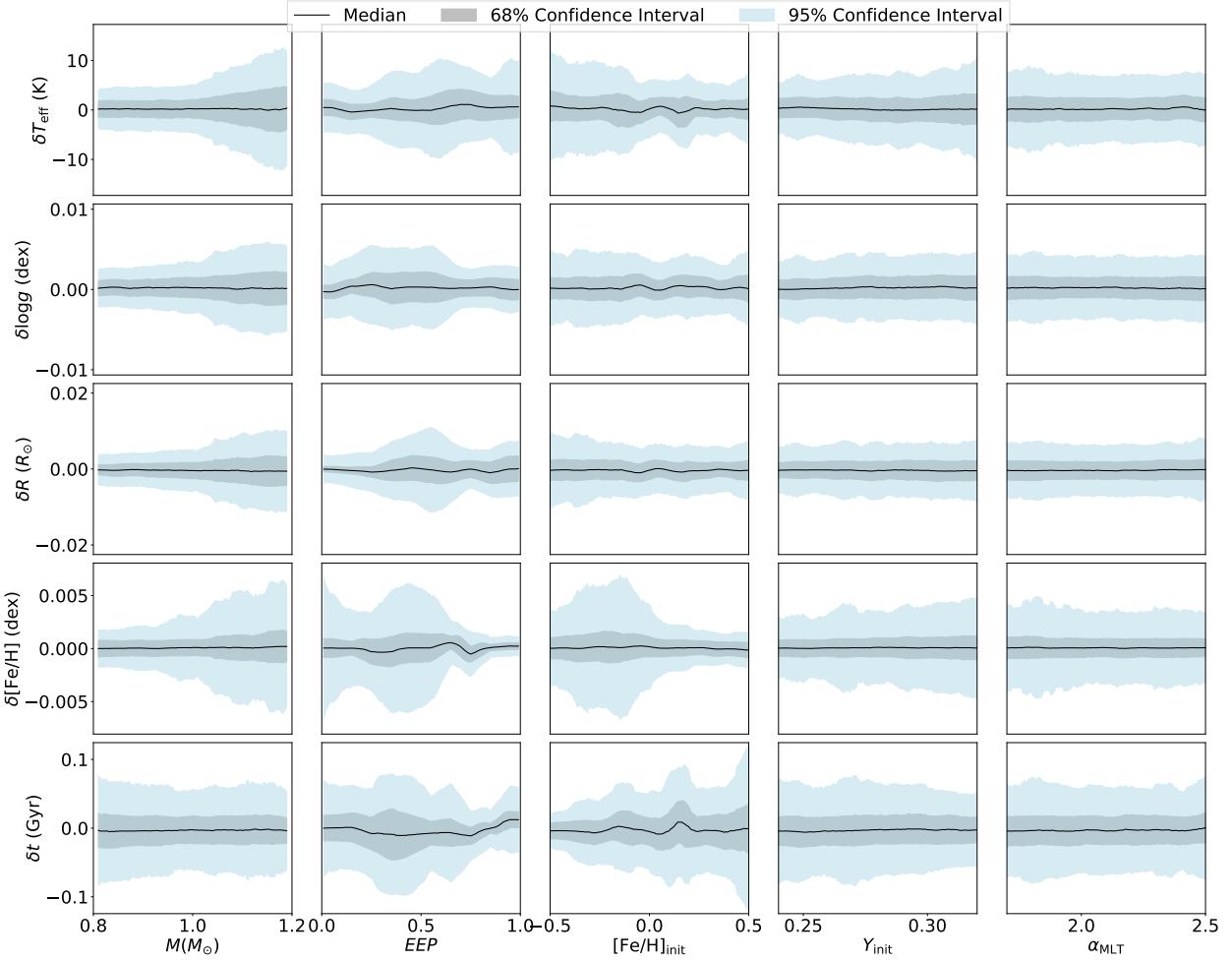
## 6 MODELLING STARS WITH GP PREDICTIONS

