



ASTROMER

A transformer-based embedding for the representation of light curves[★]

C. Donoso-Oliva^{1,3,4}, I. Becker^{2,4,5} , P. Protopapas³, G. Cabrera-Vives^{1,4} , M. Vishnu³, and H. Vardhan⁵

¹ Department of Computer Science, Universidad de Concepcion, Concepcion 4070386, Chile
e-mail: cridonoso@inf.udec.cl

² Department of Computer Science, Pontificia Universidad Catolica de Chile, Macul, Santiago 7820436, Chile

³ Inst. for Applied Computational Science, Harvard University, Cambridge, MA 02138, USA

⁴ Millennium Institute of Astrophysics (MAS), Nuncio Monsenor Sotero Sanz 100, Providencia, Santiago, Chile

⁵ Univ. AI, Singapore, 050531, Singapore

Received 2 May 2022 / Accepted 31 October 2022

ABSTRACT

Taking inspiration from natural language embeddings, we present ASTROMER, a transformer-based model to create representations of light curves. ASTROMER was pre-trained in a self-supervised manner, requiring no human-labeled data. We used millions of R-band light sequences to adjust the ASTROMER weights. The learned representation can be easily adapted to other surveys by re-training ASTROMER on new sources. The power of ASTROMER consists in using the representation to extract light curve embeddings that can enhance the training of other models, such as classifiers or regressors. As an example, we used ASTROMER embeddings to train two neural-based classifiers that use labeled variable stars from MACHO, OGLE-III, and ATLAS. In all experiments, ASTROMER-based classifiers outperformed a baseline recurrent neural network trained on light curves directly when limited labeled data were available. Furthermore, using ASTROMER embeddings decreases the computational resources needed while achieving state-of-the-art results. Finally, we provide a Python library that includes all the functionalities employed in this work.

Key words. methods: statistical – stars: statistics – techniques: photometric

1. Introduction

Over the past few decades, efforts have been made to develop machine learning tools to analyze and discover variable phenomena in the sky. These tools will face a challenge with the construction of the new generation of telescopes, which will generate substantially more data than the old generation (Kremer et al. 2017). Moreover, the observations will be deeper and more precise than ever.

With the upcoming telescopes such as the *Vera C. Rubin* Observatory (Ivezić et al. 2019), periodic observations of a significant portion of the sky every few days will become the norm. They will produce light curves for all the objects in the observed fields.

Traditional machine learning methods rely on features to explore the variable behavior of light curves. They are based on quantities such as amplitude, period, and color information, among others (Sánchez-Sáez et al. 2021). The expected volume of data will present a significant challenge to feature-based methods if applied to every measured object.

Deep learning techniques have an advantage over traditional machine learning as they do not need to pre-calculate features, as they extract informative representations of the data automatically (LeCun et al. 2015). Also, deep learning methods leverage GPU parallelization, which accelerates information extraction from the light curves.

In recent years, some of these methods have been developed with promising results. Naul et al. (2018) presented a

recurrent autoencoder to extract features from folded periodical light curves. After training, the authors fitted a random forest classifier on the embedded space, outperforming models trained on predefined features. This model was one of the first approaches in astronomy that used unlabeled light curves to extract characteristics from an unsupervised deep learning model. A similar idea was proposed a year later, by Tsang & Schultz (2019), who used an autoencoder not only for classification, but also for novelty detection. These approaches remain supervised, which is problematic for training on small datasets (Charnock & Moss 2017).

Self-supervised models are ideal for overcoming the limitations of labeled datasets, as they can generate the target values automatically without human-generated labels (Liu et al. 2021). Most of these methods learn representations by solving auxiliary tasks that do not require any label, such as infilling of time series. The learned representations can be used downstream to solve other tasks, such as the classification of variable stars or the regression of physical parameters, such as temperature or redshift. This learning scheme has achieved impressive results in natural language processing (NLP), with Bidirectional Encoder Representations from Transformers (BERT, Devlin et al. 2018) being one of the most notable examples. BERT-like models learn to extract contextual representations using two auxiliary tasks from a dataset of raw text. In this stage, known as pre-training, the model learns a general representation of the data, without being specific to any particular task. This process is the most resource-intensive part that usually takes weeks to complete, as the models, as well as the datasets, are large. After this stage is complete, the model only needs a few hours of further training to

[★] The library, main code, and pre-trained weights are available at <https://github.com/astromer-science>

adjust the weights to a more specific domain, in a stage known as fine-tuning. Once finalized, BERT is used as the initial stage of a new model. This scheme enables BERT-based models to achieve a state-of-the-art performance.

