

Closing the stellar labels gap: Stellar label independent evidence for $[\alpha/M]$ information in *Gaia* BP/RP spectra

ALEXANDER LAROCHE ^{1,2} AND JOSHUA S. SPEAGLE (沈佳士) ^{3,1,2,4}

¹*David A. Dunlap Department of Astronomy & Astrophysics, University of Toronto, 50 St George St, Toronto, ON M5S 3H4, Canada*

²*Dunlap Institute for Astronomy & Astrophysics, University of Toronto, 50 St George Street, Toronto, ON M5S 3H4, Canada*

³*Department of Statistical Sciences, University of Toronto, 9th Floor, Ontario Power Building, 700 University Ave, Toronto, ON M5G 1Z5, Canada*

⁴*Data Sciences Institute, University of Toronto, 17th Floor, Ontario Power Building, 700 University Ave, Toronto, ON M5G 1Z5, Canada*

ABSTRACT

Data-driven models for stellar spectra which depend on stellar labels suffer from label systematics which decrease model performance: the “stellar labels gap”. To close the stellar labels gap, we present a stellar label independent model for *Gaia* BP/RP spectra. We develop a novel implementation of a variational auto-encoder, which learns to generate an XP spectrum and accompanying ‘scatter’ without relying on stellar labels. We demonstrate that our model achieves competitive XP spectra reconstructions in comparison to stellar label dependent models. We find that our model learns stellar properties directly from the data itself. We then apply our model to XP/APOGEE giant stars to study the $[\alpha/M]$ information in *Gaia* XP. We provide strong evidence that the XP spectra contain meaningful $[\alpha/M]$ information by demonstrating that our model learns the α -bimodality, without relying on stellar label correlations for stars with $T_{\text{eff}} < 5000$ K, while also being sensitive to the anomalous abundances of *Gaia*-Enceladus stars. We publicly release our trained model, codebase and data. Importantly, our stellar label independent model can be implemented for any/all XP spectra because our model performance scales with training object density, not training label density.

Keywords: methods: data analysis - techniques: spectroscopic - stars: abundances - stars: fundamental parameters

1. INTRODUCTION

The advent of increasingly large-scale spectroscopic surveys, such as APOGEE (Majewski et al. 2017), LAMOST (Yan et al. 2022), GALAH (Buder et al. 2021), and most recently *Gaia* BP/RP (XP; *Gaia* Collaboration 2022), has motivated astronomers to develop data-driven models to cope with the massive influx of data (e.g. Ness et al. 2015; Ting et al. 2019; O’Brien et al. 2021; Leung & Bovy 2019; Zhang et al. 2023; Leung & Bovy 2023; Li et al. 2023). These data-driven techniques often seek to discern stellar properties from stellar spectra and/or generate stellar spectra from stellar properties.

Typically, data-driven methods are differentiated from physics-driven methods in the following way: A model which relies on synthetic stellar spectra is physics-driven, in the sense that it explicitly attempts to match theoretical spectra to observations. Conversely, a model is data-driven if it does not rely on theoretical spectra.

Generally, data-driven models incorporate some sort of machine learning algorithm to shed theoretical modeling.

Data-driven models have begun to be favoured over their physics-driven predecessors due to a disconnect between theory and practice known as the “synthetic gap.” The synthetic gap is a combination of theoretical systematics and instrumental effects which in tandem produce discrepancies between synthetic and observed stellar spectra, including but not limited to: one-dimensional modeling, assumption of hydrostatic equilibrium and local thermal equilibrium, and telluric lines. Although significant progress has been achieved in stellar atmosphere modeling, e.g. three-dimensional non-local thermal equilibrium models (Bergemann 2014), these simplified modeling assumptions are still frequently adopted when determining stellar properties of large-scale surveys. For a more comprehensive overview, see the Introduction of O’Brien et al. (2021).

However, it is seldom emphasized that any data-driven model which relies on stellar labels, a term which typically refers to effective temperature T_{eff} , surface gravity $\log g$, metallicity $[M/H]$ (and occasionally α -abundance $[\alpha/M]$), is *implicitly* physics-driven. Indeed, stellar labels utilized during training of such a data-driven model are typically estimated from theoretical stellar spectra. This is not strictly the case: e.g. T_{eff} and $\log g$ for *Gaia* FGK benchmark stars were determined with angular diameter measurements and bolometric fluxes (Heiter et al. 2015). However, such samples contain nowhere near the requisite number of stars for training, although useful for validation (34 *Gaia* FGK benchmark stars). Furthermore, this ab initio physics-driven estimation is often times several generations removed from the data-driven model being implemented, as it is becoming increasingly more common to ‘train machine learning on machine learning.’ In this context: training a machine learning model on stellar labels which were themselves obtained from a machine learning model. It has been demonstrated that training large language models on synthetic data produces a decrease in output diversity which worsens with each successive ‘ML on ML’ iteration (Guo et al. 2023). A similar concern for stellar properties is therefore warranted.

We therefore introduce the concept of an additional gap: the “stellar labels gap,” encompassing stellar label systematics which negatively impact performance of data-driven models which rely on stellar labels (stellar label dependent models). The stellar labels gap includes:

- (i) Poorly estimated stellar labels
- (ii) Regions of stellar label space where labels are insufficient summary statistics for a spectrum
- (iii) Regions of stellar label space with a dearth of labels to train on
- (iv) Stellar multiplicity¹ (binaries, triples, etc.)

These stellar label systematics then lead to systemic bias in stellar label dependent model predictions. In order to close the stellar labels gap, astronomers should develop purely data-driven models which do not rely on stellar labels (stellar label independent models).

To that end, this work presents a stellar label independent model for *Gaia* XP spectra: an unsupervised learning model which applies data compression. We develop a novel implementation of a variational auto-encoder (VAE); a *scatter* VAE, which learns to generate an XP spectrum while simultaneously estimating intrinsic

scatter for individual XP spectra. We demonstrate certain advantages of our stellar label independent approach by comparing our model performance to stellar label dependent models. Subsequently, we interpret the behaviour of model by contrasting stellar label space to our latent space. We then apply our model to the high- and low- α sequences to provide stellar label independent evidence that the *Gaia* XP spectra contain meaningful $[\alpha/M]$ information.

The subsequent Sections of this paper are organized as follows: Section 2 briefly reviews the data used in this work, namely *Gaia* XP spectra and APOGEE stellar labels. Section 3 presents our stellar label independent model architecture as well as our model training procedure. Section 4 compares our trained model performance to stellar label dependent models and interprets what our model has learned about the XP spectra. We then use our model, in Section 5 to conclusively demonstrate that the *Gaia* XP spectra contain meaningful $[\alpha/M]$ information, without the well known issue of α -abundance correlations with stellar labels, for stars with $T_{\text{eff}} < 5000$ K. Finally, in Section 6 we conclude by discussing the implications of our work in the context of stellar label in/dependent modeling of stellar spectra, its limitations, and promising future applications.

2. DATA

In this Section, we review the data used in this work: the *Gaia* XP spectra, as well as the APOGEE derived stellar labels we use in order to (i) compare our stellar label independent model to stellar label dependent models and (ii) interpret the behaviour of our model.

2.1. *Gaia* low-resolution BP/RP spectra

The *Gaia* low-resolution XP spectra in *Gaia* Data Release 3 (GDR3, Gaia Collaboration 2022) is comprised of 220+ million flux-calibrated, low-resolution spectra (De Angeli et al. 2022; Montegriffo et al. 2022). These spectra combine measurements from the Blue Photometer (BP) and Red Photometer (RP) *Gaia* instruments which span 330-680 and 640-1050 nm, respectively.

We use the XP coefficient spectra in this work. XP spectra deviate from traditional spectra because they are typically reported in Hermite polynomial space. Specifically, the BP and RP spectra are transformed from *discrete* wavelength-space into *continuous* coefficient-space: 110 coefficients which weight a set of Hermite polynomials (55 blue and 55 red coefficients). One can transform coefficients into fluxes, and vice-versa, ideally without loss of information. Previous machine learning models have been successfully implemented when using data in both XP wavelength space

¹ While in principle stellar label models can account for binarity, in practice this is rarely implemented.

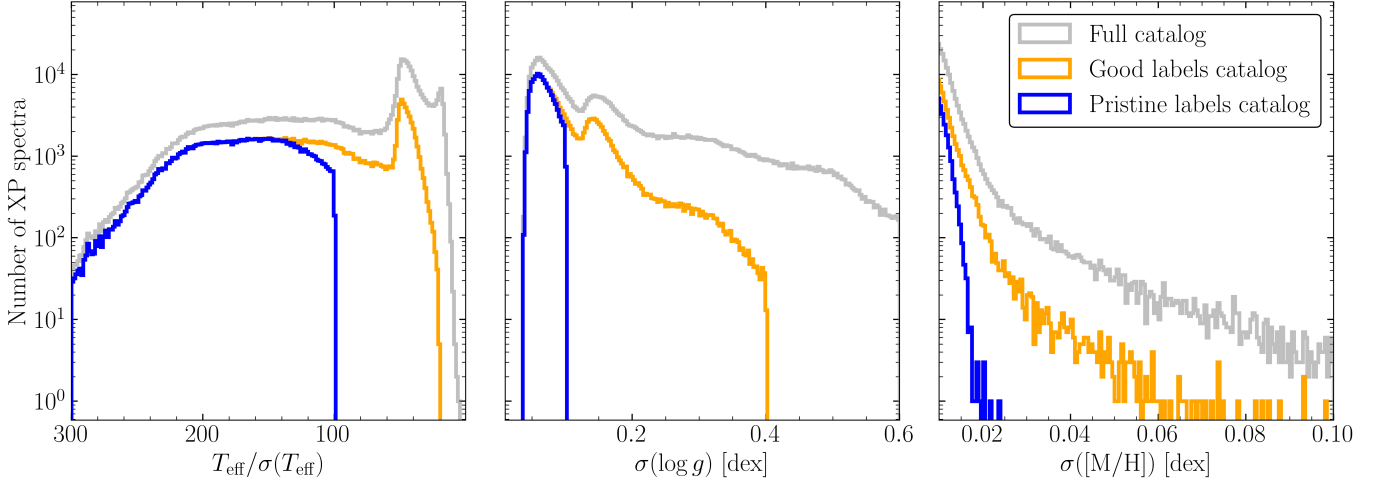


Figure 1. Signal-to-noise ratio/error distributions over stellar labels for the full, good labels and pristine labels catalogs defined in Section 2.2. T_{eff} SNR, $\log g$ uncertainty and $[M/H]$ uncertainty from the *Astro-NN* value-added APOGEE catalog are presented in left, middle and right panels, respectively. From the full, to the good labels, to the pristine labels catalog, quality cuts become increasingly restrictive. The full catalog is used to train our model, the good labels catalog is used to interpret our model behaviour, and the pristine labels catalog is used to investigate α -abundance information in the XP spectra.

and coefficient space. For instance, [Leung & Bovy \(2023\)](#) train their large astronomy model in XP coefficient space, whereas [Zhang et al. \(2023\)](#) train their deep stellar label model in XP wavelength-space.

2.2. Catalogs

To construct our train/test datasets, we perform a cross-match between the GDR3 XP spectra and the *astroNN*² value-added stellar label catalog for APOGEE Data Release 17 (DR17, [Leung & Bovy 2019](#); [Abdurro'uf 2022](#)). In principle, our model can be trained on the entire XP dataset, since our model does not require stellar labels. However, a primary focus of this work is to compare our stellar label independent model to stellar label dependent models, as well as compare generative spaces: stellar label space versus latent space. As such, we restrict ourselves to the cross-matched XP/APOGEE dataset. For the purposes of this work, we construct three different catalogs from the XP/APOGEE cross-match: the full catalog, the good labels catalog and the pristine labels catalog:

- The *full catalog* is the complete XP/APOGEE cross-match with no quality cuts whatsoever, which contains 502,311 stars. This is the dataset we use to train our model. Stellar label dependent models almost exclusively require quality cuts on training labels, which inherently limits the amount of training data. However, our stellar label inde-

pendent model can train on XP spectra which have poor APOGEE stellar labels.

- The *good labels catalog* is the XP/APOGEE cross-match restricted to APOGEE labels with high signal-to-noise ratio (SNR), which contains 202,970 stars (approximately 40% of the full catalog). We obtain the good labels catalog with the following quality cuts:
 - (i) $T_{\text{eff}}/\sigma(T_{\text{eff}}) > 30$
 - (ii) $\sigma(\log g) < 0.4$
 - (iii) $\sigma([M/H]) < 0.2$
 - (iv) $0 < BP - RP < 4$
 - (v) $6 < G < 17.5$
 - (vi) No bit set in **STARFLAG**
 - (vii) No bit 19 (**M.H.BAD**) or bit 23 (**STAR.BAD**) in **ASPCAPFLAG**
- The *pristine labels catalog* is created from analogous cuts to the good labels catalog, except the stellar label cuts are even more restrictive:
 - (i) $T_{\text{eff}}/\sigma(T_{\text{eff}}) > 100$
 - (ii) $\sigma(\log g) < 0.1$
 - (iii) $\sigma([M/H]) < 0.05$

The pristine stellar label cuts, in combination with the same color, magnitude and flag cuts from the good labels catalog, yield 123,804 stars (approximately 25% of the full catalog).

We depict the stellar label cuts which define the full, good labels and pristine labels catalogs in Figure 1. For the good labels catalog, the effective temperature (surface gravity, metallicity) cuts remove 16,028 (43,705,

² <https://github.com/henrysky/astroNN>

32,314) stars. For the pristine labels catalog, the effective temperature (surface gravity, metallicity) cuts remove 273,466 (286,191, 33,769) stars.