### 6.1 Augmenting the model grid

Now we are able to use learned GP models to augment the original stellar grid. We randomly sample 5,000,000 data points with uniform distributions for five fundamental inputs ( $M$ ,  $EEP$ ,  $[Fe/H]_{init}$ ,  $Y_{init}$ , and  $\alpha_{MLT}$ ). We then predict output quantities using GP models and their systematic uncertainties using GP-SYS models. This GP-based model dataset can be downloaded following the instruction at <https://github.com/litanda/GPGird>. In Figure 7, we

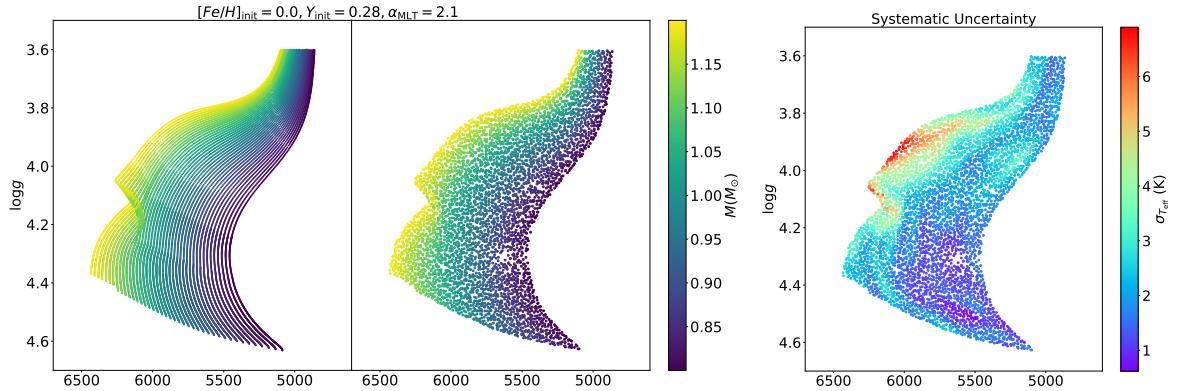
**Table 3.** Training and validating errors for GPR Models

Model Type	Inputs	NTraining	Sampling rate	Testing Errors (at 68/95/99.7%)				
				$T_{\text{eff}}$ (K)	$\log g$ ( $10^{-3}$ dex)	$R$ ( $10^{-3} R_{\odot}$ )	$[Fe/H]_{\text{surf}}$ ( $10^{-3}$ dex)	$\tau$ ( $10^{-2}$ Gyr)
GP	2D	20,000 x 1	96%	1/5/11	1/3/8	2/6/14	0.5/2/12	1/3/9
GP with 10 sections	3D	20,000 x 1	5%	2/6/16	1/4/10	3/7/17	2/6/22	2/7/22
	3D	20,000 x 10	50%	2/5/15	1/4/11	2/7/17	1/3/20	2/6/19
GP	5D	20,000 x 1	0.2%	3/9/34	2/5/18	4/11/36	2/7/30	3/9/27
GP with 3 sections	5D	20,000 x 3	0.6%	3/8/27	2/5/18	3/7/26	1/4/24	3/7/22
GP with 5 sections	5D	20,000 x 5	1%	2/7/25	1/4/15	3/7/24	1/4/21	2/6/22
GP with 10 sections	5D	20,000 x 10	2%	2/7/27	1/4/14	2/7/26	1/4/20	2/6/21
GP with 20 sections	5D	20,000 x 20	4%	2/7/26	1/4/14	2/7/27	1/3/18	2/6/22
GP with 100 sections	5D	20,000 x 100	20%	2/7/25	1/4/14	2/7/26	1/3/17	2/6/18

