

PART III ESSAY 2024-25
FACULTY OF MATHEMATICS, DAMTP

Applying Contrastive Learning to Stellar Spectra

May 8, 2025

Abstract

We present STARCLIP

Key words: methods: data analysis – techniques: spectroscopic – stars: fundamental parameters

Contents

1	Introduction	1
1.1	Related Works	3
1.2	Outline	3
1.3	Notations	3
2	Machine Learning Methodology	4
2.1	Supervised Learning	4
2.1.1	1D-CNNs: Feedforward Stage	4
2.1.2	1D-CNNs: Backpropagation Stage	5
2.2	Self-supervised Learning	6
2.2.1	Contrastive Learning	6
2.2.2	Role of Temperature in Contrastive Learning	7
3	Data	7
4	StarCLIP Model Implementation	7
4.1	Spectral Encoder	7
4.2	StarCLIP Model	8
4.2.1	8
5	Results	8
5.1	StarCLIP Model	8
5.1.1	Example Retrieval by Similarity Search	8
5.1.2	Stellar Property Estimation	9
5.1.3	Manifold Learning	9
6	Conclusion	9
7	Acknowledgments	9
8	Data Availability	10

1 Introduction

Absorption features exhibited on a star’s spectrum encode its physical and chemical properties, which in turn provides a ‘fossil record’ of the formation history of the host galaxy in which that star is located (Sandford et al., 2020). The James Webb Space Telescope (JWST), deployed in 2021, exhibits unique spectroscopic capabilities, enabling efficient, high-quality resolved star spectroscopy in neighbouring galaxies which have thus far been too distant, faint or crowded for previous instruments to detect (Sandford et al., 2020). These capabilities have motivated the launch of several recent astronomical programmes, including a number of completed surveys observing ~ 100 red giant branch (RGB) stars in M31 (Nidever et al., 2023) and ~ 300 RGB stars in three globular clusters (Scalco et al., 2024), an ongoing survey – at the time of this paper – observing ~ 300 RGB stars across three isolated dwarf galaxies (Space Telescope Science Institute (STScI), 2023) and a planned programme to observe ~ 200 more M31 disc stars (STScI, 2024).

However, due to its novelty and uniqueness, the domain of JWST resolved-star spectroscopy remains largely unexplored, and many unanswered questions remain regarding optimal strategies for collecting, reducing and analysing its valuable data. In recent years, a growing subset of astronomical surveys have employed data-driven methodologies from machine learning (ML) to extract insights from optical data (Parker et al., 2024). To date, these algorithms are generally divided into two groups:

- **Supervised** algorithms use input-output pairs, provided by a human expert, to learn the mapping from a set of features to target variables. Once obtained, the mapping can be used to perform predictive *classification* or *regression* tasks on unseen data. These methods have achieved tremendous success in data-rich settings, proving pivotal in image recognition tasks and natural language processing (NLP) (Fabbro et al., 2018). However, requiring extensive labelled training data inhibits their application to data-poorer domains, where labels are scarce, imbalanced or expensive to obtain. Additionally, supervised networks are generally *task-specific*, with each new task typically requiring a bespoke model trained or designed from scratch (Parker et al., 2024). This may lead to high computational costs, limited generalisability, and higher risk of overfitting.
- **Unsupervised** algorithms take in as inputs only the set of measured features, without accompanying labels. Their output is a non-linear, non-invertible transformation of the input features, consisting of (i) *clusters*, which are groups of objects sharing similar features, (ii) lists of detected *anomalies*, which are unusual objects in the sample, or (iii) low-dimensional representations of the objects (Baron, 2019). These algorithms can extract new knowledge from datasets that we did not know existed and, therefore, could not have directly searched for (Reis et al., 2021). To this extent, they may be arguably more useful for ‘AI for science’ research, where human-annotated labels are not always available in quality or quantity. However, in practice, unsupervised networks typically underperform compared to their supervised counterparts, especially in computer vision (CV), so they are seldom implemented in real-world applications (Zhuang et al., 2019).

Large-scale surveys such as Apache Point Observatory Galactic Evolution Experiment (APOGEE), Galactic Archaeology with HERMES (GALAH) and Large Sky Area Multi-object Fiber Spectroscopic Telescope (LAMOST) provide large ($n \sim 10^5$ to 10^7), homogeneous databases of labelled spectra that are virtually ideal for supervised ML applications (Fabbro et al., 2018). By contrast, the JWST survey provides a considerably smaller ($n \sim 100$) database, which limits its eligibility for supervised analysis. What is more, large differences in field-of-view and sensitivity with ground-based surveys make it difficult to acquire overlapping stars with previous datasets (Gardner et al., 2006). This poses challenges in exploring ways to use informative results from large surveys to facilitate analysis of JWST data, motivating the use of synthetic spectra to force overlap. What results is a database of complementary or *cross-modal* observations of the same optical spectra.

One promising ML approach to analyse cross-modal realisations is **contrastive learning**, a form of **self-supervised learning** (SSL). In this approach, a model simultaneously embeds cross-modal data into a shared latent space, creating pairs of high-quality embeddings, *i.e.*, low-dimensional representations of objects that preserve much of their physical information (Parker et al., 2024). With minimal processing, these embeddings can serve as a robust ‘foundation’ for downstream tasks,

particularly, zero- and one-shot learning¹. Hence, these models are often referred to as **foundation models** (Bommasani et al., 2022). Notably, in computationally constrained settings where training bespoke supervised models from scratch is impractical, foundation models consistently demonstrate superior performance on zero- and one-shot tasks compared to their fully supervised counterparts (Bommasani et al., 2022).