The main components of BERT are the self-attention layers (Vaswani et al. 2017), which codify similarities between words in the input sentence. These layers can be computed efficiently in parallel, in contrast to the sequential behavior of recurrent neural networks. Recently, some attention-based works have been proposed in astronomy, showing competitive results with state-of-the-art models (Allam Jr & McEwen 2021; Pimentel et al. 2023; Pan et al. 2022).

Following the advances in NLP to treat sequential data, representation learning, and self-supervised training strategies, we present ASTROMER, a self-supervised model for extracting a general representation of astronomical light curves. In this work we show that using light curve embeddings to train other models, such as classifiers of variable stars, we can match or decrease the number of epochs needed to get the best evaluation metrics. We empirically demonstrate the benefits of using pre-trained representations when training classifiers on small datasets. Using fewer than 100 samples per class, we significantly overcome a long short-term memory (LSTM) trained on light curves directly. Additionally, we provide a python library that includes pre-trained models and the necessary infrastructure to fine-tune and obtain domain specific embeddings. Sharing pre-trained models by the community helps save computational resources and decreases CO2 emissions while improving the performance of automatic learning models. Furthermore, we aim to create different versions of ASTROMER, just as BERT has many variations depending on the target (Polignano et al. 2019; Liu et al. 2019; de Vries et al. 2019; Moradshahi et al. 2019; Masala et al. 2020; Vunikili et al. 2020).

2. Problem statement

Let $X \in \mathbb{R}^{L \times d_x}$ be a single-band light curve where d_x is the number of features with L observations over time. In this case, every observation consists of $d_x = 2$: the magnitude and the modified Julian date (MJD), where the observations occur. Naturally, the number of observations and times vary between different stars, and it strongly depends on the survey science goals. In this representation, the maximum number of points in the light curves remains fixed, even if some of them are shorter than L , in which case we perform zero padding and masking.

The main objective is to train a model $f(X, \mathcal{D}_A, \theta)$ on a massive set of light curves \mathcal{D}_A from a given survey A , with trainable parameters θ . In particular, we propose using learned representations of a transformer-based encoder to create embeddings $Z \in \mathbb{R}^{L \times d_k}$ representing the objects' variability in d_k -dimensional space. We can fine-tune the model's weights to adapt to other surveys and use it to solve downstream tasks, such as classification or regression.

3. Methods

In this section, we describe the main components of ASTROMER, which belongs to the transformer neural network proposed by Vaswani et al. (2017). In particular, we focus on two processes: the self-attention block and the positional encoding (PE). Both methods are effective for capturing relationships between observations to encode light curves.

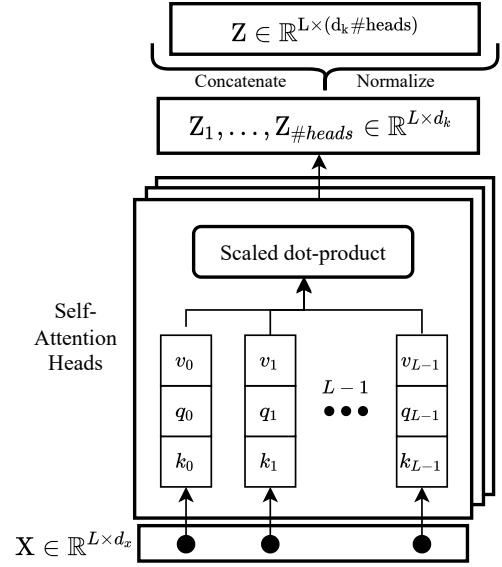


Fig. 1. Self-attention block diagram. Each observation vector $\mathbf{x}_i \in X$ (denoted by solid circles) is projected into three vectors: (k)ey, (q)uery, and (v)alue. Then their similarities are computed by the scaled dot-product between vectors according to Eq. (1).