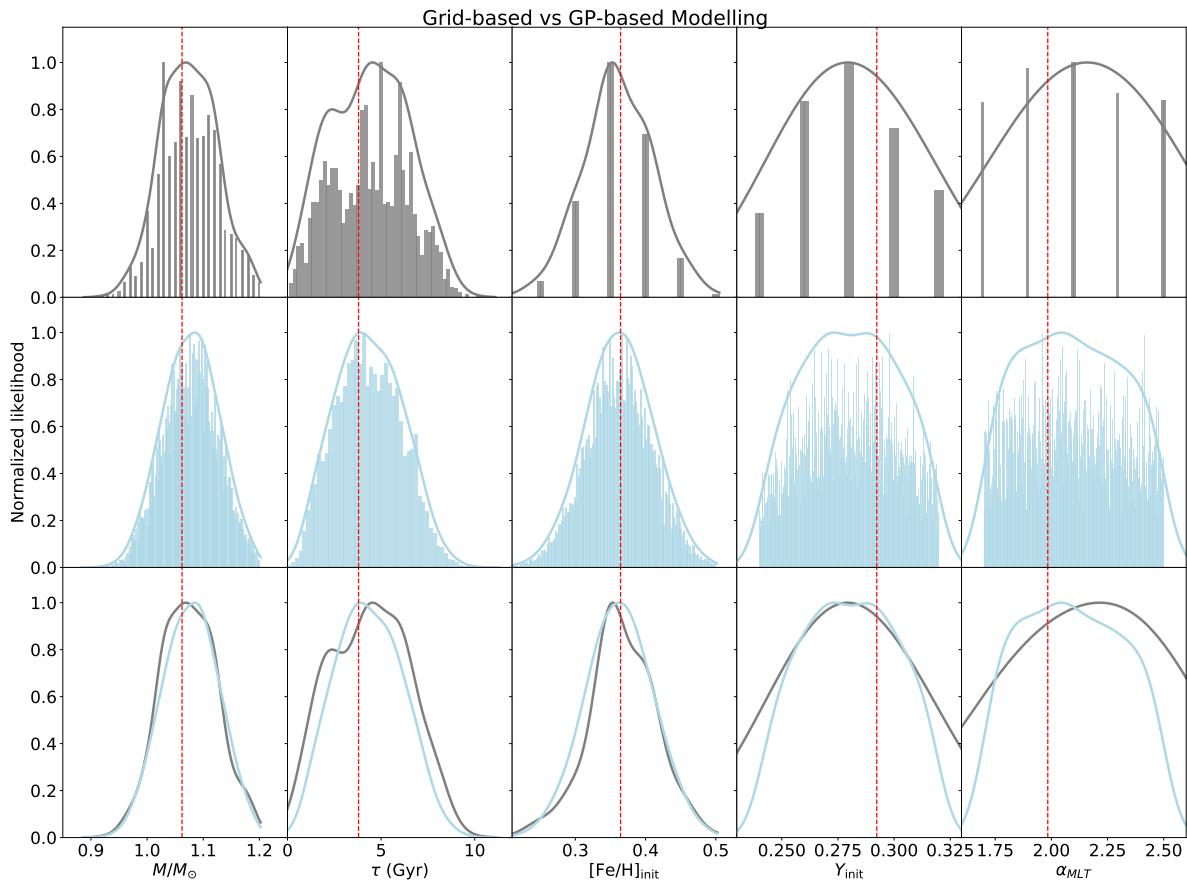


**Figure 5.** Roll medians and 68/95% confidential intervals of testing errors against GP model inputs. Black solid lines indicate the median value; grey and blue shadows represent the 68% and 95% confidential interval. Testing errors of  $T_{\text{eff}}$ ,  $\log g$ , and  $R$  mainly depend on  $M$  and  $EEP$ . Metallicity error strongly depends on  $M$ ,  $EEP$ , and  $[Fe/H]_{\text{init}}$ , and age error has a significant correlation to  $EEP$  and  $[Fe/H]_{\text{init}}$ . However, testing errors do not obviously relate to  $Y_{\text{init}}$  or  $\alpha_{\text{MLT}}$ .