In this work, we present STARCLIP, and its variant, MOCKSTARCLIP, cross-modal foundation models designed for stellar spectroscopy. Our approach consists of two main steps. First, we train one-dimensional **convolutional neural networks** (1D-CNNs) on single-modal, labelled spectral datasets. Next, we deploy the CNNs as pre-trained encoders to embed spectra in a unified latent space. We apply contrastive learning to align embeddings around shared physical semantics: embeddings corresponding to the same underlying star are brought closer together, while those corresponding to different stars are pushed farther apart. The resulting embeddings can then be seamlessly used for a variety of downstream tasks, including semantic searches and zero- and one-shot predictive tasks, which ultimately facilitates analysis of JWST data.

Spectroscopic data consist of 19,001 pairs of APOGEE² and synthetic JWST spectra generated from *atlas12* and *synthe* codes maintained by R. Kurucz (Kurucz, 1970, 1993, 2005, 2013, 2017; Kurucz and Avrett, 1981). Stellar labels associated to each pair consist of atmospheric parameters, namely, effective temperature (T_{eff}) and surface gravity ($\log g$), and 18 elemental abundances ($[X/H]$). To emulate the real JWST data, we apply two distinct sets of stochastic transformations, resulting in ‘mock’ and ‘extended-mock’ versions. A small sample ($n = 59$) of available real data on missing gaps of JWST signals, collected from JWST Cycle 3³, is used in this data transformation pipeline. STARCLIP and MOCKSTARCLIP result from applying our methodology to APOGEE–mock JWST pairs and downsampled mock–extended-mock JWST pairs, respectively.

A schema of the pipeline is depicted in Figure 1.

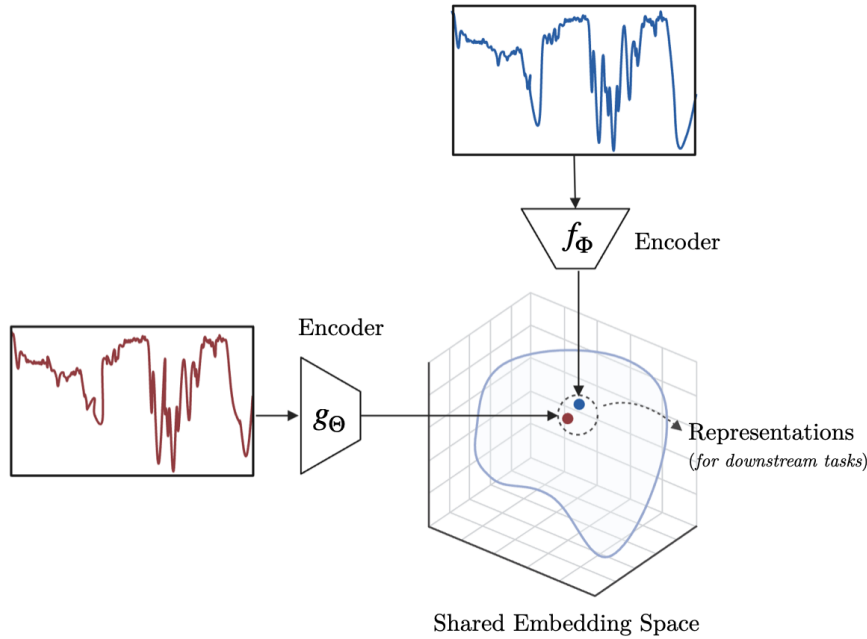


Figure 1: Illustration of the architecture of the STARCLIP and MOCKSTARCLIP models. Two encoders, f_Φ and g_Θ , separately embed pairs of cross-modal spectra into a shared embedding space, where they are aligned under contrastive loss. Learned embeddings then contain high-level physical information that allow them to be seamlessly transferred for downstream tasks. Modalities are distinguished by colour; here, red and blue.

¹In zero- and one-shot learning, a model leverages its learned representations to identify a new class of interest, having previously seen zero or one exemplar of that class during training, respectively.

²<https://www.sdss4.org/dr17/>

³<https://www.stsci.edu/jwst/science-execution/approved-programs/general-observers/cycle-3-go>

We summarise the contributions of this paper, as follows:

- We develop one of the first self-supervised foundation ML models for analysing JWST data, with a pipeline put into place to ensure seamless integration with *real* JWST data.
- We apply supervised learning to train robust CNNs to predict stellar parameters with high precision and accuracy.
- We apply cross-modal contrastive training to align pre-trained encoders around shared physical properties, creating a discriminative latent space.
- With minimal processing, we enable zero- and one-shot transfer of the models to downstream tasks, including (i) in-modal and cross-modal cosine similarity searches and (ii) stellar property estimation from spectra. These tasks empirically show that latent embeddings capture key physical properties of underlying stars.
- We apply non-linear manifold learning on the latent space, visualising its intrinsic local geometry in a lower-dimensional setting.

1.1 Related Works

This paper builds on two recent works. [Fabbro et al. \(2018\)](#) present a deep CNN architecture, called STARNET, capable of determining stellar parameters – T_{eff} , $\log g$, and $[Fe/H]$ – with high accuracy and precision from real APOGEE and synthetic APOGEE ASSET spectra. Elsewhere, [Parker et al. \(2024\)](#) present a cross-modal foundation model, called ASTROCLIP, which aligns representations of galaxy images and spectra in a shared, unified latent space. Pre-trained transformer-based encoders are used to separately embed the modalities into latent space, where they are aligned under contrastive loss. The embeddings are then applied to a variety of downstream tasks, including semantic similarity searches, photometric redshift estimation, galaxy property estimation, and morphology classification, in which they demonstrate remarkable performance, even relative to supervised baselines. Both articles provide accompanying code in publicly accessible GITHUB repositories at the following links: [STARNET](#) ([Fabbro et al., 2018](#)) and [ASTROCLIP](#) ([Parker et al., 2024](#)), respectively.