3.1. Self-attention block

Vaswani et al. (2017) introduced the self-attention mechanism as an alternative to the classical attention technique based on recurrent neural networks (RNNs). The idea was to quantify relationships between observations without conditioning the operations to follow a sequential order. Thus, unlike RNNs, self-attention blocks can be executed in parallel, being more efficient and faster to train.

An attention block consists of multiple self-attention heads that compute the similarities of every input within the sequence to each other. In other words, every head measures the relationship between pairs of observations; the more an observation affects the other, the more attention the model pays. As shown in Eq. (1), we perform a weighted sum over the input values V_i to obtain the attention matrix Z_i using learnable query-key compatibilities ($Q_i K_i^T$):

$$Z_i = \text{softmax} \left(\frac{Q_i K_i^T}{\sqrt{d_k}} \right) V_i, \quad (1)$$

Queries, keys, and values (Q_i, K_i, V_i) are input transformations such that,

$$Q_i = X W_i^q, \quad K_i = X W_i^k, \quad \text{and} \quad V_i = X W_i^v, \quad (2)$$

with W_i^q , W_i^k , and W_i^v being the trainable weight matrices of the i th head, and d_k a hyper-parameter specifying the embedding size of each self-attention head. The normalization factor $\sqrt{d_k}$ in the denominator scales the variance of the product, while the softmax bounds the values to represent a multinomial probability distribution.

The final attention vector is the normalized concatenation of $\{Z_0, \dots, Z_i, \dots, Z_{\#heads}\}$, the output from every i th head in the block. Figure 1 illustrates the self-attention block.

3.2. Positional encoding

Light curves are sequences of brightness measurements as a function of time, generally not equally spaced. The cadence of a

survey is given by the time interval between observations, which is often not periodic.

Temporal information is essential for characterizing light curves since they bring information about the variability of a particular object. For example, RR Lyrae and Delta Scuti stars have short periods, ranging from hours to days, while the periods of type II Cepheids and long-period variable stars can be months or years.

By construction, self-attention layers do not take into account the temporal information of light curves when learning similarities between the observations. In other words, any re-ordering of the magnitudes within a sequence would produce the same attention vector.

A traditional way to include temporal information during self-attention learning consists in modifying the input to incorporate the relative positions into the same vector. We create a representation for the time and add that information to the brightness representation. It should be noted that the positions range from 0 to $L - 1$, where L is the length of the light curve.

In this work, we adjust a classical positional encoder used by Vaswani et al. (2017), to work with the time domain of the light curves¹. The positional encoding (PE) projects an embedding that encodes relative positions and it is consistent to any transformation of them, no matter the number of observations or the first day in MJD where the object began to be observed. It consists of trigonometric functions that codify each observational time t_l directly, from the l th observation at different w_j frequencies.

$$PE_{j,t_l} = \begin{cases} \sin(t_l \cdot \omega_j) & j \text{ is even} \\ \cos(t_l \cdot \omega_j) & j \text{ is odd} \end{cases} \quad (3)$$

In Eq. (3), $j \in [0, \dots, d_{pe} - 1]$, where d_{pe} is the PE dimensionality and w_j is the angular frequency defined as

$$\omega_j = \frac{1}{1000^{2j/d_{pe}}}. \quad (4)$$

Trigonometrical functions are a natural choice for capturing periodic behavior while bringing unique values within $[-1, 1]$, which work efficiently on GPUs (Patel 2020). The sine and cosine functions take $d/2$ different angular frequencies, spanning wavelengths from 2π to $2\pi \times 1000^2$. We make the same remarks as in the original paper regarding the function used. A non-trainable PE is enough to achieve reasonable results, simplifying the network. The value of 1000 gave the best results, rather than the original 10000. The angular frequencies are indexed by the PE dimension number. Larger frequencies are given by the smaller PE dimensions, while the smaller frequencies are given by the larger PE dimensions. Figure 2 shows the resulting PE given the sequence of times of the MACHO light curve F_1.4175.3433. At high frequencies, the PE shows high variability across all the time steps, which show dynamic behavior. In contrast, lower frequencies show low variability.