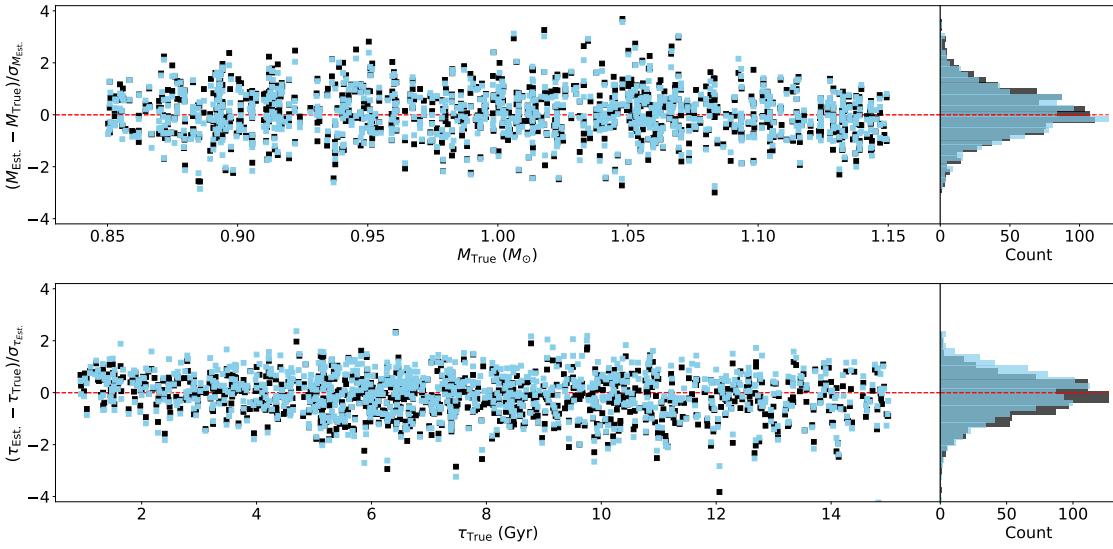
**Figure 7.** Left: Comparing the original grid and GP predictions on Kiel diagram. Right: The systematical uncertainties for  $T_{\text{eff}}$  given by GP-SYS models.



**Figure 8.** Probability distributions of five fundamental parameters from grid-based (top row) and GP-based modelling (middle row) for a fake star. Grey solid lines in the top row and blue solid lines in the middle row are the kernel density of probability distributions. The bottom row demonstrates comparisons between kernel densities based on the two methods. True fundamental parameters, indicated by red dashed lines, are  $M = 1.062 M_\odot$ ,  $\tau = 3.79 \text{ Gyr}$ ,  $[Fe/H]_{\text{init}} = 0.364$ ,  $Y_{\text{init}} = 0.292$ , and  $\alpha_{MLT} = 1.984$ . Observed constraints for this fake star are  $T_{\text{eff}} = 5652 \pm 50 \text{ K}$ ,  $\log g = 4.424 \pm 0.005$ ,  $[Fe/H]_{\text{surf}} = 0.31 \pm 0.05$ , and  $R = 1.047 \pm 0.031 R_\odot$ .

We hence do a number of preliminary studies for the best option. A limitation of the training step for the GP is the size of the training data set and the requirement to perform matrix inversion and calculate the matrix determinant. For training a normal GP model, the practical training data size upper limit is around 20,000 (the number depends on the capacity of GPU device) which is apparently not big

enough for the high-demission grid including ~10,000,000 models. To overcome this issue, we section the training data according to their evolutionary stage (*EEP*) and train GP models for each section. This section scenario significantly improves overall accuracy of GP predictions and we find no edge effects at the boundary of sections. When inspecting the systematic uncertainty of GP models, we notice



**Figure 9.** Differences between true and estimated stellar masses and ages over their estimated uncertainty of 1,000 fake stars. Black and blue symbols represent inferences with the original grid and with GP predictions. Count distributions of offsets are demonstrated on the right side.

that GP predictions are very over-confident in the high-demission problem: the uncertainties given by GP models are mostly smaller than testing errors by an order of magnitude. Because of this, we use the testing errors to estimate systematic uncertainties across the parameter space. Eventually, we provide a GP-based model dataset including 5,000,000 models with randomly sampled fundamental inputs.