We build our own foundation models – STARCLIP and its variant MOCKSTARCLIP – in a two-step process involving carefully modified versions of STARNET and ASTROCLIP.

1.2 Outline

Code for our models is publicly available on a GITHUB repository [here](#). Our paper is organised, as follows. In §2, we provide theoretical background on supervised and self-supervised methods relevant to this paper, addressing their merits and disadvantages. In §3, we provide the astronomical datasets used for analysis. In §4, we discuss in detail the implementation and training process of the STARCLIP and MOCKSTARCLIP models. We present our results on in-modal and cross-modal similarity searches, stellar property prediction, and non-linear manifold learning in §5. Finally, §6 contains some concluding remarks and further extensions to our work.

1.3 Notations

In this paper, we use $[n]$, where $n \in \mathbb{N}$, to denote the set $\{1, \dots, n\}$. We use $|\cdot|$ to denote the dimension of vectors and $\|\cdot\|$ to represent their Euclidean ℓ_2 -norm. We write $a \gg b$ whenever there exists a sufficiently small constant c such that $b/a < c$ holds. We denote by $a \otimes b$ the convolution operation between vectors a and b .

2 Machine Learning Methodology

In this section, we discuss the supervised and self-supervised methodologies underlying our models.

2.1 Supervised Learning

Let $X \in \mathcal{X} \subseteq \mathbb{R}^n$ and $Y \in \mathcal{Y} \subseteq \mathbb{R}^m$ be random variables, where $Y = f(X)$, for some unknown f . In supervised learning (SL), the learning algorithm is given a set of training examples $\{(x_i, y_i)\}_{i=1}^T$, drawn from the joint distribution of X and Y . Its objective is to learn a mapping $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$, which minimises the expected loss, as measured by a specific *loss function* $L : \mathcal{Y}^2 \rightarrow \mathbb{R}$. That is,

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}[L(Y, f(X))], \quad (2.1)$$

over some space of hypotheses functions \mathcal{F} (Koushik, 2016). In regression tasks, choosing the mean squared error (MSE) loss gives: $\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}[\|Y - f(X)\|^2]$.

2.1.1 1D-CNNs: Feedforward Stage

1D-CNNs, introduced by Boser et al. (1990), are a class of deep learning algorithms that solve equation (2.1) by adaptively learning spatial hierarchies of features in the input data (Yamashita et al., 2018). As shown in Figure .., 1D-CNNs work by passing X through a series of computational *nodes*, arranged in five types of layers: (i) input, (ii) convolutional, (iii) pooling, (iv) fully connected and (v) output. To illustrate the sequence of operations performed, we consider a single training example $x \in X$ fed to the network. In what follows, denote by $\zeta^{(l)}$ the output vector of the l th hidden layer and (\mathbf{W}, \mathbf{b}) the collection of learnable weights/filters and biases associated to the network.

- **Input layer.** The input layer is simply the identity map: $\zeta_{\text{in}} := x$.
- **Convolutional layers.** Following Cacciari and Ranfagni (2024), consider a collection of p filters $\{W_j^{(1)}\}_{j=1}^p \in \mathbf{W}$ of size $s^{(1)}$. Each filter slides across the single axis of the input vector $\zeta_{\text{in}} := x$ of length n , aggregating local information from sections of $s^{(1)}$ consecutive elements. Assuming each filter iterates one element at a time (*i.e.*, stride = 1), the discrete *convolution* between its filter weights and the input data can be computed as:

$$C_j^{(1)}(i) = (x \otimes W_j^{(1)})(i) = \sum_{u=1}^{s^{(1)}} x(i+u-1)W_j^{(1)}(u), \quad (2.2)$$

where $C_j^{(1)}, j \in [p]$, has length $n - s^{(1)} + 1$. As shown in Figure ..., we then have a 2D *multi-filter convolution* matrix $C^{(1)}$ with size $(n - s^{(1)} + 1) \times p$, whose columns are given by (2.2). Adding bias $b_j^{(1)} \in \mathbf{b}$ and applying an activation map $\phi : \mathbb{R} \rightarrow \mathbb{R}$ to the ij th element of $C^{(1)}$ then gives the ij th convolved feature at the output of the first convolutional layer:

$$\zeta_j^{(1)}(i) = \phi \left(C_j^{(1)}(i) + b_j^{(1)} \right). \quad (2.3)$$

In general, the feature map $\zeta^{(1)}$ captures localised, low-level features from the input data. To perform an additional convolution, consider a third-order tensor $W^{(2)} \in \mathbf{W}$ of dimension $s^{(2)} \times p \times q$, where $s^{(2)}$ is the filter size and $p \times q$ denotes the number of filters. We define a 2D multifilter convolution matrix $C^{(2)}$ with size $(n - s^{(2)} + 1) \times q$, whose columns $C_j^{(2)}$ are given by the sum of p convolutions:

$$C_j^{(2)}(i) = \sum_{l=1}^p \left(\zeta_l^{(1)} \otimes W_{l,j}^{(2)} \right) (i). \quad (2.4)$$

As before, adding bias $b_j^{(2)} \in \mathbf{b}$ and applying an activation ϕ to the ij th element of $C^{(2)}$ gives the ij th convolved feature at the output of the second convolutional layer:

$$\zeta_j^{(2)}(i) = \phi \left(C_j^{(2)}(i) + b_j^{(2)} \right). \quad (2.5)$$

Obtaining the second feature map $\zeta^{(2)}$ by performing a convolution across the first $\zeta^{(1)}$ enables the model to extract higher-order feature representations.