We compare the Kiel diagrams of the three catalogs in Figure 2. The full catalog appears to have poorly estimated labels at both ends of the main sequence, where surface gravities are systematically overestimated, and the metallicity gradient along the giant branch is somewhat obfuscated. Conversely, the good labels catalog has both a ‘cleaner’ main sequence and ‘sharper’ giant branch metallicity gradient. The pristine catalog is then exclusively composed of giant stars with very accurate stellar labels. Finally, apart from the stellar label cuts defining the pristine catalog removing all main-sequence stars, we do not observe major differences between the apparent magnitude and *Gaia* color distributions across our catalogs, beyond the color and magnitude cuts we impose.

2.3. Preprocessing

Before training our model on the full catalog, we perform the following pre-processing on the XP coefficient spectra. We first adopt the now common practice of concatenating the 55 BP and 55 RP coefficients into a single vector of length 110. Second, we normalize each individual spectrum by its corresponding *Gaia* *G*-band mean flux. Finally, we standard normalize the XP spectra coefficient by coefficient to zero mean and unit variance. Recently, Zhang et al. (2023) implemented median/quantile normalization for XP preprocessing to reduce the negative impact of outliers during training. However, one of the strengths of our model is the ability to incorporate outliers through our star-by-star scatter estimation, described in Section 3.2. As such, we opt for standard normalization.

3. METHOD: STELLAR LABEL INDEPENDENT GENERATIVE MODEL

This work presents a stellar label independent model which generates XP spectra. Specifically, we develop a novel implementation of a variational auto-encoder (VAE): *scatter* variational auto-encoder (*sVAE*). Before describing our *sVAE* architecture, we briefly review the concept of a (V)AE.

3.1. Variational auto-encoder review

An AE, which needs not be variational, can be thought of as a non-linear generalization of Principal Component Analysis. An AE begins with an encoder which compresses input data, in our case a stellar spectrum, to a low-dimensional latent representation. A decoder then attempts to reconstruct the stellar spectrum from the

latent representation. Ideally, this will allow the latent representation to learn key features which are shared across a set of stellar spectra. The variational nature of a VAE is added to an AE by upgrading the latent space from a collection of discrete points to a latent distribution \mathcal{Z} . The most popular VAE methodology is that of Kingma & Welling (2013), who encode input data onto an independent multivariate Gaussian distribution. The latent space can therefore be entirely characterized by a latent mean vector μ and variance vector σ' , with latent space dimension $n_{\mathcal{Z}}$.

3.2. Scatter variational auto-encoder

We present the high-level architecture of our model in Figure 3. Our *sVAE* differs from a traditional VAE since, in addition to a decoder (green) which estimates stellar spectra, we introduce a second ‘decoder’ which estimates intrinsic scatter *on a star-by-star basis* (purple). Importantly, both the XP reconstruction and XP scatter estimate are generated from the same latent space (red). After discarding the encoder (orange) post-training, new XP coefficients (and scatter) can be generated given an arbitrary latent space vector.

The input of our encoder is a pre-processed XP spectrum. The XP coefficients are fed through five intermediate layers, composed of 90, 70, 50, 30 and 10 neurons, respectively. All intermediate layers are activated by the gaussian error linear unit (GELU). The latent parameters are then given by a linear transformation of the final intermediate layer. In this work, we have fixed the latent space to 6 dimensions, meaning the encoder produces 12 outputs (6 means and variances). This choice was informed by initial experimentation with latent dimensions ranging from 1 to 20 dimensions. We found that a 6 dimensional latent space produced an optimal balance between interpretability and reconstruction error. Note that as is typical with VAEs, the reparametrization trick randomly samples $\epsilon \sim \mathcal{N}(0, 1)$ (standard normal distribution) to transform the latent mean μ and variance σ vectors into a latent space vector via $z = \epsilon\mu + \sigma$. Our decoder is the mirror image of our encoder, and takes as input a vector drawn from the latent distribution. We then reconstruct an XP spectrum, by feeding the latent vector through five intermediate layers analogous to the encoder, except in reverse order. Finally, a reconstruction of the 110 XP coefficients is produced with a linear transform. The scatter estimator has the same architecture as the decoder, except we enforce positivity with a final Sigmoid activation. Importantly, weights and biases of the scatter estimator are entirely disconnected from the decoder.

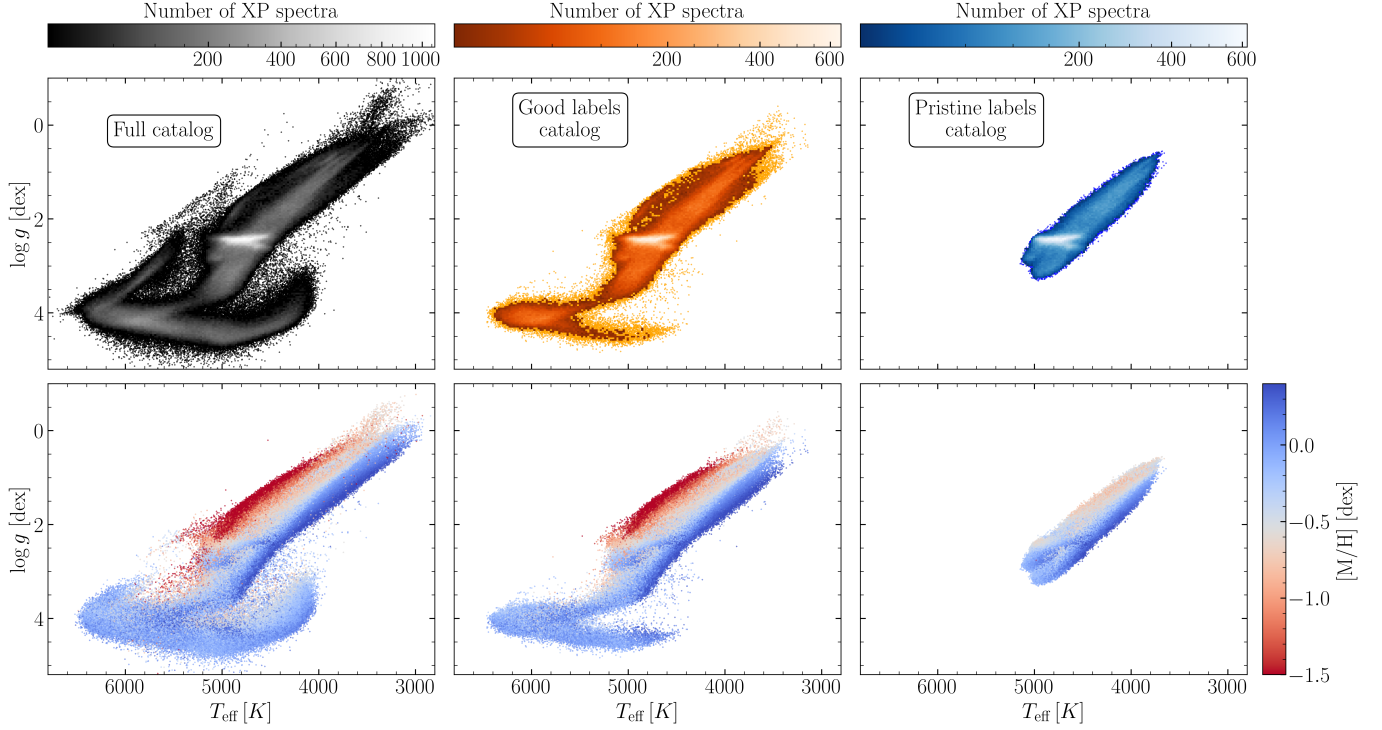


Figure 2. Kiel diagrams for the full (left), good labels (middle) and pristine labels (left) catalogs, colored by XP spectra number density (top row) and metallicity (bottom row). The increasingly restrictive quality cuts for the good labels and pristine labels catalog relative to the full catalog shrink the stellar label space they cover.

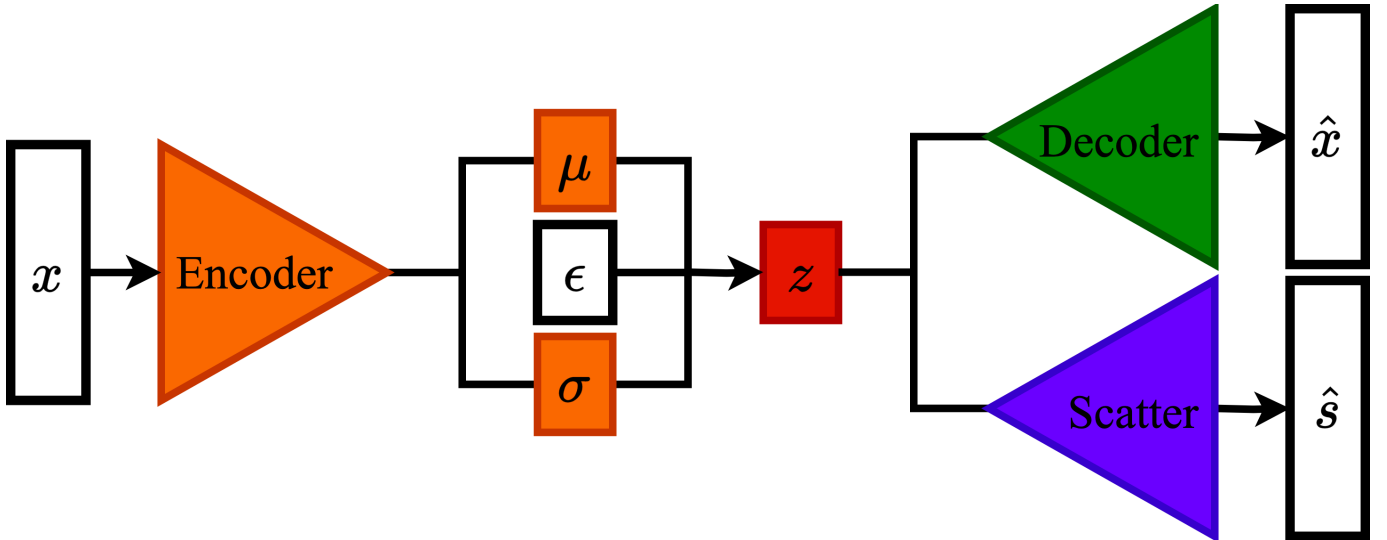


Figure 3. The *scatter* VAE architecture: An XP spectrum x is compressed into latent means μ and latent variances σ . The re-parametrization trick through ϵ then projects the encoding into the latent space z . Subsequently, the decoder reconstructs an estimate of the XP spectrum \hat{x} , and the scatter decoder produces a scatter estimate for the spectrum \hat{s} . During training, x and \hat{x} are encouraged to be as similar as possible, while simultaneously enforcing Gaussian structure in the latent space.

The interpretation of what the scatter estimator represents is by no means straightforward. The scatter estimator does not produce an estimate of intrinsic scatter in the traditional sense: variance of the entire XP dataset assuming zero measurement error. Rather, the

scatter could be both ‘intrinsic’ to an individual star: stellar variability arising from stellar spots, pulsations, etc. or simply be anomalous relative to the training set, but also ‘extrinsic’: arising from either (i) information loss through our data compression procedure or

(ii) *Gaia* systematics, poor observations or underestimation of uncertainties. As such, it is more accurate to think of the scatter estimate as an error term which includes traditional intrinsic scatter, systematics, outliers, etc. Empirically, we noticed a marked increase in performance when including the scatter estimator in our model. Therefore, despite not being able to pinpoint precisely which of the above aspects the scatter estimate mitigates, we can assert that our *sVAE* architecture does improve spectral modeling.

3.3. Training

After data pre-processing (see Section 2.3), we randomly split the full catalog into 90% for training, and allocate the remaining 10% for test data. We adopt a similar training approach to Leung & Bovy (2023), and use the AdamW optimizer (Kingma & Ba 2014; Loshchilov & Hutter 2017) with Cosine Annealing with Warm Restarts (Loshchilov & Hutter 2016) as our learning-rate scheduler. Cosine Annealing is a non-monotonic learning rate scheduler with cosine-shaped cycles. We select an initial cycle learning rate of 10^{-4} and a final learning rate of 10^{-10} . We train our *sVAE* for 5000 epochs, with 10 cycles of 500 epochs. As in Leung & Bovy (2023), we find that a batch size of 1024 is ideal for balancing training speed and model performance.

We denote the j^{th} coefficient of the i^{th} XP spectrum as x_{ij} , where $j = 1, \dots, 110$ with accompanying uncertainties σ_{ij} , and i spans the size of our training set. Furthermore, we denote the *sVAE* estimate of a coefficient (from the decoder) for a given XP spectrum as \hat{x}_{ij} , with accompanying scatter \hat{s}_{ij} (from the scatter estimator). During training, we aim to minimize the following loss function to optimize our *sVAE* model parameters:

$$\mathcal{L} = \tilde{\chi}^2(x, \hat{x}, \sigma, \hat{s}) + D_{\text{KL}}(\mu, \sigma'). \quad (1)$$

$\tilde{\chi}^2$ is the reconstruction loss between the input XP spectrum x and the reconstructed spectrum \hat{x} , whilst incorporating observational uncertainties σ and intrinsic scatter \hat{s} . In Eq. (1), the $\tilde{\chi}^2$ term is given by

$$\tilde{\chi}^2 = \chi^2(x, \hat{x}, \sigma, \hat{s}) + P(\sigma, \hat{s}). \quad (2)$$

The first term in Eq. (2) is the traditional (reduced) χ^2 , given by

$$\chi^2(x, \hat{x}, \sigma, \hat{s}) = \frac{1}{110N} \sum_{i=1}^N \sum_{j=1}^{110} \frac{(x_{ij} - \hat{x}_{ij})^2}{\sigma_{ij}^2 + \hat{s}_{ij}^2}, \quad (3)$$

and the second is a penalty term P given by

$$P(\sigma, \hat{s}) = \frac{1}{110N} \sum_{i=1}^N \sum_{j=1}^{110} \log(\sigma_{ij}^2 + \hat{s}_{ij}^2), \quad (4)$$

The second term in Eq. (1) is the latent space structure loss, for which we select the KL divergence (Kullback & Leibler 1951), and is given by

$$D_{\text{KL}}(\mu, \sigma') = \frac{1}{6N} \sum_{i=1}^N \sum_{j=1}^6 \left[\mu_{ij}^2 + \sigma'_{ij}{}^2 - \left(1 + \log \sigma'_{ij}{}^2 \right) \right], \quad (5)$$

where μ (σ') are the latent means (variances) and we are summing over the 6 latent space dimensions. In the above form, D_{KL} is not particularly intuitive. For the purposes of latent space structure loss, it can be thought as a distance between probability distributions. $D_{\text{KL}} = 0$ if two distributions are identical and $D_{\text{KL}} > 0$ otherwise. Hence, assuming a multivariate normal latent distribution, the KL divergence is a measure of the Gaussianity of the *sVAE* latent space. Note that no z variables directly appear in Eq. (5) due to marginalizing out z from the joint probability distribution $P(z|\mu, \sigma')$.