We use GP-based models to characterise 1,000 fake stars to examine whether truths of stellar properties can be recovered. We find that GP-based masses and ages are consistent with the injected truth values. The uncertainties are dominated by observational noise, saying that, the systematic uncertainty due to the GP approximation does not obviously affect the modelling on interferences. Comparing with the probability distributions of original sparse grid, GP models are fully sampled in the input range and hence improve the accuracy and precision of inferred parameters. The improvement is remarkable for the stellar age (by 7%). Moreover, the continuous sampling makes it possible to properly estimate some fundamental parameters which are sparse in the grid, e.g., the helium fraction. These results indicate that the method demonstrated in this work is reliable and efficient for interpolating an established model grid and it can improve the modelling solutions because of the statistically-sound sampling.

## ACKNOWLEDGEMENTS

This work has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (CartographY GA. 804752). R.A.G. acknowledges funding from the PLATO CNES grant. Development of GPYTORCH is supported by funding from the Bill and Melinda Gates Foundation, the National Science Foundation, and SAP.

## APPENDIX A: SETUP OF GP MODEL TRAINING

This section includes detailed discussions about the selections of mean function, likelihood and loss function, optimiser, and early stopping.

### A1 Mean Function

We first investigate the mean function. As discussed above, the data distribution is generally smooth but complex in some regions of parameter space (e.g., the sub-giant hook). Although the choice of mean function is not crucial for training GP models, we find that using a constant or a linear mean function leads to a significantly long training time. Hence, we apply a neural network mean function which is flexible enough to manage both simple and complex features to accelerate the training. We adopt an architecture based on that of Lytle et al. (2021) comprising 6 hidden layers and 128 nodes per layer. All layers are fully-connected and the output of each layer, except for the last, is transformed by the Exponential Linear Unit (ELU) activation function Clevert et al. (2015). The ELU activation function provides a smooth function from inputs to outputs, which is preferred over its more common, faster counterpart, the Rectified Linear Unit (RELU).

### A2 Likelihood and Loss Function

Our training dataset is a theoretical model grid, hence there is no observed uncertainty for each data point, but a tiny random uncertainty exists due to the approximations in the MESA numerical method. We model this using  $\sigma_w$ . This noise model is then assumed to be a Gaussian function with a very small variance. A likelihood specifies the mapping from input values  $f(X)$  to observed labels  $y$ . We adopt the standard likelihood for regression which assumes a standard homoskedastic noise model whose conditional distribution is

$$p(y|f(x)) = f + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2), \quad (\text{A1})$$

where  $\sigma$  is a noise parameter. We used a small and fixed noise parameter and run a few tests. However, the strict noise parameter makes GP models hard to converge. When this noise parameter is set as free, it reduces to a small value anyway in the training progress because it is data-driven. For this reason, we did not put strict constraint for or prioritise this noise parameter. In practice, we only set up a loose upper limit ( $\sigma < 0.1$ ) to speed up the training. One thing should be noted that a GP model with a large noise parameter is not a proper description for the stellar grid. Because of this, we only adopt GP models with  $\sigma \lesssim 10^{-4}$ . We train GP models using the negative logarithm of the likelihood function as the loss function.

### A3 Optimiser

We compared two optimisers named SGD and Adam. Here SGD refers to Stochastic Gradient Descent, and Adam is a combination of the advantages of two other extensions of stochastic gradient descent, specifically, Adaptive Gradient Algorithm and Root Mean Square Propagation. The SGD optimiser in the GPyTorch package involves the formula given by Sutskever et al. (2013). The formula makes it possible to train using stochastic gradient descent with momentum thanks to a well-designed random initialisation and a particular type of slowly increasing schedule for the momentum parameter. The application of momentum in SGD could improve its efficiency and make it less likely to stuck in local minimums. On the other hand, the Adam optimiser includes the 'AMSGrad' variant developed by Reddi et al. (2018) to improve its weakness in the convergence to an optimal solution. With these new developments, the two optimisers give very similar results. We finally choose Adam because it works relatively efficiently and stable. We adaptive learning rate in the training process. Our training starts with a learning rate of 0.01 and decreases by a factor of 2 when the loss value does not reduce in previous 100 iterations.