- **Pooling layer.** Using a 1D pooling operation (here, max-pooling), the pooling layer down-samples the feature map by a factor of w , thereby compressing its size and extracting only the strongest features learnt in the convolutional layers. To perform the max-pooling operation, a window of size w moves in strides of the same length along each column of the feature map, extracting the maximum value from each considered sub-region. In particular, its operation can be represented by:

$$\zeta_j^{(3)}(i) = \max\{\zeta_j^{(2)}(i_s)\}_{k=1,2,\dots,w}, \quad (2.6)$$

where $i_s = wi - k + 1$ (Cacciari and Ranfagni, 2024). This produces an output vector $\zeta^{(3)}$ of size $w \times q$.

- **Fully connected layers.** A flattening operation is performed, re-shaping the 2D output $\zeta_j^{(3)}$ into a 1D vector f , as follows:

$$f(j + n(i - 1)) = \zeta_j^{(3)}(i). \quad (2.7)$$

This produces an output vector f of length wq , which can be fed into the fully-connected network. Assuming the first and second fully-connected layers admit weights and biases $(W^{(4)}, b^{(4)})$ and $(W^{(5)}, b^{(5)})$ in (\mathbf{W}, \mathbf{b}) , then

$$\zeta^{(4)}(i) = \phi \left(\sum_{j=1}^{wq} W^{(4)}(i, j) f(j) + b^{(4)}(i) \right) \quad (2.8)$$

$$\zeta^{(5)}(i) = \phi \left(\sum_{j=1}^{|\zeta^{(4)}|} W^{(5)}(i, j) \zeta^{(4)}(j) + b^{(5)}(i) \right). \quad (2.9)$$

- **Output layer.** The output vector is simply given by $\zeta_{\text{out}} = \zeta^{(5)}$, corresponding to the network's prediction of the true label (*target*) associated to x .

2.1.2 1D-CNNs: Backpropagation Stage

Altogether, the combination of these layers transform the input x into a non-linear output $\hat{f}(x; \mathbf{W}, \mathbf{b})$. For each batch of T training examples $\{(x_i, y_i)\}_{i=1}^T$, \mathbf{W} and \mathbf{b} can be estimated by minimising the *empirical loss* between predicts and targets:

$$\hat{\mathbf{W}}, \hat{\mathbf{b}} = \operatorname{argmin}_{\mathbf{W}, \mathbf{b}} \frac{1}{T} \sum_{i=1}^T \|y_i - \hat{f}(x_i; \mathbf{W}, \mathbf{b})\|^2, \quad (2.10)$$

where we neglect regularisation for simplicity (Fabbro et al., 2018). During the training phase, a suitable choice of gradient-based optimiser is deployed to minimise empirical losses of the form (2.10). Using the *backpropagation algorithm*, the algorithm propagates the loss function backwards through the layers of the network, from the output to the input layer, successively computing its gradients with respect to each parameter in the network using the chain rule. As these gradients indicate how much each parameter must be changed to reduce the loss, the algorithm can make informed and iterative adjustments to those weights until reaching a minimum (Cacciari and Ranfagni, 2024). Of note, the *Adam* optimiser is particularly suitable to optimising loss functions, where the underlying datasets are large and/or the parameter spaces are high-dimensional (Kingma and Ba, 2017). With few memory requirements, Adam can efficiently minimise the objective function, using adaptive estimates of low-order gradient moments to adjust learning rates (Kingma and Ba, 2017).

Although 1D-CNNs have achieved state-of-the-art performance, particularly in CV tasks, their computational cost poses a challenge. Having a high number of learnable parameters, 1D-CNNs must be trained on a large volume of labelled training data, which limits their applicability to problems where data sizes are small (Keshari et al., 2018) or where labels are limited, costly or even unavailable.

2.2 Self-supervised Learning

An alternative to SL, SSL has drawn attention for its efficient training of large-scale models without the need for explicit labels. Intuitively, SSL leverages the input data itself as supervision: it uses encoders to find complex structures and co-occurrence relationships in the data in order to generate high-quality, discriminative representations of the data. To illustrate, given the incomplete sentence ‘I like _____ football’, a well-trained NLP model might predict ‘playing’ as the appropriate blank in this Cloze Test, because ‘playing’ and ‘football’ frequently co-occur in text corpora (Liu et al., 2023). These representations can then be applied to almost all types of downstream tasks, where their performance often exceeds even supervised learning (Chuang et al., 2020).

As summarised in Liu et al. (2023), SSL methods can be broadly classified into three categories:

- **Generative:** trains an encoder-decoder pair to map input x to an embedding z and then reconstruct x from z .
- **Contrastive:** trains an encoder to map input pairs (x_1, x_2) to embeddings (z_1, z_2) to measure the similarity between them.
- **Adversarial:** trains an encoder-decoder pair to generate false samples and a discriminator to distinguish them from real samples.