Finally, we emphasize that Eq. (1) does not account for the full *Gaia* XP coefficient covariance matrix. Our decision to neglect covariances during training is based on computational feasibility. Specifically, the inclusion of the full covariance matrix would have lead to an increased computational cost due to both loss function evaluation and the scatter output necessarily being augmented to two dimensions (to match the covariance matrix) if we took this approach.

3.4. Data and code

Our *sVAE* model, implemented in PyTorch, can be trained on a single NVIDIA RTX4070 Ti GPU in ~ 6 hours, with our current XP/APOGEE cross-match of $\sim 450,000$ training objects³. Our model can project $\sim 10^5$ XP spectra into the latent space, and simulate XP spectra from the latent space in a few seconds. Our trained stellar label independent model and codebase are publicly available at GitHub⁴ for others to reproduce our results, build upon our existing model and apply our model to XP spectra beyond the XP/APOGEE cross-match. A copy of AlexLaroche7/xp-vae was deposited to Zenodo: doi:10.5281/zenodo.14041978 (Laroche 2024). Furthermore, all data associated with this work is available on Zenodo: doi:10.5281/zenodo.14041773.

³ Linearly extrapolating this training time to the full *Gaia* DR3 XP dataset of 220 million spectra, we roughly estimate that training on all currently available XP spectra would take ~ 120 GPU days (on the same GPU, with the same batch size, architecture, etc.). To give some perspective, a foundation language model such as LLaMA (Touvron et al. 2023) takes 21 days \times 2048 GPUs = $\sim 43,000$ GPU days to train.

⁴ <https://github.com/AlexLaroche7/xp-vae>

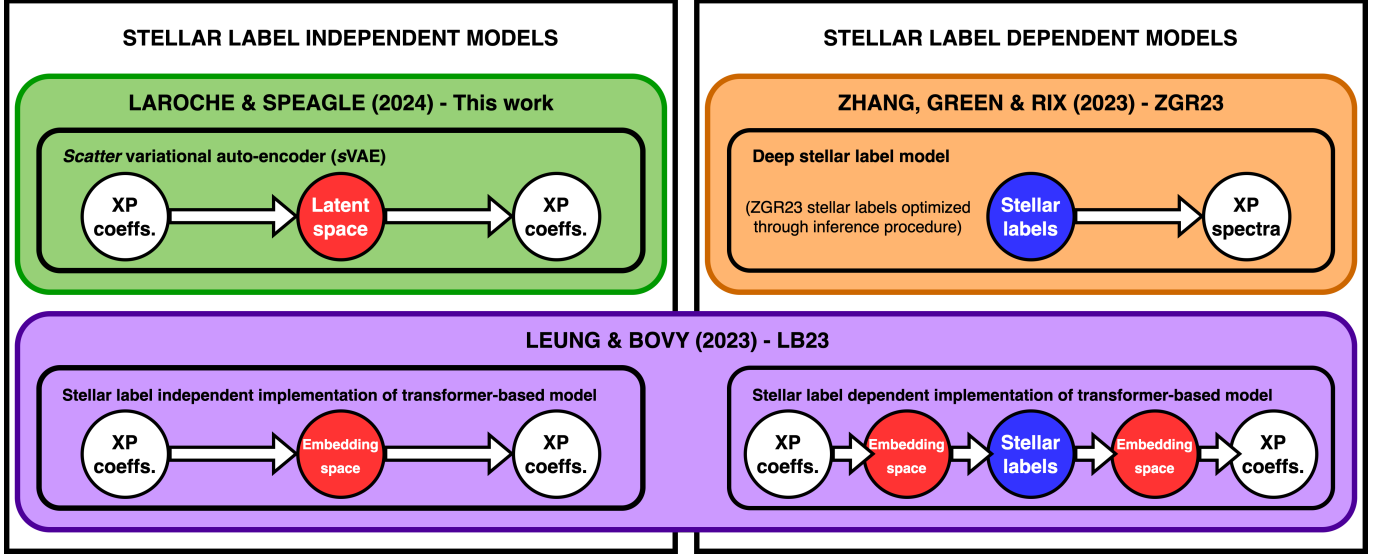


Figure 4. A high-level overview of the models used in this work. We classify stellar independent models (*left*) and stellar dependent models (*right*) by whether or not stellar labels are used to represent an XP spectrum. Our sVAE (green) projects an XP coefficient spectrum into the latent space, and subsequently reconstructs the coefficients. The stellar label independent implementation of LB23 (purple, *left*) projects an XP coefficient spectrum into the embedding space, and subsequently reconstructs the coefficients. On the other hand, ZGR23 (orange) generates an XP spectrum, in wavelength space, from stellar labels which were optimized through their inference procedure. Finally, the stellar label dependent implementation of LB23 (purple, *right*) projects an XP coefficient spectrum into the embedding space, estimates stellar labels, then re-projects the labels into the embedding space to ultimately reconstruct the coefficients.

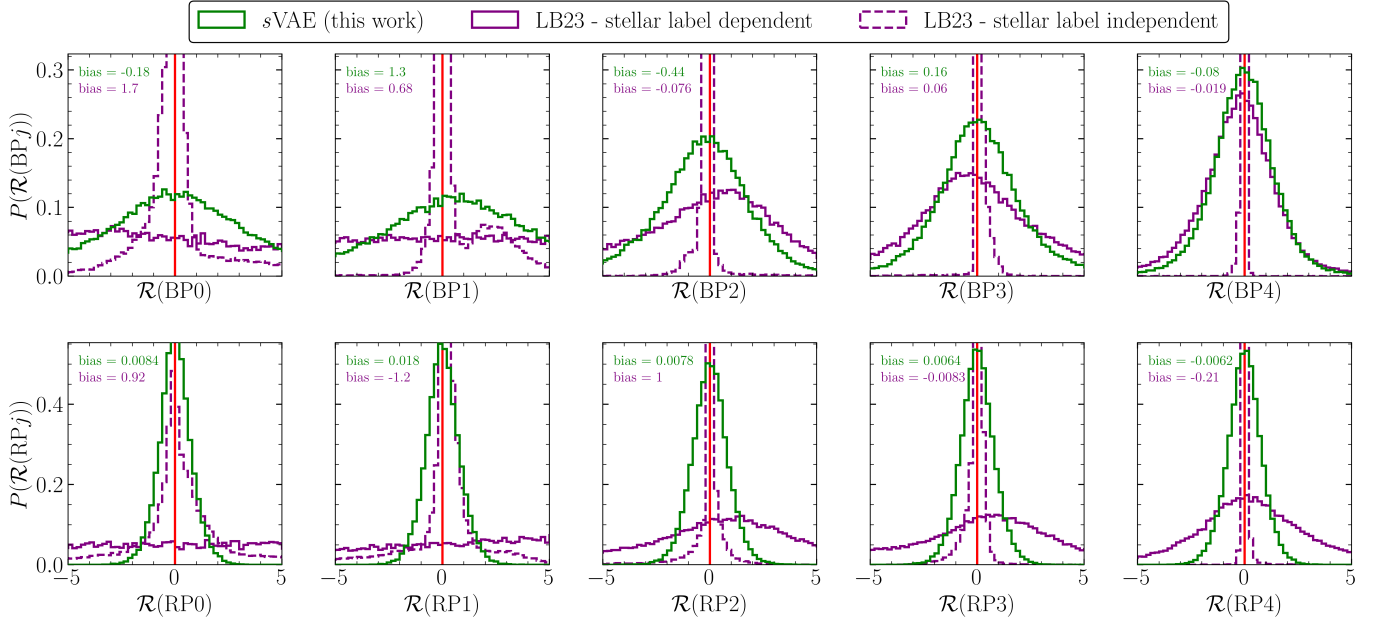


Figure 5. Relative coefficient reconstruction error distributions, $\mathcal{R}(j)$ from Eq. (6), for the stellar label dependent (solid purple) and independent (dashed purple) implementations of Leung & Bovy (2023) in comparison to our stellar label independent model (green), over test data in the full catalog for the first 5 BP and RP coefficients. Bias is presented for each coefficient for our sVAE and the stellar label dependent LB23 model. Our sVAE outperforms the stellar label dependent LB23 model, but underperforms relative to the stellar label independent LB23 model, which suffers from less information loss due to its larger embedding space (64 tokens in comparison to 6 latent variables).

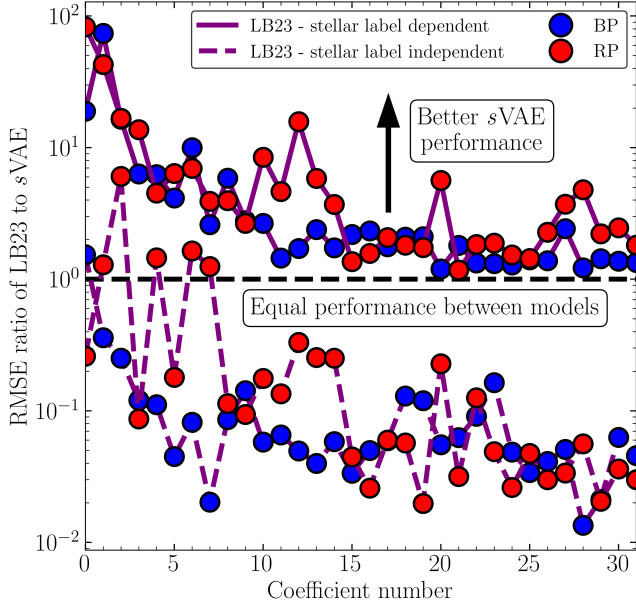


Figure 6. Ratio of spectrum reconstruction error as a function of BP/RP coefficient, $\text{RMSE}(j)$ from Eq. (7), for the stellar label independent and dependent implementations of LB23 relative to our stellar label independent model. Our $s\text{VAE}$ produces more accurate reconstructions of the lower order coefficients, and decays to negligible improvement for noise-dominated higher order coefficients, relative to the stellar label dependent LB23 model. Conversely, the stellar label independent LB23 model largely outperforms our $s\text{VAE}$ due to less information loss.

4. RESULTS

We now present the results of our trained $s\text{VAE}$. First, we compare our model performance to two models from the literature in Section 4.1. Second, we demonstrate the behaviour of our model by observing stellar label trends in the $s\text{VAE}$ latent space in Section 4.2.

4.1. Model performance

To compare the performance of our stellar label independent approach to stellar label dependent models, we compare our model to two generative models for XP spectra (see Figure 4):

1. *Large astronomy model* (Leung & Bovy 2023, LB23): A transformer-based model which transfers the methods of Large Language Models (LLMs) to astronomical data: a ‘large astronomy model’ (LAM). LB23 can perceive/predict any combination of XP data or stellar labels. As such, we compare with LB23 by implementing their model in both a stellar independent *and* dependent way. Specifically:

- (i) *Stellar label independent*: We project XP coefficients into the embedding space, and subsequently reconstruct the coefficients.
- (ii) *Stellar label dependent*: We project the XP coefficients into the embedding space, estimate stellar labels from the embedding space, re-project the labels into the embedding space, and finally reconstruct the coefficients. The stellar labels for LB23 are: T_{eff} , $\log g$ and $[M/H]$ and $J - H$ and $J - K$ colours from 2MASS photometry (Skrutskie et al. 2006). Near-infrared photometry is included to provide a proxy for extinction by breaking the temperature-extinction degeneracy.

Note that the LB23 model has a context window length of 64: the maximum number of tokens their model can handle. As such, we only project/reconstruct the first 32 BP and RP coefficients (64 total) with their model.

2. *Deep stellar label model* (Zhang et al. 2023, ZGR23): A data-driven model which estimates stellar parameters. ZGR23 can also generate XP spectra from T_{eff} , $\log g$, $[M/H]$, extinction and distance. Note that because ZGR23 produce their own stellar parameter estimates, we use the ZGR23 stellar parameters catalog to generate XP spectra with their model. Additionally, ZGR23 advise that their results are only reliable for stars which pass their reliability cut, which we enforce by requiring `quality_flags` < 8 (see Section 4 for further discussion on the implications of their reliability cut).