## APPENDIX B: STATE-OF-THE-ART IMPLEMENT FOR LARGE DATASET

In section 4.3, we investigate the strategy for large dataset. We test two State-of-the-art approaches that designed for training big data whose data size is more than the limit of a GP model. Here we mention some details about the two methods and the results.

We first consider the Stochastic Variational GP (SVGP) approach based on the GPyTorch APPROXIMATEGP module. We train our data based on the SVGP example on [https://docs.gpytorch.ai/en/v1.1.1/examples/04\\_Variational\\_and\\_Approximate\\_GPs/SVGP\\_Regression\\_CUDA.html](https://docs.gpytorch.ai/en/v1.1.1/examples/04_Variational_and_Approximate_GPs/SVGP_Regression_CUDA.html). SVGP is an approximate scheme rely on the use of a series of inducing points which can be selected in the parameter space. It trains using mini batches and hence is able to deal with large data size. The other key point of SVGP is the number of inducing points. Because the kernel is only built on these points, the number determines the complexity of kernel. When the underline function is simple, for instance, a power law, a small number of inducing points is enough. For our application, a large number of inducing points are required. Underline principles and detailed descriptions of this approach can be found in Hensman et al. (2014). In our tests on 3D problem, we find a practical issue with the SVGP approach. When we load in a large training sample which takes a lot GPU memory, the rest memory can capacitate only 10,000 inducing points. This is to say, we use more training data but sacrifice the kernel complexity.

The result shows that using SVGP model does not improve the GP predictions compared with normal GP model. For instance, the 68th, 95th, and 99.7th testing errors for  $T_{\text{eff}}$  are 2.2, 6.8, and 15.1K (EI = 24.1 K) for the SVGP and 2.0, 5.8, and 15.7 K (EI = 23.5K) for the normal GP model. This is because the evolutionary feature are complex across multiple demissions. Reducing the kernel complexity is not ideal. We conclude that the SVGP is suitable for training large data which have relatively simple variations but not a good choice for training the model grid.

We also investigate another approach designed for large dataset named Structured Kernel Interpolation (SKI GP). SKI GP was introduced by Wilson & Nickisch (2015). It produces kernel approximations for fast computations through kernel interpolation and is a great way to scale a GP up to very large datasets (100,000+ data points). We follow the example on [https://docs.gpytorch.ai/en/stable/examples/02\\_Scalable\\_Exact\\_GPs/KISSGP\\_Regression.html](https://docs.gpytorch.ai/en/stable/examples/02_Scalable_Exact_GPs/KISSGP_Regression.html) to develop our script. We run a few tests to train a 3D SKI GP model with 100, 000 training data. Compare with the Normal GP and SVGP, its testing errors for  $T_{\text{eff}}$  are slightly improved to 2.0, 6.1, and 14.8K (EI = 22.9K). However, the further test on the 5-demission data is not ideal: a SKI GP model using 100, 000 training data performs much worse than a normal model with only 20,000 training data. The poor behaviour consists with what has been discussed by Wilson & Nickisch (2015): the SKI GP poorly scale to data with high dimensions, since the cost of creating the grid grows exponentially in the amount of data. We attempt to make some additional approximations with the GPyTorch ADDITIVESTRUCTUREKERNEL module. It makes the base kernel to act as one-dimension kernels on each data dimension and the final kernel matrix will be a sum of these 1D kernel matrices. However, the testing errors are not significantly improved.