The PE relates directly to the objects under study. Since the periods of variable stars typically range from 0.1 to 1000 days (Catelan & Smith 2015), the PE should be sensible to periods in this range. To capture short periods, the minimum wavelength should be 2×0.1 . In practice, the survey's cadence conditions the shortest period we can study. Moreover, most objects will not show variability for periods larger than 1000 days. As shown in Fig. 2, the PE does not show any change across the time steps.

¹ We based our implementation on the tensorflow code (Tensorflow 2022) as it achieved a better validation RMSE after pre-training than using the original implementation.

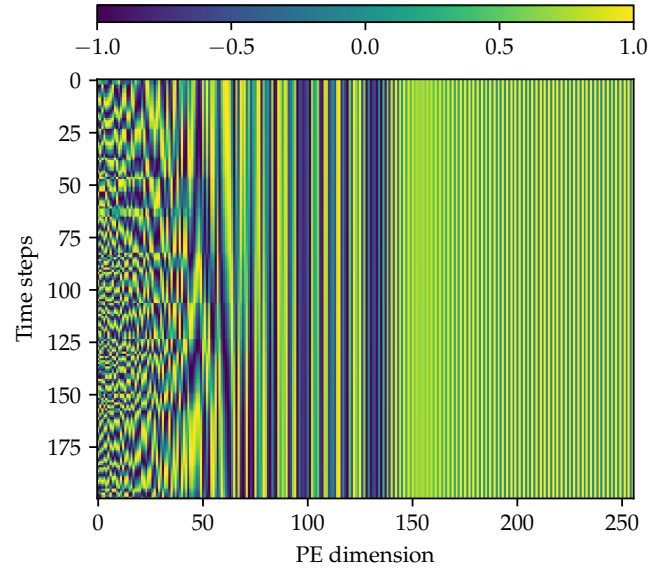


Fig. 2. Resulting positional embedded space from the F_1.4175.3433 MACHO object.

4. Proposed solution

Here we introduce and describe ASTROMER, a transformed-based model capable of producing embeddings from light curves. Embeddings summarize photometric data according to a learned representation that can be easily adapted to new data by re-training (or fine-tuning) the model's weights. In particular, we take inspiration from BERT (Devlin et al. 2018), a self-supervised trained model that does not need human-based annotations to learn representations of the input data.

In the following sections, we explain the main modules and objectives of the proposed model. Before describing the component of the model, we motivate the architecture by describing the target task in Sect. 4.1. Then; Sect. 4.2 presents a general overview of the architecture, where we introduce the ASTROMER modules. Finally, the last sections detail the inner mechanisms of encoding (Sect. 4.3) and decoding (Sect. 4.4) light curves.

4.1. Target task: Modeling of masked light curves

ASTROMER aims to learn representations that summarize and characterize the light curve's domain. It does by reconstructing the input light curve from a subsample of observations. The embedding should be as general as possible to help downstream tasks, such as classification or regression. In this sense, the model needs to adjust its weights to create a vector that is informative enough to recover the input light curve.

In order to force the model to learn contextual information, we hide a random subset of magnitudes within the sequence, and reconstruct them using the remaining part of the light curve. Hidden elements are informed to the model by a mask vector m_i that is different for every i th light curve sample. The loss function then takes the form

$$\mathcal{L}_{oss} = \sqrt{\frac{1}{N-1} \sum_{i=0}^{N-1} \sum_{l=0}^{L-1} m_{il} (x_{il} - \hat{x}_{il})^2}, \quad (5)$$

where N is the number of examples, and L is the length of the input sequences, which i equal for all examples. We notice