We compare our model reconstruction errors over test data (10% of the full catalog) to errors from both stellar label independent models. LB23 and ZGR23 are similar in the sense that they both generate XP spectra from labels, but dissimilar in their output: LB23 predicts XP spectra in coefficient space, whereas ZGR23 predicts XP spectra in wavelength space. In order to compare to ZGR23, we convert our model XP coefficient predictions to wavelength space with `GaiaXPpy`.⁵

To compare our model to LB23, we also define the *relative coefficient reconstruction error* $\mathcal{R}(j)$ as the relative error for a single coefficient for a single star:

$$\mathcal{R}(j) \equiv \frac{x_{ij} - \hat{x}_{ij}}{\sigma_{ij}}, \quad (6)$$

⁵ <https://gaia-dpci.github.io/GaiaXPpy-website/>, version 2.1.0 (Ruz-Mieres & Kostrzewa-Rutkowska 2023)

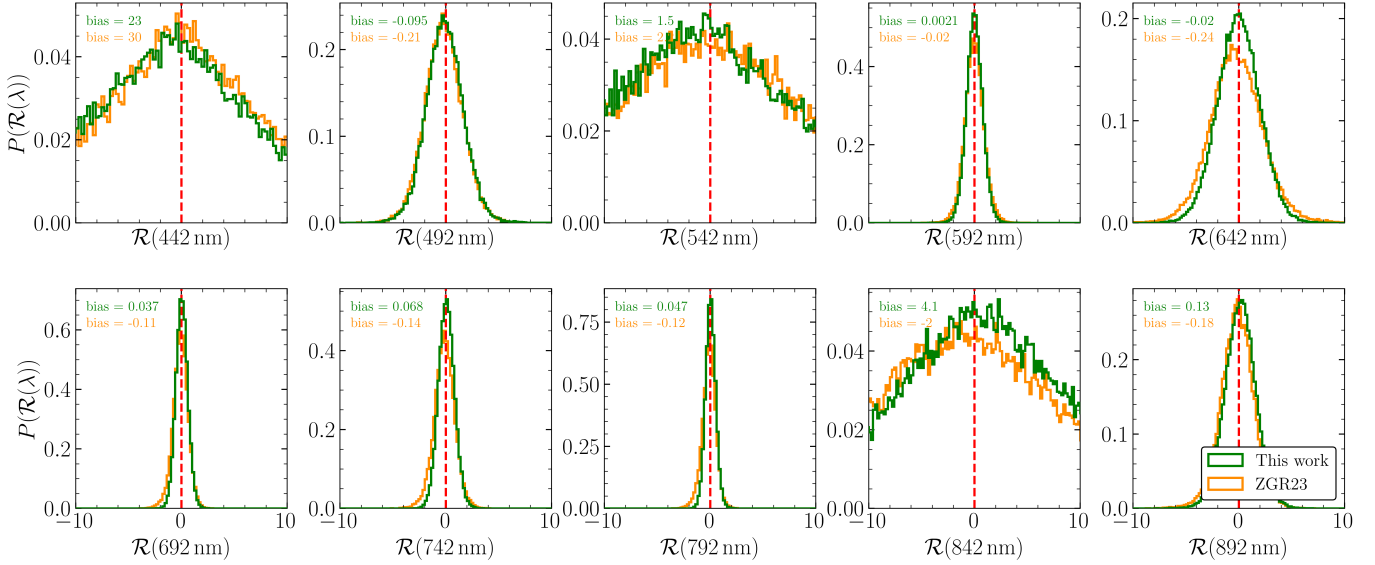


Figure 7. Relative wavelength reconstruction error distributions, $\mathcal{R}(\lambda)$ from Eq. (8), for the deep stellar label model of Zhang et al. (2023) (ZGR23, orange) in comparison to our stellar label independent model (green), over test data in the full catalog for 10 wavelengths uniformly distributed across the XP wavelengths. Bias is presented at each wavelength. Here, we apply the reliability cut of ZGR23 to both sets of reconstruction errors. From 600-800 nm, our *s*VAE produces less bias than ZGR23.

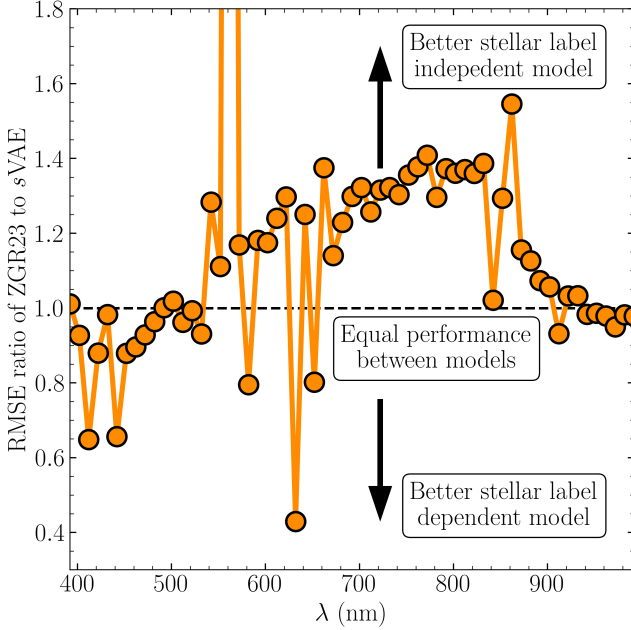


Figure 8. Ratio of spectrum reconstruction error as a function of BP/RP wavelength, $\text{RMSE}(\lambda)$ from Eq. (9), for ZGR23 relative to our stellar label independent model. Similar to Figure 7, we apply the ZGR23 reliability cut to both sets of reconstruction errors. From 550-900 nm, our *s*VAE produces more accurate wavelength reconstructions than ZGR23 over most wavelengths, by 20-40%.

where $\mathcal{R} \equiv \mathcal{R}(j)$ because the error is computed in XP coefficient space, over j . We also define the *spectrum reconstruction error* $\text{RMSE}(j)$ as the root mean square

error over coefficients for a single star:

$$\text{RMSE}^2(j) \equiv \frac{1}{64} \{ (x_{BP} - \hat{x}_{BP})^T C_{BP}^{-1} (x_{BP} - \hat{x}_{BP}) + (x_{RP} - \hat{x}_{RP})^T C_{RP}^{-1} (x_{RP} - \hat{x}_{RP}) \}, \quad (7)$$

where $\text{RMSE}(j)$ computes the reconstruction error while accounting for the XP covariance matrices. Specifically, x_{BP} and x_{RP} are vectors containing the first 32 BP and RP coefficients due to the length 64 context window of LB23. Similarly, \hat{x}_{BP} and \hat{x}_{RP} are the model reconstructions for the first 32 BP and RP coefficients, and C_{BP} (C_{RP}) are the 32×32 sub-matrices of the full BP (RP) covariance matrices.

Next, to compare our model to ZGR23, we define analogous quantities to the former two errors, but in XP wavelength space. We define the *relative wavelength reconstruction error* $\mathcal{R}(\lambda)$ as the relative error for a single wavelength for a single star:

$$\mathcal{R}(\lambda) = \frac{x_i(\lambda) - \hat{x}_i(\lambda)}{\sigma_i(\lambda)}, \quad (8)$$

where $\mathcal{R} \equiv \mathcal{R}(\lambda)$ because the error is computed in wavelength space, over λ . Lastly, we define the *spectrum reconstruction error* $\text{RMSE}(\lambda)$ as the root mean square error over wavelengths for a single star:

$$\text{RMSE}(\lambda) \equiv \sqrt{\frac{1}{61} \sum_{j=1}^{61} \mathcal{R}^2(\lambda)}, \quad (9)$$

where $\text{RMSE}(\lambda)$ is integrating the square of $\mathcal{R}(\lambda)$ over wavelengths. Note that $j \in [1, 61]$ because ZGR23 predict the flux at 61 XP wavelengths⁶.

In Figure 5, we compare the relative coefficient error distributions for our model to LB23 across the first five BP and RP coefficients, Eq. (6). Our relative error distributions are Gaussian with negligible bias, relative to the stellar label dependent LB23 model. Conversely, the stellar label independent LB23 model evidently outperforms our model over most coefficients. Additionally, in Figure 6, we compare the ratio of coefficient reconstruction errors over the full catalog for LB23 relative to our trained *s*VAE, Eq. (7), for both the stellar label independent and dependent implementations. We find that, for the lowest order XP coefficients, our model outperforms the stellar label dependent LB23 model by 1-2 orders of magnitude. Furthermore, the relative performance increase for our model decays to unity for higher order coefficients, which are the most noise dominated. Conversely, we observe that the stellar label independent LB23 model outperforms our *s*VAE by an order of magnitude for all but the lowest order coefficients.

We conclude that our *s*VAE outperforms the stellar label dependent model, but not the stellar label independent model, of LB23. This is expected, given the length 64 context window of LB23. In other words, the stellar label independent LB23 model is effectively learning a one-to-one mapping between coefficients, without data compression. In contrast, our model compresses the XP spectrum into 6 latent variables, and as such loses XP information relative to the LB23 stellar independent model.

Analogously, in Figure 7, we compare the relative wavelength error distributions for our model to ZGR23 across 10 wavelengths uniformly distributed throughout the XP spectra wavelength range, Eq. (8). Then, in Figure 8, we compare the ratio of wavelength reconstruction errors for ZGR23 relative to our model, Eq. (9). Here, we only compare errors for stars in the full catalog which satisfy the ZGR23 reliability cut, which decreases the test sample from 50,232 stars down to 37,302 stars. Our stellar label independent model produces less bias from 600-800 nm (in the red half of the XP spectra). Both the wavelength reconstruction errors and relative wavelength error distributions suggest that our model better reconstructs XP spectra for most wavelengths beyond the blue end, relative to ZGR23. By inspecting the model wavelength error ratios in Figure 8, we observe

that our model outperforms ZGR23 in flux reconstruction for most of the wavelengths from 550-900 nm by approximately 20-40%. However, for wavelengths at the extremities ZGR23 achieves equivalent or better performance than our *s*VAE, concentrated towards the blue end (< 500 nm). Here, it is informative to compare the APOGEE Kiel diagram in Figure 2 to the training sample of ZGR23 (Figure 1 in Zhang et al. 2023). We speculate that the larger number of hot stars in the ZGR23 training sample may contribute to our *s*VAE under-performing at bluer wavelengths. Lastly, there are certain anomalous wavelengths which do not appear to follow the general error ratio trend as a function of wavelength.

Finally, in Figure 9 we assess our model performance relative to LB23 and ZGR23 over test data in the full catalog, as a function of stellar labels. Importantly, here we do *not* apply the reliability cut of ZGR23 to our model error distributions, but rather only to ZGR23. Since our model is fully independent of stellar labels, we would expect to only observe error trends coming from training object density (panels (i)-(iii)), as opposed to genuine trends coming from the labels themselves, which is indeed the case. Our model errors increase towards the wings of the stellar label distributions (panels (iv)-(ix)) as the number of training objects decreases. First, our XP reconstruction errors in coefficient space are typically 0.5-1 orders of magnitude smaller than the stellar label dependent LB23 model across stellar label space, until training object density decays at the extremities. Conversely, the stellar label independent LB23 model outperforms our model across much of stellar label space. Second, we observe that our model is particularly more robust when applied to cool stars (panels (iv) and (vii)) and low surface gravity stars (panels (v) and (viii)). The ZGR23 error distributions below $T_{\text{eff}} \approx 4000$ K and below $\log g \approx 1$ dex are effectively meaningless because almost all stars in these regimes are removed by their quality cut. This can be observed by comparing the ZGR23 stellar label distributions (orange) to the XP/APOGEE cross-match stellar label distributions (green). We emphasize that this is an important advantage of our stellar label independent approach: not relying on stellar labels to simulate spectra allows our model to extend into regions of stellar parameter space where labels are unreliable. As such, with an appropriate training sample our model has the potential accurately reproduce both M-dwarf (low temperature) and giant (low surface gravity) BP/RP spectra.

To summarize, our stellar label independent model outperforms the stellar label dependent transformer-based model of LB23, but underperforms in comparison

⁶ GaiaXPy does not currently have the functionality to account for the full XP covariance matrices, and as such we only consider univariate errors in Eq. (9).