## REFERENCES

- 93 Anzai Y., 2012, Pattern recognition and machine learning. Elsevier
- 94 Asplund M., Grevesse N., Sauval A. J., Scott P., 2009, *Annual Review of*  
*Astronomy and Astrophysics*, **47**, 481
- 95 Bellinger E. P., Angelou G. C., Hekker S., Basu S., Ball W. H., Guggenberger  
E., 2016, *ApJ*, **830**, 31
- 96 Bishop C. M., 1994
- 97 Choi J., Dotter A., Conroy C., Cantiello M., Paxton B., Johnson B. D., 2016,  
*ApJ*, **823**, 102
- 98 Clevert D.-A., Unterthiner T., Hochreiter S., 2015, arXiv e-prints, [p. arXiv:1511.07289](https://arxiv.org/abs/1511.07289)
- 99 Dotter A., 2016, *ApJS*, **222**, 8
- 100 Ferguson J. W., Alexander D. R., Allard F., Barman T., Bodnarik J. G.,  
Hauschildt P. H., Heffner-Wong A., Tamai A., 2005, *ApJ*, **623**, 585
- 101 Gardner J. R., Pleiss G., Bindel D., Weinberger K. Q., Wilson A. G., 2018,  
in Advances in Neural Information Processing Systems.
- 102 Goodfellow I., Bengio Y., Courville A., 2016, Deep learning. MIT press
- 103 Hendriks L., Aerts C., 2019, *PASP*, **131**, 108001
- 104 Hensman J., Matthews A., Ghahramani Z., 2014, arXiv preprint  
[arXiv:1411.2005](https://arxiv.org/abs/1411.2005)
- 105 Ho T. K., 1995, in Proceedings of 3rd international conference on document  
analysis and recognition. pp 278–282
- 106 Hon M., Stello D., Yu J., 2018, *MNRAS*, **476**, 3233
- 107 Hon M., Bellinger E. P., Hekker S., Stello D., Kuszlewicz J. S., 2020,  
*MNRAS*, **499**, 2445
- 108 Kingma D. P., Ba J., 2014, arXiv e-prints, [p. arXiv:1412.6980](https://arxiv.org/abs/1412.6980)
- 109 Lytle A. J., et al., 2021, *MNRAS*,
- 110 Mombarg J. S. G., Van Reeth T., Aerts C., 2021, arXiv e-prints, [p. arXiv:2103.13394](https://arxiv.org/abs/2103.13394)
- 111 Morel P., Thévenin F., 2002, *A&A*, **390**, 611
- 112 Paquette C., Pelletier C., Fontaine G., Michaud G., 1986, *ApJS*, **61**, 177

- 923 Paxton B., Bildsten L., Dotter A., Herwig F., Lesaffre P., Timmes F., 2011,  
 924     [The Astrophysical Journal Supplement Series, 192, 3](#)  
 925 Paxton B., et al., 2013, [The Astrophysical Journal Supplement Series, 208,](#)  
 926     4  
 927 Paxton B., et al., 2015, [The Astrophysical Journal Supplement Series, 220,](#)  
 928     15  
 929 Paxton B., et al., 2018, [ApJS, 234, 34](#)  
 930 Paxton B., et al., 2019, [ApJS, 243, 10](#)  
 931 Reddi S., Kale S., Kumar S., 2018, in International Conference on Learning  
 932     Representations.  
 933 Rendle B. M., et al., 2019, [MNRAS, 484, 771](#)  
 934 Rogers F. J., Nayfonov A., 2002, [ApJ, 576, 1064](#)  
 935 Schölkopf B., Smola A., Müller K.-R., 1997, in International conference on  
 936     artificial neural networks. pp 583–588  
 937 Sutskever I., Martens J., Dahl G., Hinton G., 2013, in International confer-  
 938     ence on machine learning. pp 1139–1147  
 939 Thoul A. A., Bahcall J. N., Loeb A., 1994, [ApJ, 421, 828](#)  
 940 Verma K., Hanasoge S., Bhattacharya J., Antia H. M., Krishnamurthi G.,  
 941     2016, [MNRAS, 461, 4206](#)  
 942 Williams C. K., Rasmussen C. E., 1996  
 943 Wilson A., Nickisch H., 2015, in International Conference on Machine  
 944     Learning. pp 1775–1784  
 945 Wu Y., et al., 2019, [MNRAS, 484, 5315](#)

946 This paper has been typeset from a  $\text{\TeX}/\text{\LaTeX}$  file prepared by the author.