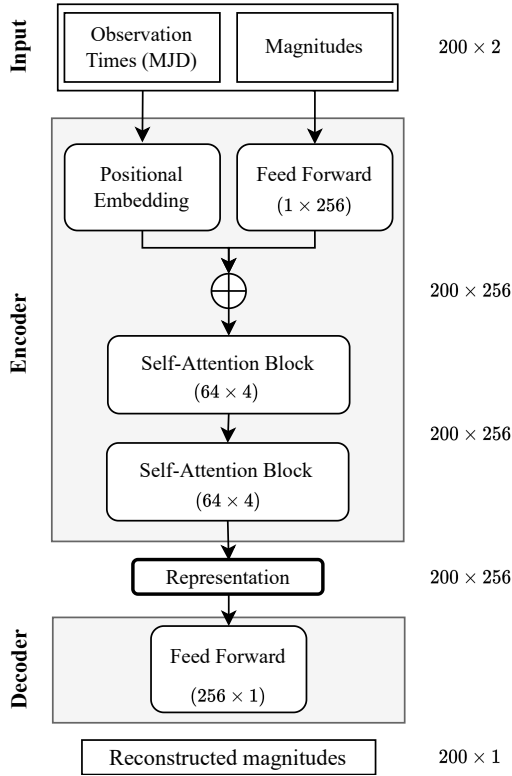


Fig. 3. ASTROMER architecture.

that the loss function containing the actual observation x_{il} and the prediction \hat{x}_{il} is equivalent to the root-mean-square error (RMSE). In Eq. (5), the mask vector m_i only permits summing over the masked values' error; nevertheless, the model reconstructs the entire sequence.

4.2. Model architecture

ASTROMER has an encoder-decoder architecture, as shown in Fig. 3. Input samples are fixed-length light curves of $L = 200$ observations with two descriptors each: observation times and magnitudes. The first part of the encoder is tasked to compute the PE, which receives the time values at each step and projects them into vectors of 256 dimensions. The next step consists in adding magnitudes to the PE without interfering with the temporal encoding. As shown in Fig. 3, we train a feed-forward network (FFN) without hidden layers that transform each magnitude to a vector of size 256. We notice the dimension of the FFN matches the size of the positional embedding. In order to not alter the temporal encoding, the FFN should learn to project brightness information in a way that can be summed with the constant part of the PE, that is to say, after the 100th dimension in Fig. 2². After adding the projected magnitudes and the PE, the new input is a matrix of dimension 200×256 .

The core of ASTROMER takes place on the two self-attention blocks. Each block contains four heads with 64 neurons each. The final representation is the normalized concatenation of the heads of the last block. Then, the resulting embedding is a collection of 200 vectors of size 256, describing the attention of every observation to each other. We notice that $d_{pe} = d_k \#heads = 256$.

² We confirm the assumption in Appendix A, where the model -after training- learned to project magnitudes on the ~ 200 th dimension.

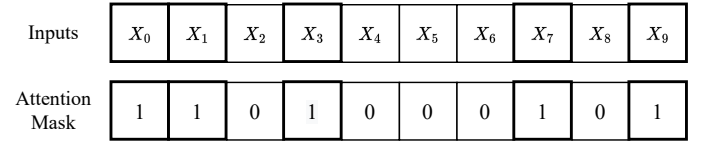


Fig. 4. Preliminary input example with 50% masking. ASTROMER aims to predict masked observations (denoted by “1”s) using only the attention associated with zero-tagged elements.

In the decoder, we take the representation to build the model's output, which during the pre-training consists of a linear FFN with no hidden layers, which reconstructs the input magnitudes. Although the decoder presented in Fig. 3 is exclusive for pre-training ASTROMER, we can use different decoding layers that focus on other downstream tasks.

4.3. Learning representations

ASTROMER uses self-attention blocks to transform light curves into embeddings via learnable parameters. As mentioned in Sect. 4.1, we train ASTROMER to predict a subset of masked observations using the concatenated attention vector at the end of the encoder. The embeddings must capture enough information to reconstruct magnitudes using only the surrounding local context of the hidden observations.

In order to exclude the masked observations from the self-attention blocks, we covert all the dimensions associated with the hidden elements to zero values. Therefore, by modifying Eq. (1), we obtain

$$Z_i = \text{softmax} \left(\frac{Q_i K_i^T + (-10^9)M}{\sqrt{d_k}} \right) V_i, \quad (6)$$

where M is a binary matrix that is one on masked positions and zero otherwise. The M matrix is dynamically created before the encoder. We notice that for the points in the light curve that are hidden (i.e., for which M is 1), the value of the softmax function will be close to zero, making Z_i close to zero. In other words, when the observation is masked, its assigned predicted value will be zero, which represents no attention.