2.2.1 Contrastive Learning

Contrastive learning (CL), the focus of this paper, aims to ‘learn to compare’ representations across different *modalities*, that is, types of data input. Its central idea is to contrast semantically similar (positive) and dissimilar (negative) pairs of possibly *augmented* views of data points, encouraging representations of similar pairs to be close together and dissimilar pairs to be far apart (Chuang et al., 2020). In this context, an *augmented view* refers to a synthetically modified version of an original data point in order to generate positive pairs for training (Lee et al., 2024). Augmentations introduce semantically meaningful variations to the data, forcing the model to learn invariant features that help it to generalise to downstream tasks (Liu et al., 2024). To illustrate, the Contrastive Language-Image Pretraining (CLIP) method, due to Radford et al. (2021), trains encoders to contrast language based-descriptions with their corresponding image pair, where the image is augmented with random square crops and the text with bi-grams with high pointwise *mutual information* (MI) for the pair.

To formalise CL, let $\mathcal{X} \subseteq \mathbb{R}^n$ and $\mathcal{Y} \subseteq \mathbb{R}^m$ be the input space of two different modalities. Consider a batch of N unlabelled samples $\{(x_1, y_1), \dots, (x_N, y_N)\} \sim \mathcal{D}_{\mathcal{X} \times \mathcal{Y}}$, where (x_i, y_i) is semantically related, e.g., a JWST spectra and its APOGEE pair. CLIP’s objective is to learn a pair of encoders $f_\Phi : \mathcal{X} \rightarrow \mathbb{R}^d$ and $g_\Theta : \mathcal{Y} \rightarrow \mathbb{R}^d$, $d \ll \min(n, m)$, that compress the two modalities from raw space into a shared latent space, such that $z_i := f_\Phi(x_i)$ and $z'_j := g_\Theta(y_j)$ are close to each other if they share ‘similar’ semantics, otherwise far away.

The Information Noise Contrastive Estimation (InfoNCE) loss used to train the encoders is specified by (van den Oord et al., 2019):

$$\mathcal{L}_{\text{InfoNCE}}(f_\Phi, g_\Theta, \tau) = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(s_{ii}/\tau)}{\sum_{j \neq i}^N \exp(s_{ij}/\tau)}. \quad (2.11)$$

Here, $\tau \in \mathbb{R}_{>0}$ denotes a trainable temperature parameter, discussed in §2.2.2, and $s_{ij} := s(z_i, z'_j)$ is a similarity measure between the features z_i and z'_j . In particular, the *cosine* similarity measure can be used:

$$s_{ij} := s(z_i, z'_j) = \frac{z_i^\top z'_j}{\|z_i\|^2 \|z'_j\|^2}. \quad (2.12)$$

By minimising (2.11), the similarity between semantically related features (z_i, z'_i) is expected to be large, whilst that between semantically unrelated features $(z_i, z'_j), j \neq i$, is expected to be small.

Of note, InfoNCE has an insightful connection to MI. A formal proof given by van den Oord et al. (2019) shows that:

$$I(z_i; z'_j) \geq \log(k) - \mathcal{L}_{\text{InfoNCE}}, \quad (2.13)$$

where k is the number of negative samples to a given sample x_i (for which $z_i = f_{\Theta}(x_i)$). This shows that minimising $\mathcal{L}_{\text{InfoNCE}}$ maximises the variational lower bound on the MI $I(z_i; z'_j)$, which is itself bounded above by $I(x_i; y_j)$ through the data-processing inequality. Moreover, the log-dependency on k implies that having more negative samples leads to improved representations, a fact formally proved by Tian et al. (2020). Altogether, InfoNCE provides a good lower bound approximation to MI, enjoying greater stability and lower variance than most other competing approaches (Song and Ermon, 2020).

2.2.2 Role of Temperature in Contrastive Learning

The temperature hyperparameter controls the strength of penalties on hard negative samples. In particular, contrastive loss with small t

As derived in Manna et al. (2023), the gradient of \mathcal{L}_i with respect to s_{ii} and s_{ij} is given by

$$\frac{\partial \mathcal{L}_i}{\partial s_{ii+}} = -\frac{1}{\tau} \sum_{k \neq i} p_{ik}, \quad \frac{\partial \mathcal{L}_i}{\partial s_{ij}} = \frac{1}{\tau} p_{ij} \quad (2.14)$$

so that

$$r(s_{ij}) = \frac{\left| \frac{\partial \mathcal{L}_i}{\partial s_{ij}} \right|}{\left| \frac{\partial \mathcal{L}_i}{\partial s_{ii}} \right|} = \frac{\exp\left(\frac{s_{ij}}{\tau}\right)}{\sum_{k \neq i} \exp\left(\frac{s_{ik}}{\tau}\right)}. \quad (2.15)$$

This shows

3 Data

Following stellar properties:

- T_{eff} : Effective Temperature.
- $\log(g)$: Surface Gravity.
- $[X/H]$: Metallicities.

4 StarCLIP Model Implementation

4.1 Spectral Encoder

Our spectral encoder is closely modelled after STARNET Fabbro et al. (2018). As illustrated in Figure ..., the encoder's is a deep CNN consisting of the following 7 layers:

- **Input layer:** Projects the input features to ..
-
- **Output layer:** Produces a 20-dimensional vector of stellar parameter estimations.

The first takes in the input data, which it passes to two 1-D convolutional layers with 4 and 16 filters, respectively, to extract local features. Each filter ... feature identifier, ..., producing a set of feature maps.

filters in the first layer detect low-level features, whilst those in the second layer convolve over ... to learn higher-order features Fabbro et al. (2018).

with ... and ... filters, then a max pooling layer with ..., followed by three successive fully connected layers with The final layer is the output layer.