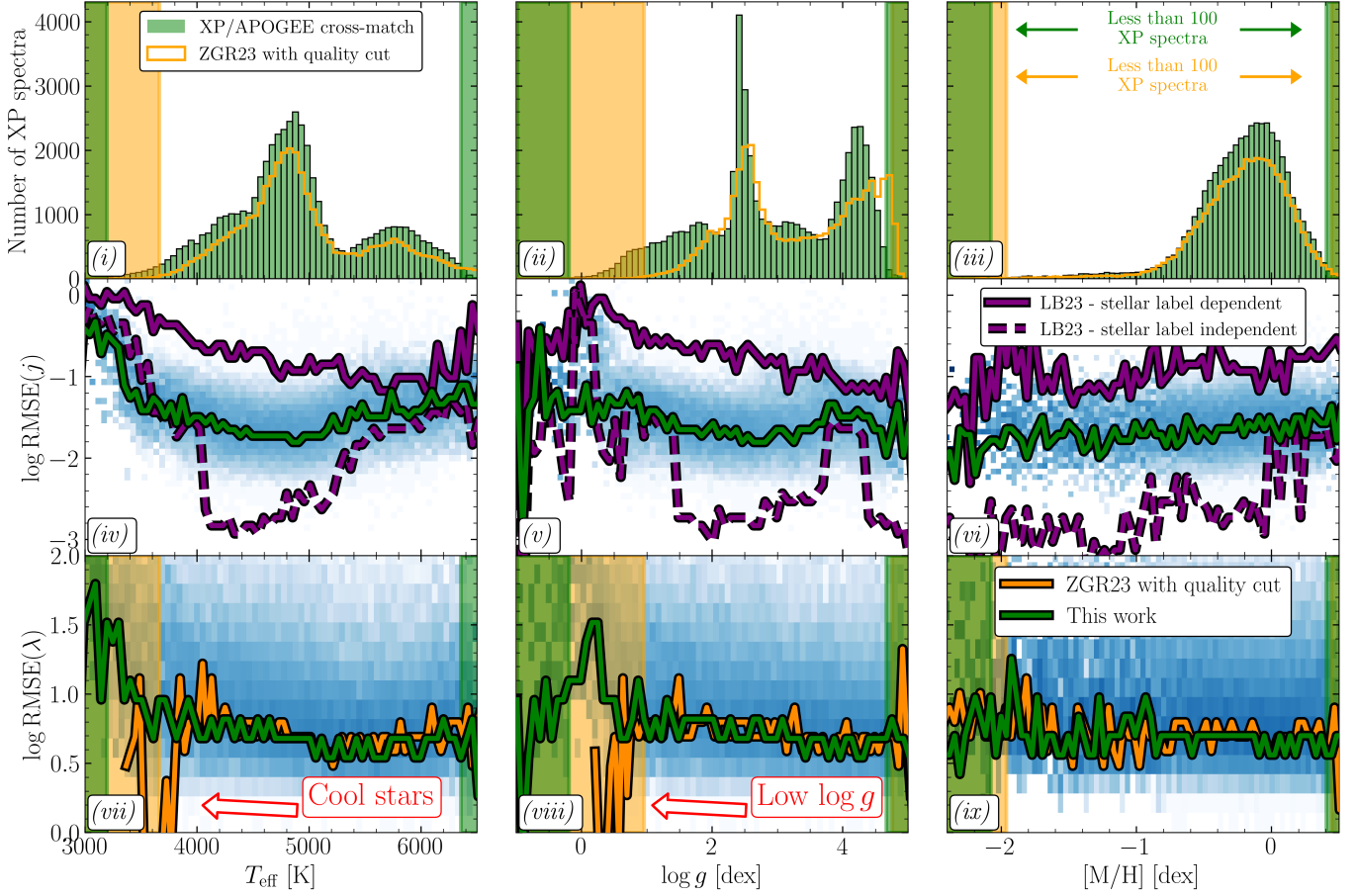


Figure 9. XP spectrum reconstruction errors for LB3 (middle row: (iv)-(vi), purple, coefficient space) and ZGR23 with their quality cut (bottom row: (vii)-(ix), orange, wavelength space) in comparison to our stellar label independent model (green), as a function of stellar labels (top row: (i)-(iii)) over test data in the full catalog. Here, we do not apply the ZGR23 quality cut to our wavelength reconstruction errors. Each model curve traces the median reconstruction error as a function of stellar labels. The blue colormap depicts our model error distributions, on log scale. Regions of stellar label space with less than 100 XP spectra are shaded in green (orange) for the test spectra (ZGR23 cross-match). Within these regions, stellar label error trends are unreliable. Our stellar label independent model can accommodate both cool (vii) and low surface gravity (viii) stars, in contrast to the stellar label dependent models (see Section 4.1 for further discussion).

to the stellar label independent LB23 model due to our comparatively harsh data compression. We also compare our model to an ‘expert’ stellar label dependent model: the deep stellar label model of ZGR23. Relative to ZGR23, we find that our stellar label independent model has specific advantages. Namely, simulating XP spectra from approximately 550-990 nm, and simulating cool stars below $T_{\text{eff}} \approx 4000$ K and low surface gravity stars below $\log g \approx 1$ dex.

4.2. Latent space vs. stellar label space

The major novelty of our *s*VAE model is its independence from stellar labels, opting instead to generate an XP spectrum from a latent space. This novelty is a strength, but also a potential weakness. On the one hand, in Section 4, we demonstrated that our stellar label independent model demonstrates increased perfor-

mance relative to stellar label dependent models. On the other hand, stellar label dependent models are inherently useful because stellar labels are grounded in astrophysical understanding. Two questions then naturally arise: ‘What astrophysical information has the latent space learned?’ and ‘Is the latent space representation of XP spectra useful?’. We will set aside the latter question and return to it in Sections 5 and 6, focusing for now on astrophysical interpretation of the latent space.

4.2.1. Kiel tracks

We will now argue that our trained *s*VAE learns a latent representation of the XP spectra which contains astrophysical information, using the good labels catalog.

First, we implement an approximate main-sequence/giant branch (MS/GB) cut in label (Kiel

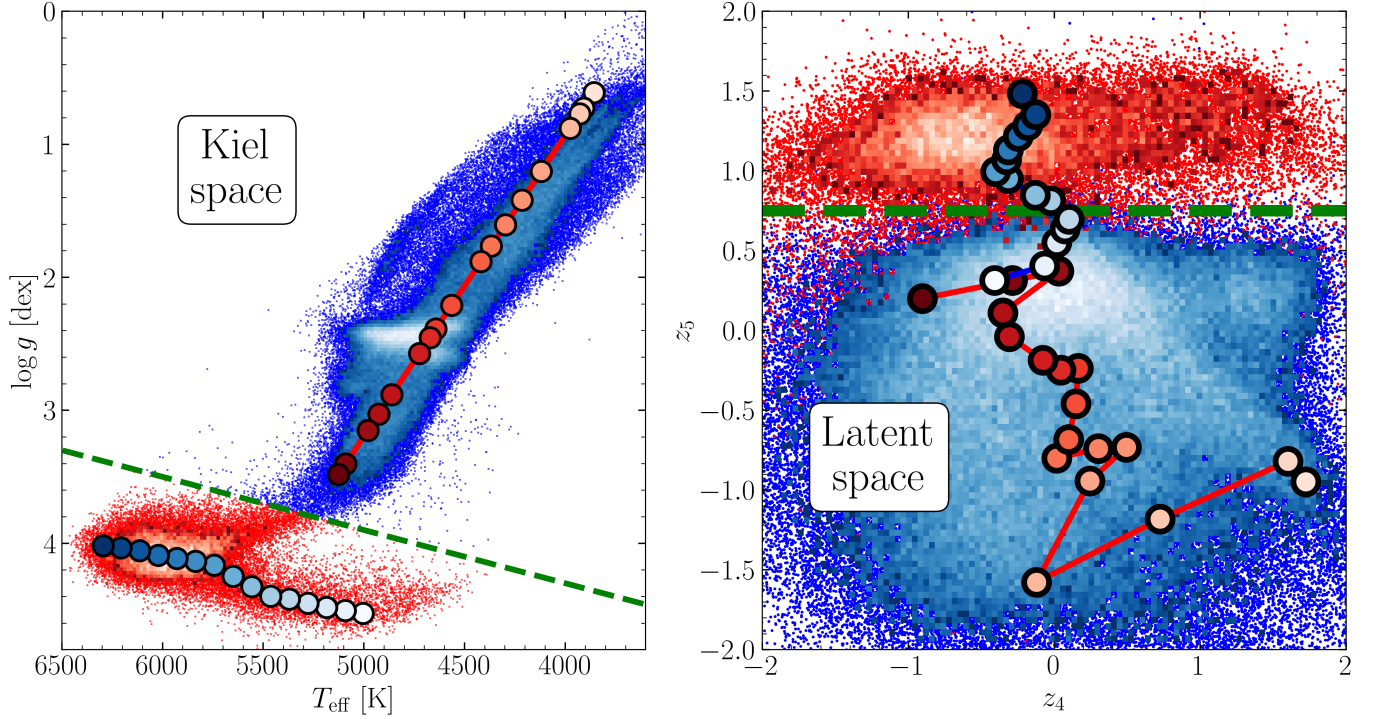


Figure 10. Main-sequence (MS, blue track) and giant branch (GB, red track) evolutionary tracks in Kiel space (left) and in our stellar label independent latent space (right). Additionally, approximate MS/GB (red color map, blue color map) boundaries are presented in both Kiel space and latent space (green dashed lines). Our stellar label independent latent space has learned the distinction between MS and GB stars, as well as the $T_{\text{eff}} - \log g$ relation, directly from the data itself.

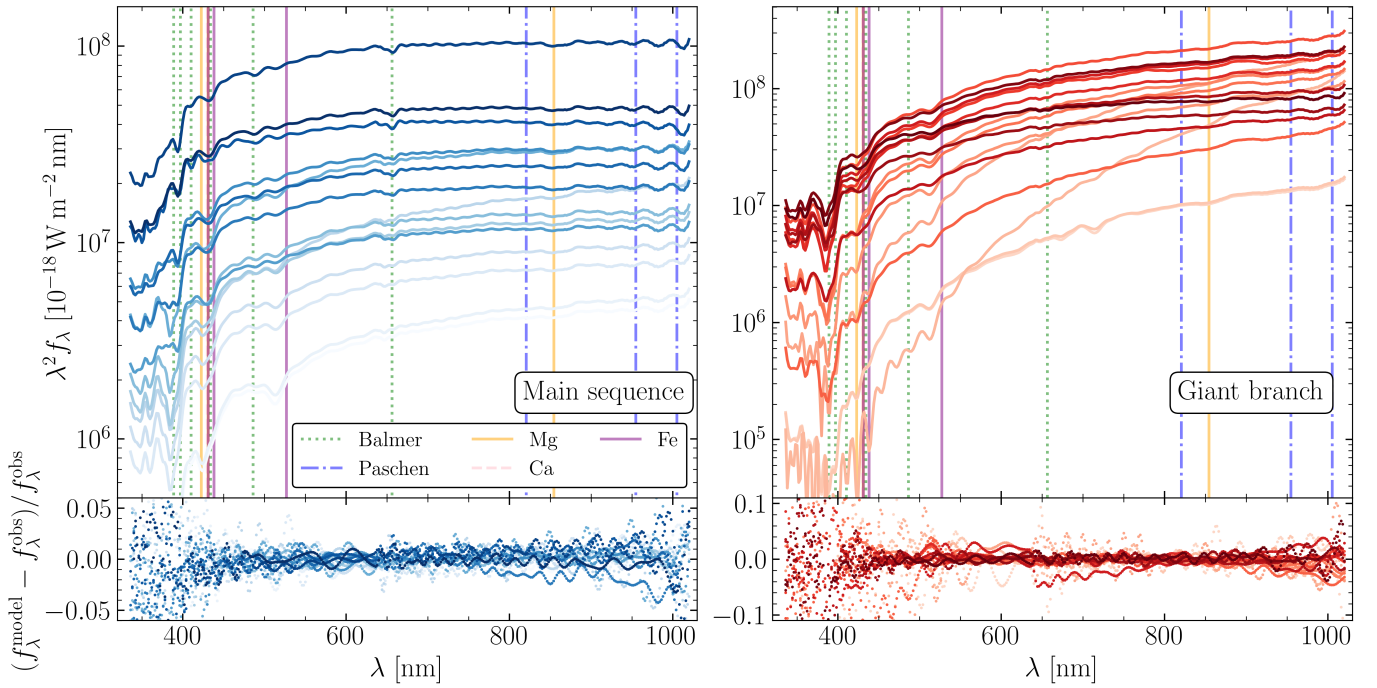


Figure 11. Model-predicted fluxes (at 1 kpc) for representative XP spectra, after de-normalization and scaling by *Gaia* G -flux, along the MS (left) and GB (right) tracks presented in Figure 10, along with percent difference relative to observations (bottom row). The Balmer and Paschen series, as well as Mg, Ca and Fe lines are shown for context.

space:

$$\log g = 5.9 - 0.4 \left(\frac{T_{\text{eff}}}{kK} \right), \quad (10)$$

visualized as the dashed green line in the left panel of Figure 10, to separate MS and GB stars. In our latent space (right panel of Figure 10), we observe that the z_5 latent dimension is effectively a MS/GB classifier: $z_5 \approx 0.75$ acts as an analogous boundary, above and below which the s VAE clusters MS and GB stars, respectively (dashed green line in the right panel of Fig 10).

Second, we produce MS and GB tracks in $T_{\text{eff}} - \log g$ (Kiel) space, at solar-metallicity. We use the Kiel relations which ZGR23, and subsequently LB23, adopted. For MS stars:

$$\log g_{\text{MS}} = \begin{cases} 4.6, & T_{\text{eff}}/kK < 5 \\ 7.1 - \left(\frac{T_{\text{eff}}}{2kK} \right), & 5 \leq T_{\text{eff}}/kK < 6.3 \\ 3.95, & T_{\text{eff}}/kK \geq 6.3, \end{cases} \quad (11)$$

and for giants:

$$\frac{T_{\text{eff,GB}}}{kK} = \begin{cases} 3.59 + 0.44 \log g & \log g < 3.65, \\ 0.09 + 1.4 \log g & \log g \geq 3.65. \end{cases} \quad (12)$$

With Eqs. (11) and (12), we bin and subsequently average XP spectra according to their APOGEE labels along the MS and GB relations in Kiel space to create tracks, shown in the left panel of Figure 10. We then project these spectra into our s VAE latent space, and average again to produce analogous tracks—*not as a function of stellar labels, but as a function of latent variables*, shown in the right panel of Figure 10. Here, we only present latent tracks in two of the most Kiel informative latent dimensions (for tracks across the full latent space, see Figure 15 in Appendix A). We observe that z_4 and z_5 function as pseudo-Kiel space, with both the MS and GB tracks exhibiting well-behaved, ordered trajectories through the latent space. That being said, it is apparent that the GB track is slightly more stochastic than the MS track. We speculate that this could arise from dust extinction which impacts the optical XP spectra, but not the APOGEE stellar labels derived from near-infrared spectroscopy. We conclude that our s VAE learns the $T_{\text{eff}} - \log g$ relation, despite the fact that the model has never ‘seen’ these stellar labels.

Additionally, in Figure 11 we present model reconstructed XP spectra in wavelength space for stars randomly drawn from within each bin used to define the MS and GB tracks presented in Figure 10. The percent difference for our reconstructions relative to observations are predominantly at the percent level, apart

from the bluest wavelengths where reconstructions deteriorate somewhat (for an in-depth discussion on model reconstruction errors, see Section 4.1).

4.2.2. Metallicity tracks

We have established that our s VAE learns classification of MS and giant stars and stellar evolutionary tracks along the MS and GB. But what about metallicity? Analogously to the tracks in Sect. 4.2.1, we first produce tracks at fixed metallicity in $T_{\text{eff}} - \log g$ space, and subsequently translate them into our s VAE latent space. In the left panel of Figure 12, we present binned Kiel tracks along both the MS and GB from high ($[M/H] \approx 0.5$) to low ($[M/H] \approx -2.2$) metallicity. Then, in the middle (right) panel of Figure 12 we project the MS (GB) metallicity tracks into our s VAE latent space, and present the three most metal-informative latent dimensions.

We observe evident metallicity gradients in the latent space, for both MS and giant stars. This is demonstrated by both the clear separation of fixed metallicity tracks in both the 2D latent space distributions and 1D cumulative distributions. Furthermore, z_1 is the most metallicity dominated latent dimension, with a strong metallicity gradient from metal-poor stars at $z_1 \ll 0$ to metal-rich stars at $z_1 \gg 0$. Interestingly, the main sequence metallicity tracks, which are crowded and overlapping in Kiel space, appear to be better separated in the latent space (along z_1).

In summary, our trained s VAE learns information about *Gaia* XP spectra which is effectively equivalent to stellar labels: T_{eff} , $\log g$ and $[M/H]$, without ever having ‘seen’ any of these labels during training. However, the latent space blends its understanding of stellar labels across multiple latent dimensions. The information our s VAE learns about *Gaia* XP spectra contains stellar label information, but can additionally learn information from stellar spectra which is not captured by stellar labels. Furthermore, our latent space can be used to probe the actual astrophysical information content of the *Gaia* XP spectra, without the ambiguity produced by stellar label correlations; an important consideration for α -abundance.

5. $[\alpha/M]$ INFORMATION IN THE GAIA XP SPECTRA

The chemical distribution of stars in the Galactic disk of our Milky Way contains important information about the formation, accretion and dynamical evolutionary histories of the Galactic disk. In particular, observations of Galactic disk stars have long been known to exhibit two components in $[M/H]$ - $[\alpha/M]$ space: the α -bimodality, comprised of the high- α and low- α sequences

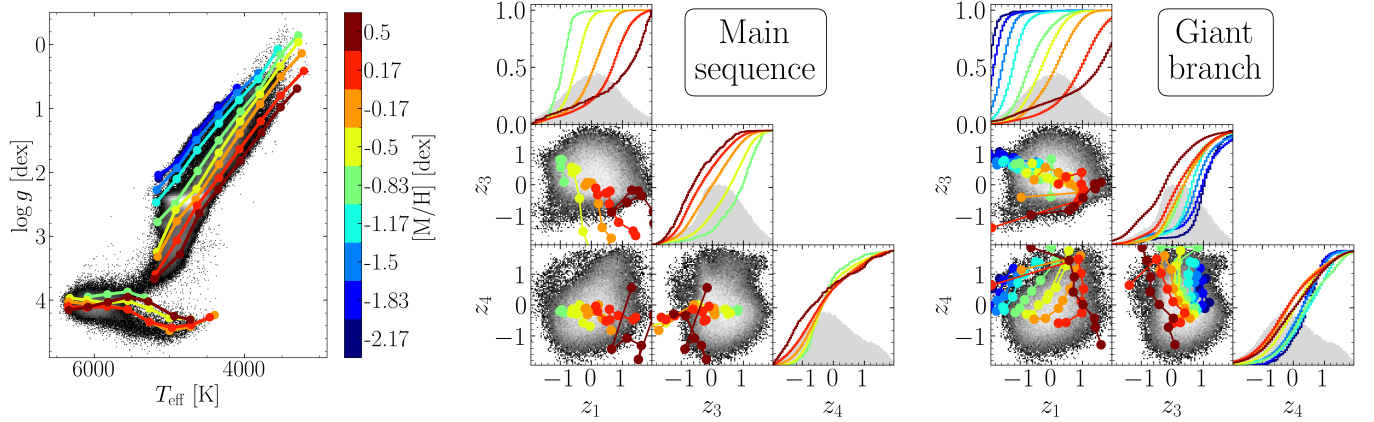


Figure 12. Fixed metallicity tracks along the MS and GB in Kiel space (left), projected into the latent space (MS middle, GB right). The three most metal informative latent dimensions are presented. In the latent space 1D marginal distributions, cumulative distribution at fixed metallicity are also presented. Metallicity gradients across the latent dimensions demonstrate that our stellar label dependent model has learned $[M/H]$ information directly from the data itself.

(Fuhrmann 1998; Bensby et al. 2003; Anders et al. 2014; Nidever et al. 2014; Hayden et al. 2015; Kordopatis et al. 2015). Initially, these sequences were used to separate the thin and thick disks (Yoshii 1982; Gilmore & Reid 1983). The thick (thin) disk is mostly comprised of the high- α (low- α) sequence, made up of old, kinematically hot (young, kinematically cold) stars. More recently, this simple picture of geometric distinction between α -sequences has been disfavored (Schönrich & Binney 2009; Bovy et al. 2012a,b,c; Bland-Hawthorn et al. 2019). However, it is nevertheless clear that there are two distinct α -sequences, now thought to be the result of chemical enrichment (Hayden et al. 2017).

The estimation of α -abundance from *Gaia* XP spectra has recently been explored in the literature. In particular, several stellar label dependent analyses reach different conclusions:

Initially, Gavel et al. (2021) used the **ExtraTrees** algorithm and argued that XP spectra α -abundance measurements for stars which *only* vary in $[\alpha/M]$ are dangerous because XP α -abundance estimates strongly depend on correlations between $[\alpha/M]$ and other stellar properties, particularly $[M/H]$. This is potentially because the *Gaia* XP spectra are so extremely low-resolution that a change in $[\alpha/M]$ is equivalent to a change in $[M/H]$ due to multiple spectral lines being blurred together. Then, Witten et al. (2022) supported these findings by demonstrating that synthetic XP spectra do not contain significant α -information. It is important to mention that both of these works emphasize that these difficulties are especially relevant for warm stars with $5000 \text{ K} > T_{\text{eff}}$. Additionally, the work of Witten et al. (2022) involved synthetic XP spectra with $G = 16$.

More recently, Li et al. (2023) developed **AspGap**: a regression model which they argue is capable of producing

precise $[\alpha/M]$ estimates for red giant branch stars with $3000 \leq T_{\text{eff}} \leq 7000 \text{ K}$, which extends into the warm regime where α -abundance estimates have been identified as problematic. Then, Hattori (2024) trained a tree-based machine learning model to estimate both $[M/H]$ and $[\alpha/M]$. Additionally, Hattori (2024) proposed specific lines which contribute to α -information: Na D lines (589 nm) and the Mg I line (516 nm).

We emphasize that the conclusions drawn by these previous analyses about α -information in *Gaia* XP are obfuscated by their reliance on stellar labels. Indirect stellar label correlations with $[\alpha/M]$ cannot be separated from a stellar label dependent model. On the other hand, our stellar label independent model can potentially assess α -information without the issue of stellar label correlations.

To contribute to answering the question of α -information in *Gaia* XP, we separate high- and low- α sequence members of the pristine labels catalog⁷ in $[M/H]$ - $[\alpha/M]$ space (left panel of Figure 13) and project them into our latent space (middle panel of Figure 13). To do so, we implement the α -bimodality split of Patil et al. (2023), who used copulas and elicitable maps to cleanly separate the high- (purple) and low- α (orange) sequences⁸. We then produce binned tracks along both sequences in $[M/H]$ - $[\alpha/M]$ space and project them into our latent space, similarly to the tracks in Section 4. We observe that our z_1 and z_3 latent variables function as

⁷ The pristine labels catalog is strictly composed of giants with $T_{\text{eff}} < 5000 \text{ K}$. As such, we are probing the cool star regime where α -abundance estimates are less problematic.

⁸ Patil et al. (2023) define their split in $[\text{Fe}/\text{H}]$ - $[\text{Mg}/\text{Fe}]$ space. As such, we do the same, and subsequently present the split in $[M/H]$ - $[\alpha/M]$ space.

pseudo-[M/H]-[α /M] space, in which the high- and low- α sequence tracks are clearly separated. However, these tracks alone are not definitive evidence of α -information, because they are functions of metallicity! Therefore, we produce a third track at fixed metallicity ([M/H] = -0.4, red), which only varies with [α /M]. Crucially, the fixed metallicity track transitioning between both sequences in [M/H]-[α /M] space is reproduced in our latent space. The high- and low- α sequence tracks, with the addition of the fixed metallicity track, is compelling evidence that the *Gaia* XP spectra contain meaningful α -information (for low-temperature stars).

To reinforce our argument that the latent space has learned important chemical information, we also present a sample of *Gaia*-Enceladus stars (e.g. Helmi et al. 2018) in both [M/H]-[α /M] space (left panel of Figure 13) and latent space (right panel of Figure 13). To select *Gaia*-Enceladus members, we implement the action diamond cuts recommended by Lane et al. (2022): $|L_z/J_{\text{tot}}| < 0.07$ and $(J_z - J_R)/J_{\text{tot}} < -0.3$, for stars in the full catalog. The evident separation of the *Gaia*-Enceladus population from both the high- and low- α sequences in the latent space (particularly across z_1) demonstrates that our model is sensitive to the anomalous abundances characteristic of *Gaia*-Enceladus.

To be even more pessimistic, one can also worry about T_{eff} and $\log g$ correlations with [α /M] which could masquerade as [α /M] trends in our latent space. Furthermore, it is also possible that the underlying extinction distributions for the high- α and low- α populations differ significantly. Therefore, it is important to demonstrate that our model actually learns the α -bimodality without relying on stellar label and extinction relationships beyond metallicity. Lastly, some work has already been undertaken to present evidence for α -information content in a similar manner to the above. For example, synthetic photometry from *Gaia* XP spectra has been used to demonstrate that [α /M] appears to correlate with the *Gaia* colour C1M395-C1M410 (see Fig. 28 of (Gaia Collaboration et al. 2022)).

As such, we train four Random Forest Classifiers (`sklearn.ensemble.RandomForestClassifier`; Pedregosa et al. 2011) on 90% of the pristine catalog to classify the high- and low- α sequences, with the intent of quantitatively demonstrating α -information in a novel way. The classifiers are as follows:

- (i) *Metal classifier*: trained on only [M/H].
- (ii) *Label classifier*: trained on all stellar labels except [α /M]; T_{eff} , $\log g$ and [M/H].
- (iii) *Label + extinction classifier*: trained on all stellar labels except [α /M], as well as 2MASS colours:

$J - H$ and $J - K$, to provide training labels which can proxy extinction.

- (iv) *Latent classifier*: trained on latent variables from our stellar label independent model.

We present the confusion matrices for the four classifiers in Figure 14 over test data (10% of the pristine labels catalog). To compare the classifiers, we use accuracy:

$$\text{Acc} = \frac{\text{T high-}\alpha + \text{T low-}\alpha}{\text{T high-}\alpha + \text{T low-}\alpha + \text{F high-}\alpha + \text{F low-}\alpha}, \quad (13)$$

where T is true and F is false. Furthermore, we define the high- α predictive value:

$$P_{\text{high-}\alpha} = \frac{\text{T high-}\alpha}{\text{T high-}\alpha + \text{F high-}\alpha}, \quad (14)$$

and the low- α predictive value:

$$P_{\text{low-}\alpha} = \frac{\text{T low-}\alpha}{\text{T low-}\alpha + \text{F low-}\alpha}. \quad (15)$$

We present Acc, $P_{\text{high-}\alpha}$ and $P_{\text{low-}\alpha}$ for each classifier in Table 1. Classifying the α -bimodality with our stellar label independent latent space yields a 16% (14%) improvement in accuracy relative to the [M/H] (labels) classifier. The latent classifier achieves 18% (14%) improvement in high- α (low- α) predictive value relative to the [M/H] classifier. Interestingly, the labels classifier achieves comparable low- α predictive value to the latent classifier. However, this is at the expense of a far worse high- α predictive value: the latent classifier achieves a 43% improvement relative to the labels classifier. This is a somewhat counter-intuitive result: including T_{eff} and $\log g$ produces worse high- α classification than [M/H] alone. This is additional evidence that our latent space is learning genuine α -information, since the inclusion of stellar labels beyond [M/H] decreases high- α classification performance as opposed to improving high- α classification performance. Finally, the inclusion of 2MASS $J - H$ and $J - K$ colours does not yield any significant classification performance, which indicates that our latent classifier is not relying on different extinction trends between α -sequences.

In summary, the striking α -bimodality tracks in our latent space, in combination with the significant improvement in α -bimodality classification with our latent classifier relative to the [M/H], label and label + extinction classifiers, is stellar label independent evidence that the *Gaia* XP spectra can and should be used to estimate [α /M] with $T_{\text{eff}} < 5000$ K. We speculate that it may be possible to estimate [α /M] above this temperature threshold if a training sample with hotter stars than APOGEE were used.