There are but a few differences:

-

$$h(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + b).$$

Whilst sigmoid, \tanh and rectified linear unit (ReLU) are the most common activation functions, *our* version of STARNET exchanges ReLU for the exponential linear unit (ELU):

$$g(z) = \begin{cases} z & \text{if } z > 0, \\ \alpha(e^z - 1) & \text{if } z \leq 0. \end{cases} \quad (4.1)$$

4.2 StarCLIP Model

The STARCLIP model receives paired observations of JWST-APOGEE data, processing the input through the appropriate spectral encoder. This results in a processed sequence of 20-dimensional embedding vectors of both types of spectra that reside in a shared, latent space. The alignment of similar embedding vectors is performed during CLIP training, detailed below.

4.2.1

In zero-shot learning, the model applies its learned representations to identify or categorise unseen data instances, without the need for additional training on those instances. In few-shot learning, the pre-trained model is

5 Results

5.1 StarCLIP Model

The (M)STARCLIP models can be easily deployed to downstream tasks, with zero-shot transfer or minimal task adaptation. To this end, we embed paired spectra in the *held-out test database*, as follows:

$$\begin{aligned} (\mathbf{x}_i^J, \mathbf{y}_j^A) &\xrightarrow{\text{StarCLIP}} (\mathbf{z}_i^J, \mathbf{z}_j^A) \in \mathbb{R}^{20} \\ (\mathbf{x}_i^J, \mathbf{y}_j^{MJ}) &\xrightarrow{\text{MStarCLIP}} (\mathbf{z}_i^J, \mathbf{z}_j^{MJ}) \in \mathbb{R}^{20}. \end{aligned} \quad (5.1)$$

Normalised embeddings are obtained on division by their l_2 -norm.

5.1.1 Example Retrieval by Similarity Search

We perform example retrieval by cosine similarity search. This sort of search has proven useful in retrieving interesting objects and performing clustering, as exemplified in [Popat et al. \(2017\)](#).

Given some queried spectra, we assign a cosine similarity score (??) between its normalised embedding $\hat{\mathbf{z}}_Q$ and all other embeddings in the held-out test set. From the pool of scores, we retrieve the spectra with the highest similarity scores to its query, completing the task; no additional transformation or processing steps are required. We call a search constrained to same-modality embeddings (e.g., $S(\hat{\mathbf{z}}_Q^J, \hat{\mathbf{z}}^J)$) an **in-modality similarity search**, and a search constrained to cross-modality embeddings (e.g., $S(\hat{\mathbf{z}}_Q^J, \hat{\mathbf{z}}^A)$) a **cross-modality similarity search**. By design, an in-modality search requires the closest neighbour to the query to be the query itself.

For all four possible pairs of modalities, we present the spectra of the ‘closest’ match retrieved by STARCLIP and MSTARCLIP in Figure ... We report the averages of the highest cosine similarity score obtained during searches in Table ... Overall, the closeness of these mean scores to 1 indicate

STARCLIP and MOCKSTARCLIP are efficient in retrieving semantically similar spectra within and across modalities.

5.1.2 Stellar Property Estimation

We evaluate STARCLIP and MSTARCLIP’s performance on stellar property estimation using JWST and mock JWST embeddings as inputs, respectively. In general, these tasks involve developing an entire pipeline and training bespoke models end-to-end from scratch [Parker et al. \(2024\)](#). However, because learned representations already capture key physical information about the input stars, we can easily use zero- and one-shot algorithms to predict stellar properties [Parker et al. \(2024\)](#). Specifically, we use zero-shot k -Nearest Neighbours (k -NN) and one-shot linear and Extreme Gradient Boosting (XGBoost) to regress ‘true’ stellar parameters against corresponding model predictions.

... some training details ...

We present our results in Figure We report the R^2 performance of all regressors in Table ...

Consistent with Parker *et al.* [Parker et al. \(2024\)](#), strong zero- and one-shot performance of the models indicate that, under CLIP alignment, JWST and mock JWST embeddings are organised in latent space around stellar properties. This is consistent with our intuition that spectra should be naturally informative about the values of their stellar parameters.

5.1.3 Manifold Learning

The results in §5.1.1 and §5.1.2 indicate ASTROCLIP and MOCKASTROCLIP are able to capture in its embeddings core physical properties of the input spectra. To further explore the structure of the embedding space, we apply a state-of-the-art **manifold learning** method: t -distributed Stochastic Neighborhood Embedding (t -SNE). Intuitively, t -SNE works by compressing high-dimensional data onto a lower dimensional manifold in such a way as to preserve much of the local geometric structure and reveal global structure, including the presence of clusters.

Low-dimensional representations are produced by minimising the K-L divergence between ...