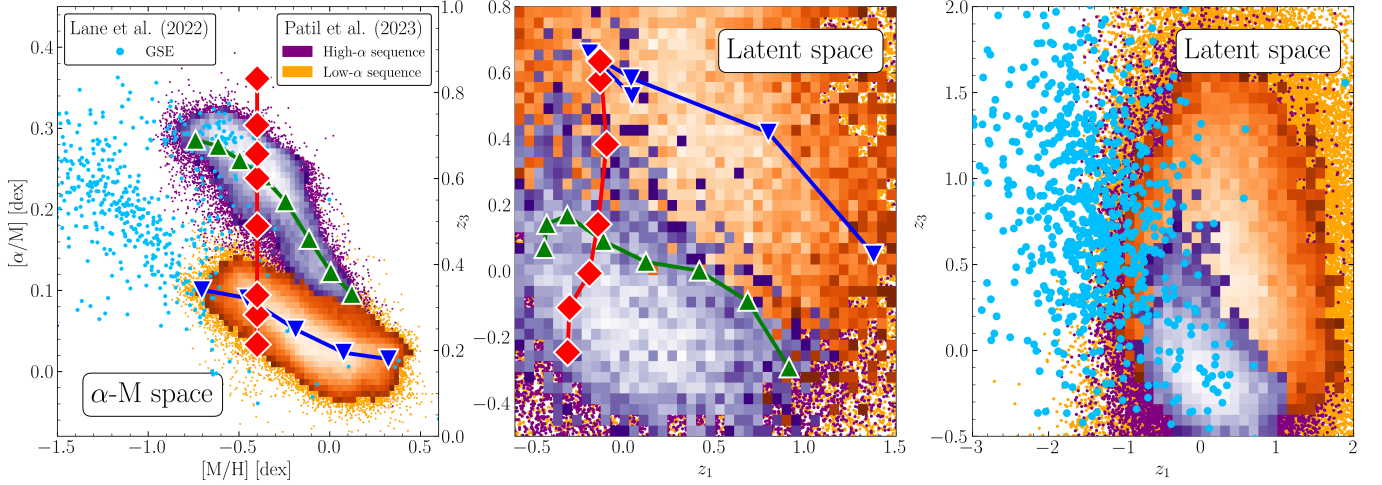


Figure 13. High- α sequence (green), low- α sequence (blue) and fixed-metallicity (red) tracks in both $[\alpha/M]$ - $[M/H]$ space (left) and latent space (center). High/low α -sequence members (purple, orange) are classified according to the α -bimodality split of Patil et al. (2023). α -information is demonstrated by the high- and low- α sequence tracks, in combination with the fixed metallicity track (which only varies in $[\alpha/M]$). Additionally, we present a sample of *Gaia*-Enceladus stars selected via action diamond cut from Lane et al. (2022) (light blue) in both $[\alpha/M]$ - $[M/H]$ space (left) and latent space (right).

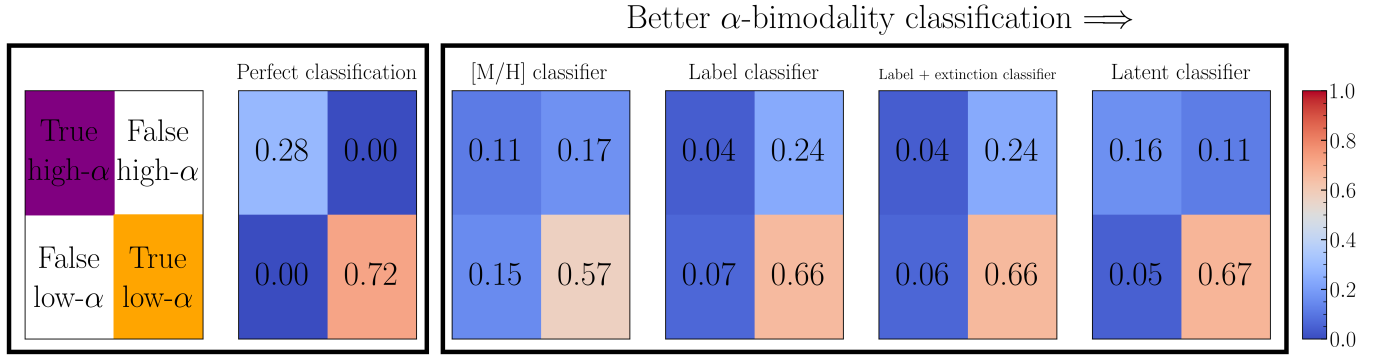


Figure 14. Confusion matrices for the label, $[M/H]$ and latent classifiers presented in Section 5. Confusion matrix classifications and perfect classification value are presented in the two left-most matrices. The latent classifier, trained on our *sVAE* latent variables, achieves better α -bimodality classification than classifiers relying on stellar labels.

Table 1. Performance metrics for the three α -bimodality classifiers presented in Section 5. A perfect classifier would yield accuracy $\text{Acc} = 1$, high- α predictive value $P_{\text{high-}\alpha} = 1$ and low- α predictive value $P_{\text{low-}\alpha} = 1$.

| | [M/H] | Label | Label + extinction | Latent |
|--------------------------|-------|-------|--------------------|--------|
| Acc | 0.68 | 0.70 | 0.70 | 0.84 |
| $P_{\text{high-}\alpha}$ | 0.40 | 0.15 | 0.15 | 0.58 |
| $P_{\text{low-}\alpha}$ | 0.79 | 0.91 | 0.92 | 0.93 |

6. SUMMARY, DISCUSSION AND OUTLOOK

We have presented a stellar label independent model for the recently released *Gaia* XP spectra from *Gaia* DR3. The primary goal of this work was to demonstrate that stellar label independent models can close the stellar labels gap, and in doing so rectify system-

atics which indirectly stem from theoretically estimated stellar labels. Our main results are as follows:

- (i) *Model performance*: Our stellar label independent model achieves competitive, and in specific regimes better, performance relative to existing stellar label dependent models. Specifically, our *scatter* VAE better reconstructs XP spectra over the redder XP wavelengths. Also, our model can be applied to cool stars and low-surface gravity stars: stellar label regions where stellar label dependent models struggle.
- (ii) *Interpretability*: We investigated stellar label trends in our stellar label independent latent space to demonstrate that the latent space is rich with astrophysical information, such as: main-sequence and giant branch evolution, and metallicity.

- (iii) *$[\alpha/M]$ information:* We provided strong evidence that the *Gaia* XP spectra contain meaningful α -abundance information without relying on stellar label correlations, for stars with $T_{\text{eff}} < 5000$, while also being sensitive to the anomalous abundances of *Gaia*-Enceladus stars. As such, we encourage the astronomy community to make use of the *Gaia* XP spectra for estimation of $[\alpha/M]$ in this temperature regime⁹.

Despite the successes of our stellar label independent approach, our model has some limitations. Specifically:

- (i) *Indirect astrophysical information:* By definition, our stellar label independent model does not produce stellar label estimates. As such, extracting astrophysical information from our model is not necessarily straightforward. Nevertheless, our model can provide important astrophysical understanding by investigating the latent space clustering of XP spectra.
- (ii) *Mileage may vary:* The variational nature of a VAE, implemented with the re-parametrization trick, means that each projection of an XP spectrum into the latent space, or simulation of an XP spectrum from the latent space, is not unique. Results and figures in this work may vary slightly if recomputed as a result. However, these variations should not be severe, since the latent variable sampling is restricted to ϵ (see Section 3).
- (iii) **Gaia* XP spectra coverage:* The first implementation of our stellar label independent model was trained on a small minority of XP spectra relative to the total number of XP spectra available ($< 1\%$). As such, we caution the application of our model to stellar populations beyond the XP/APOGEE cross-match, such as white dwarfs. This does not mean that our model cannot eventually accommodate these populations. Future implementations of our stellar label independent model can be trained on any/all XP spectra because our model does not require stellar labels to train on.

The potential of our stellar label independent approach was certainly not fully explored in this initial work. We identify several promising areas for future work:

- (i) *Denoising:* Data compression into the latent space can remove noise from input XP spectra, because

our model is designed to learn the most important features shared by the training data and neglect observational noise. Imposing restrictions on our scatter estimator could produce de-noised XP spectra with model error estimates smaller than observational uncertainties.

- (ii) *Global XP model:* In principle, our stellar label independent model could be trained to accurately simulate the entire XP dataset. This would require a sophisticated training procedure due to the gigantic amount of training data required for a single model to learn all of the stellar populations in XP spectra ($\sim 10^8$ training objects).
- (iii) *Outlier detection:* A particularly promising application of our stellar label independent model is the detection of rare sub-populations in the *Gaia* XP spectra via the latent space. Recently, [Lucey et al. \(2022\)](#) released a catalog of carbon-enhanced metal poor stars candidates (CEMPs). A preliminary application of our model to their CEMP catalog has uncovered a population of outliers. We suspect that a fraction of these outliers could be binary systems containing CEMP stars, since a major formation channel for CEMPs is binary mass transfer (e.g. [Goswami et al. 2021](#)). We will soon present a detailed analysis of this outlier population in *Laroche et al. 2025 (in prep.)*.

In conclusion, the *Gaia* DR3 XP spectra is by far the largest spectroscopic survey to date. Stellar label dependent models are incapable of exploiting the entirety of the astrophysical information in the XP dataset, because they are limited by the availability of stellar labels to train on. Novel data-driven techniques must be developed to tackle this big data problem. Furthermore, as the sizes of future large-scale surveys grow, the stellar labels gap for stellar label dependent models will only widen as the relative stellar label coverage decreases (e.g. SPHEREx: [Unwin et al. 2016](#); Roman: [Zellem et al. 2022](#); Vera Rubin: [Ivezić et al. 2019](#)). The big-data era in astronomy is both a blessing, for discovery, and a curse, for data analysis. Our stellar label independent model is an important step towards *fully* data-driven modeling which can confront the intimidating amount of data which present and future large-scale surveys provide.

ACKNOWLEDGMENTS

We are thankful to the anonymous referee for providing a constructive report that helped improve the clarity and quality of our manuscript. AL and JSS would like to thank Maria Drout, Jo Bovy, Henry Leung, Gregory M.

⁹ The extension of $[\alpha/M]$ estimates to hotter stars in the *Gaia* XP spectra was not addressed in this work.

Green, Jiadong Li and Jeff Shen for valuable feedback. AL also acknowledges helpful conversations with James Lane, Nolan Koblishke and Amy Prickett which greatly improved this work. AL acknowledges support from the Natural Sciences and Engineering Research Council of Canada (NSERC) and is (was) partially funded through a NSERC Canada Graduate Scholarship - Doctoral (Master's). AL is also supported by the Data Sciences Institute at the University of Toronto through grant number DSI-DSFY3R1P02. JSS acknowledges support from the Natural Sciences and Engineering Research Council of Canada (NSERC) funding reference #RGPIN-2023-04849. The Dunlap Institute is funded through an endowment established by the David Dunlap family and the University of Toronto.

This work has made use of data from the European Space Agency (ESA) mission *Gaia* (<https://www.cosmos.esa.int/gaia>), processed by the *Gaia* Data Processing and Analysis Consortium (DPAC, <https://www.cosmos.esa.int/web/gaia/dpac/consortium>). Funding for the DPAC has been provided by national institutions, in particular the institutions participating in the *Gaia* Multilateral Agreement. Additionally, this work made use of the Python package *GaiaXPy*, developed and maintained by members of the *Gaia* DPAC, and in particular, Coordination Unit 5 (CU5), and the Data Processing Centre located at the Institute of Astronomy, Cambridge, UK (DPCI).

APPENDIX

A. COMPLETE LATENT SPACE REPRESENTATION

By inspecting the latent tracks across all six dimensions in Figure 15, it is clear that the Kiel information our model has learned is not exclusively encoded into z_4 and z_5 (see Figure 2). Rather, the stellar label information is shared amongst all latent dimensions, with varying degrees of sensitivity to stellar labels. This additional level of complexity, relative to stellar label dependent models which do not blend stellar label information across summary statistics, can lead to better performance (see Section 4.1). Simultaneously, the blending of astrophysical information across the latent space makes the assessment of our model behaviour less straightforward. This also means that our latent dimensions can learn information beyond stellar labels. As such, our stellar label independent framework is very much optimized for outlier detection.

REFERENCES

- Abdurro'uf, e. 2022, *ApJS*, 259, 35, doi: [10.3847/1538-4365/ac4414](https://doi.org/10.3847/1538-4365/ac4414)
- Anders, F., Chiappini, C., Santiago, B. X., et al. 2014, *A&A*, 564, A115, doi: [10.1051/0004-6361/201323038](https://doi.org/10.1051/0004-6361/201323038)
- Bensby, T., Feltzing, S., & Lundström, I. 2003, *A&A*, 410, 527, doi: [10.1051/0004-6361:20031213](https://doi.org/10.1051/0004-6361:20031213)
- Bergemann, M. 2014, in *Determination of Atmospheric Parameters of B*, ed. E. Niemczura, B. Smalley, & W. Pych, 187–205, doi: [10.1007/978-3-319-06956-2_17](https://doi.org/10.1007/978-3-319-06956-2_17)
- Bland-Hawthorn, J., Sharma, S., Tepper-Garcia, T., et al. 2019, *MNRAS*, 486, 1167, doi: [10.1093/mnras/stz217](https://doi.org/10.1093/mnras/stz217)
- Bovy, J., Rix, H.-W., & Hogg, D. W. 2012a, *ApJ*, 751, 131, doi: [10.1088/0004-637X/751/2/131](https://doi.org/10.1088/0004-637X/751/2/131)
- Bovy, J., Rix, H.-W., Hogg, D. W., et al. 2012b, *ApJ*, 755, 115, doi: [10.1088/0004-637X/755/2/115](https://doi.org/10.1088/0004-637X/755/2/115)
- Bovy, J., Rix, H.-W., Liu, C., et al. 2012c, *ApJ*, 753, 148, doi: [10.1088/0004-637X/753/2/148](https://doi.org/10.1088/0004-637X/753/2/148)
- Buder, S., Sharma, S., Kos, J., et al. 2021, *MNRAS*, 506, 150, doi: [10.1093/mnras/stab1242](https://doi.org/10.1093/mnras/stab1242)
- De Angeli, F., Weiler, M., Montegriffo, P., et al. 2022, arXiv e-prints, arXiv:2206.06143, doi: [10.48550/arXiv.2206.06143](https://doi.org/10.48550/arXiv.2206.06143)
- Fuhrmann, K. 1998, *A&A*, 338, 161
- Gaia Collaboration, Montegriffo, P., Bellazzini, M., et al. 2022, arXiv e-prints, arXiv:2206.06215, doi: [10.48550/arXiv.2206.06215](https://doi.org/10.48550/arXiv.2206.06215)
- Gaia Collaboration, e. 2022, arXiv e-prints, arXiv:2208.00211, doi: [10.48550/arXiv.2208.00211](https://doi.org/10.48550/arXiv.2208.00211)
- Gavel, A., Andrae, R., Fouesneau, M., Korn, A. J., & Sordo, R. 2021, *A&A*, 656, A93, doi: [10.1051/0004-6361/202141589](https://doi.org/10.1051/0004-6361/202141589)
- Gilmore, G., & Reid, N. 1983, *MNRAS*, 202, 1025, doi: [10.1093/mnras/202.4.1025](https://doi.org/10.1093/mnras/202.4.1025)
- Goswami, P. P., Rathour, R. S., & Goswami, A. 2021, *A&A*, 649, A49, doi: [10.1051/0004-6361/202038258](https://doi.org/10.1051/0004-6361/202038258)
- Guo, Y., Shang, G., Vazirgiannis, M., & Clavel, C. 2023, arXiv e-prints, arXiv:2311.09807, doi: [10.48550/arXiv.2311.09807](https://doi.org/10.48550/arXiv.2311.09807)

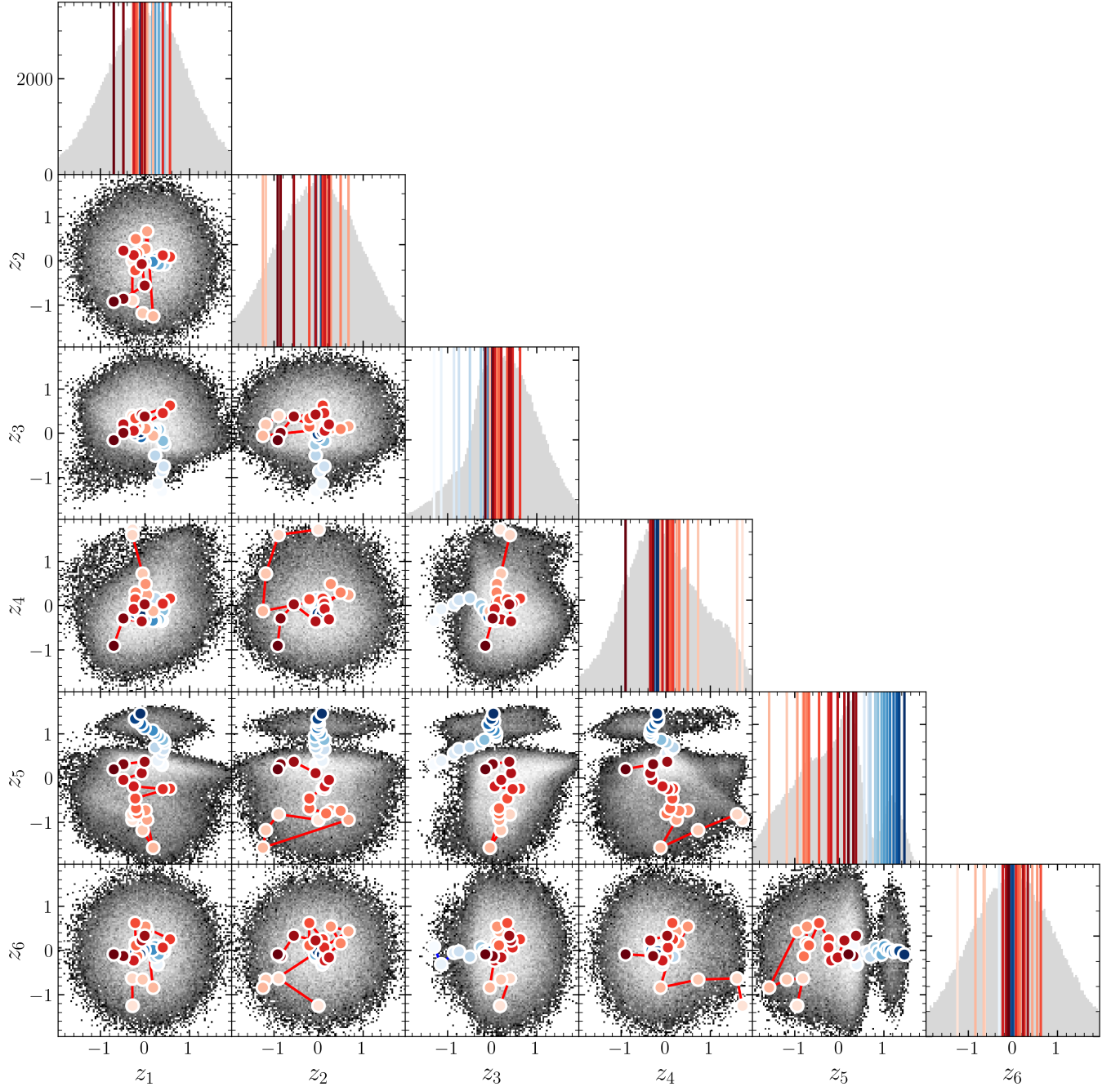


Figure 15. Main-sequence (MS, blue track) and giant branch (GB, red track) evolutionary tracks as a function all latent variables (the same tracks presented in Figure 2). 1D and 2D marginal distributions are presented over the entire 6D latent space. In the 1D marginal distributions, evolutionary track points are presented as vertical lines. Stellar label information is shared across all latent dimensions.

- Hattori, K. 2024, arXiv e-prints, arXiv:2404.01269, doi: [10.48550/arXiv.2404.01269](https://doi.org/10.48550/arXiv.2404.01269)
- Hayden, M. R., Recio-Blanco, A., de Laverny, P., Mikolaitis, S., & Worley, C. C. 2017, *A&A*, 608, L1, doi: [10.1051/0004-6361/201731494](https://doi.org/10.1051/0004-6361/201731494)
- Hayden, M. R., Bovy, J., Holtzman, J. A., et al. 2015, *ApJ*, 808, 132, doi: [10.1088/0004-637X/808/2/132](https://doi.org/10.1088/0004-637X/808/2/132)
- Heiter, U., Jofré, P., Gustafsson, B., et al. 2015, *A&A*, 582, A49, doi: [10.1051/0004-6361/201526319](https://doi.org/10.1051/0004-6361/201526319)
- Helmi, A., Babusiaux, C., Koppelman, H. H., et al. 2018, *Nature*, 563, 85, doi: [10.1038/s41586-018-0625-x](https://doi.org/10.1038/s41586-018-0625-x)
- Ivezić, Ž., Kahn, S. M., Tyson, J. A., et al. 2019, *ApJ*, 873, 111, doi: [10.3847/1538-4357/ab042c](https://doi.org/10.3847/1538-4357/ab042c)
- Kingma, D. P., & Ba, J. 2014, arXiv e-prints, arXiv:1412.6980, doi: [10.48550/arXiv.1412.6980](https://doi.org/10.48550/arXiv.1412.6980)
- Kingma, D. P., & Welling, M. 2013, arXiv e-prints, arXiv:1312.6114, doi: [10.48550/arXiv.1312.6114](https://doi.org/10.48550/arXiv.1312.6114)
- Kordopatis, G., Wyse, R. F. G., Gilmore, G., et al. 2015, *A&A*, 582, A122, doi: [10.1051/0004-6361/201526258](https://doi.org/10.1051/0004-6361/201526258)
- Kullback, S., & Leibler, R. A. 1951, *The Annals of Mathematical Statistics*, 22, 79.
<http://www.jstor.org/stable/2236703>
- Lane, J. M. M., Bovy, J., & Mackereth, J. T. 2022, *MNRAS*, 510, 5119, doi: [10.1093/mnras/stab3755](https://doi.org/10.1093/mnras/stab3755)
- Laroche, A. 2024, AlexLaroche7/xp_vae: First release, v1.0.0, Zenodo, doi: [10.5281/zenodo.14041979](https://doi.org/10.5281/zenodo.14041979)
- Leung, H. W., & Bovy, J. 2019, *MNRAS*, 483, 3255, doi: [10.1093/mnras/sty3217](https://doi.org/10.1093/mnras/sty3217)
- . 2023, *MNRAS*, doi: [10.1093/mnras/stad3015](https://doi.org/10.1093/mnras/stad3015)
- Li, J., Wong, K. W. K., Hogg, D. W., Rix, H.-W., & Chandra, V. 2023, arXiv e-prints, arXiv:2309.14294, doi: [10.48550/arXiv.2309.14294](https://doi.org/10.48550/arXiv.2309.14294)
- Loshchilov, I., & Hutter, F. 2016, arXiv e-prints, arXiv:1608.03983, doi: [10.48550/arXiv.1608.03983](https://doi.org/10.48550/arXiv.1608.03983)
- . 2017, arXiv e-prints, arXiv:1711.05101, doi: [10.48550/arXiv.1711.05101](https://doi.org/10.48550/arXiv.1711.05101)
- Lucey, M., Al Kharusi, N., Hawkins, K., et al. 2022, arXiv e-prints, arXiv:2206.08299, doi: [10.48550/arXiv.2206.08299](https://doi.org/10.48550/arXiv.2206.08299)
- Majewski, S. R., Schiavon, R. P., Frinchaboy, P. M., et al. 2017, *AJ*, 154, 94, doi: [10.3847/1538-3881/aa784d](https://doi.org/10.3847/1538-3881/aa784d)
- Montegriffo, P., De Angeli, F., Andrae, R., et al. 2022, arXiv e-prints, arXiv:2206.06205, doi: [10.48550/arXiv.2206.06205](https://doi.org/10.48550/arXiv.2206.06205)
- Ness, M., Hogg, D. W., Rix, H. W., Ho, A. Y. Q., & Zasowski, G. 2015, *ApJ*, 808, 16, doi: [10.1088/0004-637X/808/1/16](https://doi.org/10.1088/0004-637X/808/1/16)
- Nidever, D. L., Bovy, J., Bird, J. C., et al. 2014, *ApJ*, 796, 38, doi: [10.1088/0004-637X/796/1/38](https://doi.org/10.1088/0004-637X/796/1/38)
- O’Brian, T., Ting, Y.-S., Fabbro, S., et al. 2021, *ApJ*, 906, 130, doi: [10.3847/1538-4357/abca96](https://doi.org/10.3847/1538-4357/abca96)
- Patil, A. A., Bovy, J., Jaimungal, S., Frankel, N., & Leung, H. W. 2023, *MNRAS*, 526, 1997, doi: [10.1093/mnras/stad2820](https://doi.org/10.1093/mnras/stad2820)
- Pedregosa, F., Varoquaux, G., Gramfort, A., et al. 2011, *Journal of Machine Learning Research*, 12, 2825, doi: [10.48550/arXiv.1201.0490](https://doi.org/10.48550/arXiv.1201.0490)
- Ruz-Mieres, D., & Kostrzewa-Rutkowska, Z. 2023, *gaia-dpci/GaiaXPy: GaiaXPy, Version 2.1.0*, Zenodo, doi: [10.5281/zenodo.8239995](https://doi.org/10.5281/zenodo.8239995)
- Schönrich, R., & Binney, J. 2009, *MNRAS*, 399, 1145, doi: [10.1111/j.1365-2966.2009.15365.x](https://doi.org/10.1111/j.1365-2966.2009.15365.x)
- Skrutskie, M. F., Cutri, R. M., Stiening, R., et al. 2006, *AJ*, 131, 1163, doi: [10.1086/498708](https://doi.org/10.1086/498708)
- Ting, Y.-S., Conroy, C., Rix, H.-W., & Cargile, P. 2019, *ApJ*, 879, 69, doi: [10.3847/1538-4357/ab2331](https://doi.org/10.3847/1538-4357/ab2331)
- Touvron, H., Lavril, T., Izacard, G., et al. 2023, *LLaMA: Open and Efficient Foundation Language Models*.
<https://arxiv.org/abs/2302.13971>
- Unwin, S. C., SPHEREx Science Team, & SPHEREx Project Team. 2016, in *American Astronomical Society Meeting Abstracts*, Vol. 228, American Astronomical Society Meeting Abstracts #228, 216.10
- Witten, C. E. C., Aguado, D. S., Sanders, J. L., et al. 2022, *MNRAS*, 516, 3254, doi: [10.1093/mnras/stac2273](https://doi.org/10.1093/mnras/stac2273)
- Yan, H., Li, H., Wang, S., et al. 2022, *The Innovation*, 3, 100224, doi: [10.1016/j.xinn.2022.100224](https://doi.org/10.1016/j.xinn.2022.100224)
- Yoshii, Y. 1982, *PASJ*, 34, 365
- Zellem, R. T., Nemati, B., Gonzalez, G., et al. 2022, in *Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series*, Vol. 12180, *Space Telescopes and Instrumentation 2022: Optical, Infrared, and Millimeter Wave*, ed. L. E. Coyle, S. Matsuura, & M. D. Perrin, 121801Z, doi: [10.1117/12.2627567](https://doi.org/10.1117/12.2627567)
- Zhang, X., Green, G. M., & Rix, H.-W. 2023, arXiv e-prints, arXiv:2303.03420, doi: [10.48550/arXiv.2303.03420](https://doi.org/10.48550/arXiv.2303.03420)