Moreover, t -SNE has $O(n^2)$ complexity, making it computationally slow for the large datasets required in the ASTROCLIP model.

to capture neighbourhood relationships that exist in the higher-dimensional latent space and

6 Conclusion

We present STARCLIP, a cross-modal foundation ML model for stellar spectroscopy. Our model is a compositional model, with CNN encoders

STARNET [Fabbro et al. \(2018\)](#)

7 Acknowledgments

We thank Drs Nathan Sandford, Jo Bovy, Joshua Speagle and Ting Li of the University of Toronto and Dr Miles Cranmer of the University of Cambridge for their insightful discussions. We thank Dr Nathan Sandford for providing the APOGEE and JWST databases, without which this work would not have been possible. Computations in this paper were run on the **Geir** server, a resource provided by the Dunlap Institute of Astronomy and Astrophysics at the University of Toronto. The analysis is supported by a fully open-source software stack, including ASTROPY ([Astropy Collaboration, 2022](#)), H5PY ([Collette, 2013](#)), MATPLOTLIB ([Hunter, 2007](#)), NUMPY ([van der Walt et al., 2011](#); [Harris et al., 2020](#)), PANDAS ([McKinney, 2010](#); [Reback et al., 2022](#)), RANDOM ([Van Rossum, 2020](#)), SCIPY ([Virtanen et al., 2020](#)), SKLEARN ([Pedregosa et al., 2011](#)) and TORCH ([Paszke et al., 2019](#)). Illustrations were created using bioRender.⁴

⁴<https://www.biorender.com/>

8 Data Availability

The data and code underlying this article are available in a GITHUB repository at this [link](#).

References

- Astropy Collaboration. Astropy collaboration, 2022. *The Astrophysical Journal*, 935:167, 2022.
- Dalya Baron. Machine learning in astronomy: a practical overview, 2019. URL <https://arxiv.org/abs/1904.07248>.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Nieves, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher R. Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the opportunities and risks of foundation models, 2022. URL <https://arxiv.org/abs/2108.07258>.
- Bernhard Boser, Yann LeCun, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Handwritten digit recognition with a back-propagation network. In *Advances in Neural Information Processing Systems*. Citeseer, 1990.
- I. Cacciari and A. Ranfagni. Hands-on fundamentals of 1d convolutional neural networks a tutorial for beginner users. *Applied Sciences*, 14(18):8500, 2024. doi: 10.3390/app14188500. URL <https://doi.org/10.3390/app14188500>.
- Ching-Yao Chuang, Joshua Robinson, Lin Yen-Chen, Antonio Torralba, and Stefanie Jegelka. De-biased contrastive learning, 2020. URL <https://arxiv.org/abs/2007.00224>.
- Andrew Collette. *Python and HDF5*. O'Reilly, 2013.
- S. Fabbro, K. A. Venn, T. O'Brian, S. Bialek, C. L. Kielty, F. Jahandar, and S. Monty. An application of deep learning in the analysis of stellar spectra. *Monthly Notices of the Royal Astronomical Society*, 475(3):2978–2993, April 2018. doi: 10.1093/mnras/stx3298. URL <https://doi.org/10.1093/mnras/stx3298>.
- J. P. Gardner, J. C. Mather, and M. et al. Clampin. The james webb space telescope. *Space Science Reviews*, 123(4):485–606, April 2006. doi: 10.1007/s11214-006-8315-7. URL <https://doi.org/10.1007/s11214-006-8315-7>.
- Charles R. Harris, K. Jarrod Millman, Stefan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernandez del Ro, Mark Wiebe, Pearu Peterson, Pierre Grard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with numpy. *Nature*, 585(7825):357362, September 2020. ISSN 1476-4687. doi: 10.1038/s41586-020-2649-2. URL <http://dx.doi.org/10.1038/s41586-020-2649-2>.
- J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3): 90–95, 2007. doi: 10.1109/MCSE.2007.55.

- Rohit Keshari, Mayank Vatsa, Richa Singh, and Afzel Noore. Learning Structure and Strength of CNN Filters for Small Sample Size Training. *arXiv e-prints*, art. arXiv:1803.11405, March 2018. doi: 10.48550/arXiv.1803.11405.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>.
- Jayanth Koushik. Understanding convolutional neural networks, 2016. URL <https://arxiv.org/abs/1605.09081>.
- R. L. Kurucz. Saors, 309. *SAOSR*, 309, 1970.
- R. L. Kurucz. *SYNTHES Spectrum Synthesis Programs and Line Data*. Smithsonian Astrophysical Observatory, Cambridge, MA, 1993.
- R. L. Kurucz. Msais, 8, 14. *MSAIS*, 8:14, 2005.
- R. L. Kurucz. Atlas12: Opacity sampling model atmosphere program. Astrophysics Source Code Library, 2013. ascl:1303.024.
- R. L. Kurucz. Atlas9: Model atmosphere program with opacity distribution functions. Astrophysics Source Code Library, 2017. ascl:1710.017.
- R. L. Kurucz and E. H. Avrett. Saors, 391. *SAOSR*, 391, 1981.
- J. Lee, I. Ham, Y. Kim, and H. Ko. Time-series representation feature refinement with a learnable masking augmentation framework in contrastive learning. *Sensors*, 24(24):7932, 2024. doi: 10.3390/s24247932. URL <https://doi.org/10.3390/s24247932>.
- Xiao Liu, Fanjin Zhang, Zhenyu Hou, Li Mian, Zhaoyu Wang, Jing Zhang, and Jie Tang. Self-supervised learning: Generative or contrastive. *IEEE Transactions on Knowledge and Data Engineering*, 35(1):857–876, 2023. doi: 10.1109/TKDE.2021.3090866.
- Z. Liu, M. Shu, and W. Zhu. Contrastive learning framework for bitcoin crash prediction. *Stats*, 7(2):402–433, 2024. doi: 10.3390/stats7020025. URL <https://doi.org/10.3390/stats7020025>.
- Siladitya Manna, Soumitri Chattopadhyay, Rakesh Dey, Saumik Bhattacharya, and Umapada Pal. Dynamically scaled temperature in self-supervised contrastive learning. *arXiv preprint arXiv:2308.01140*, 2023. URL <https://doi.org/10.48550/arXiv.2308.01140>. Last revised 10 May 2024 (v2).
- W. McKinney. Data structures for statistical computing in python. In S. van der Walt and J. Millman, editors, *Proceedings of the 9th Python in Science Conference*, pages 56–61, 2010.
- David L. Nidever, Karoline Gilbert, Erik Tollerud, Charles Siders, Ivanna Escala, Carlos Allende Prieto, Verne Smith, Katia Cunha, Victor P. Debattista, Yuan-Sen Ting, and Evan N. Kirby. The prevalence of the α -bimodality: First jwst α -abundance results in m31, 2023. URL <https://arxiv.org/abs/2306.04688>.
- Liam Parker, Francois Lanusse, Siavash Golkar, Leopoldo Sarra, Miles Cranmer, Alberto Bietti, Michael Eickenberg, Geraud Krawezik, Michael McCabe, Rudy Morel, Ruben Ohana, Mariel Pettee, Bruno Rgaldo-SaintBlancard, Kyunghyun Cho, and Shirley Ho. Astroclip: a cross-modal foundation model for galaxies. *Monthly Notices of the Royal Astronomical Society*, 531(4):49905011, June 2024. ISSN 1365-2966. doi: 10.1093/mnras/stae1450. URL <http://dx.doi.org/10.1093/mnras/stae1450>.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.

- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830, 2011.
- Shraddha K. Popat, Pramod B. Deshmukh, and Vishakha A. Metre. Hierarchical document clustering based on cosine similarity measure. In *2017 1st International Conference on Intelligent Systems and Information Management (ICISIM)*, pages 153–159, 2017. doi: 10.1109/ICISIM.2017.8122166.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- Jeff Reback, Jbrockmendel, Wes McKinney, Joris Van Den Bossche, Matthew Roeschke, Tom Augspurger, Simon Hawkins, Phillip Cloud, Gfyoung, Sinhrks, Patrick Hoeffler, Adam Klein, Terji Petersen, Jeff Tratner, Chang She, William Ayd, Shahar Naveh, JHM Darbyshire, Richard Shadrach, Marc Garcia, Jeremy Schendel, Andy Hayden, Daniel Saxton, Marco Edward Gorelli, Fangchen Li, Torsten Wörtwein, Matthew Zeitlin, Vytautas Jancauskas, Ali McMaster, and Thomas Li. pandas-dev/pandas: Pandas 1.4.3, June 2022.
- I. Reis, M. Rotman, D. Poznanski, J.X. Prochaska, and L. Wolf. Effectively using unsupervised machine learning in next generation astronomical surveys. *Astronomy and Computing*, 34:100437, 2021. ISSN 2213-1337. doi: 10.1016/j.ascom.2020.100437. URL <https://www.sciencedirect.com/science/article/pii/S2213133720300913>.
- Nathan R. Sandford, Daniel R. Weisz, and Yuan-Sen Ting. Forecasting chemical abundance precision for extragalactic stellar archaeology. *The Astrophysical Journal Supplement Series*, 249(2):24, July 2020. ISSN 1538-4365. doi: 10.3847/1538-4365/ab9cb0. URL <http://dx.doi.org/10.3847/1538-4365/ab9cb0>.
- M. Scalco, M. Libralato, R. Gerasimov, L. R. Bedin, E. Vesperini, D. Nardiello, A. Bellini, M. Griggio, D. Apai, M. Salaris, A. Burgasser, and J. Anderson. Jwst imaging of the closest globular clusters - iii. multiple populations along the low-mass main-sequence stars of ngc 6397. *Astronomy Astrophysics*, 689:A59, 2024. doi: 10.1051/0004-6361/202450589.
- Jiaming Song and Stefano Ermon. Understanding the limitations of variational mutual information estimators, 2020. URL <https://arxiv.org/abs/1910.06222>.
- Space Telescope Science Institute (STScI). Jwst cycle 2 general observer abstracts, 2023. URL https://www.stsci.edu/files/live/sites/www/files/home/jwst/science-execution/approved-programs/general-observers/cycle-2-go/_documents/jwst-cycle2-general-observer-abstracts.pdf. Accessed: 2024-12-15.
- STScI. Jwst cycle 3 general observer abstracts, 2024. URL https://www.stsci.edu/files/live/sites/www/files/home/jwst/science-execution/approved-programs/general-observers/cycle-3-go/_documents/jwst-cycle3-general-observer-abstracts.pdf. Accessed: 2024-12-15.
- Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding, 2020. URL <https://arxiv.org/abs/1906.05849>.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2019. URL <https://arxiv.org/abs/1807.03748>.
- S. van der Walt, S. C. Colbert, and G. Varoquaux. The numpy array: A structure for efficient numerical computation. *Computing in Science & Engineering*, 13(2):22–30, 2011. doi: 10.1109/MCSE.2011.37.
- Guido Van Rossum. *The Python Library Reference, release 3.8.2*. Python Software Foundation, 2020.
- P. Virtanen et al. Scipy 1.0: Fundamental algorithms for scientific computing in python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.

- R. Yamashita, M. Nishio, R.K.G. Do, and K. Togashi. Convolutional neural networks: an overview and application in radiology. *Insights into Imaging*, 9:611–629, Aug 2018. doi: 10.1007/s13244-018-0639-9. URL <https://doi.org/10.1007/s13244-018-0639-9>.
- Chengxu Zhuang, Alex Lin Zhai, and Daniel Yamins. Local aggregation for unsupervised learning of visual embeddings, 2019. URL <https://arxiv.org/abs/1903.12355>.