In this work, we mask 50% of the total number of observations per light curve. It is important to note that M should be reshaped to match the square matrix $Q_i K_i^T \in \mathbb{R}^{L \times L}$. Figure 4 shows an example of a ten-observation input.

During training, the mask M provides the model with information about the observations to be predicted. The model can only focus on the target magnitudes and forgets contextual information about unmasked inputs. To avoid learning biased representations conditioned by the masking information only, Devlin et al. (2018) proposed replacing a subset of masked values with random and real elements. In practice, we turn 20% of the masked values from one to zero, while replacing their associated magnitudes with random values from the same sequence. Similarly, we make another 20% of the masked values visible but this time without changing their magnitudes. The unmasked values are still considered in the loss function as they are part of 50% of the initial masked elements. Figure 5 shows the final composition of the input.

Having defined the mask vector, we initiate the forward pass of the encoder by transforming input times and magnitudes to a 256 vector that mixes temporal and brightness information. Finally, we pass the combined input sequences via two self-attention blocks that capture similarities between observations and then output the attention-based embeddings.

Inputs	X_0	X_6	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9
Attention Mask	1	0	0	0	0	0	0	1	0	1
Loss Mask	1	1	0	1	0	0	0	1	0	1

Fig. 5. Final input composition. Following the example in Fig. 4, 20% of the masked values are replaced by random magnitudes while changing the “1”s in the attention mask vector to “0”s. Similarly, we make another 20% of the masking visible, keeping the actual observations. Doted-line squares indicate both random and real observations. At this point, we should keep a second mask containing the initial 50% masked values to evaluate the loss function.

4.4. Decoding embeddings

The decoder is essential to adjust the encoder weights. The output conditions the purpose of the embedded representation. Despite considering only masked values in the loss function, the decoder reconstructs the entire sequence of magnitudes. As we explained in Sect. 4.1, using a mask vector, we give a value of zero to the observations that are not part of the loss.

The embedding is a collection of 200 vectors (the length of the input sequence) of size 256 (model dimensionality). Every vector describes the relationship between a particular observation and the rest of the sequence. In order to recover the original magnitude, we apply a FFN without hidden layers and no activation on all the vectors of size 256 separately. Although we transform every vector separately, the network weights remain the same.

Finally, it is important to mention that the decoder only works during the pre-training and fine-tuning of ASTROMER. After that, we usually use the encoder along with a new decoder specializing in another specific task, such as the classification of variable stars. At this step, the encoder can also be trained together with the new decoder. However, the encoder cannot change drastically, as it would forget the learned representations.

5. Data description

This section describes the data involved in the pre-training and fine-tuning of the ASTROMER model. Every light curve contains magnitudes and observation times in MJD.

ASTROMER takes advantage of a massive set of light curves, which are not necessarily labeled. However, to evaluate a downstream task (i.e., classification), we use catalogs of variable stars with their corresponding labels. This work simulates the science case of having a small subset of labeled samples and employs them to fine-tune the model, even when the class information is unnecessary.

5.1. Unlabeled dataset

R-band light curves of the Massive Compact Halo Object survey (MACHO; Alcock et al. 2000) were collected from the Galactic bulge (fields 1 and 10), and the Large Magellanic Cloud (LMC, Fields 101–104). The photometry was taken as it was provided, without any additional processing.

Given the nature of photometric observations, the observed flux of every star will exhibit some kind of variability, such as the intrinsic observational noise. As such, we cannot expect every star to be variable. Even though non-variable and noisy

Table 1. MACHO catalog distribution.

Tag	Class name	# of sources
Cep_0	Cepheid type I	1182
Cep_1	Cepheid type II	683
EC	Eclipsing binary	6824
LPV	Long period variable	3046
RRab	RR Lyrae type ab	7397
RRc	RR Lyrae type c	1762
Total		20 894

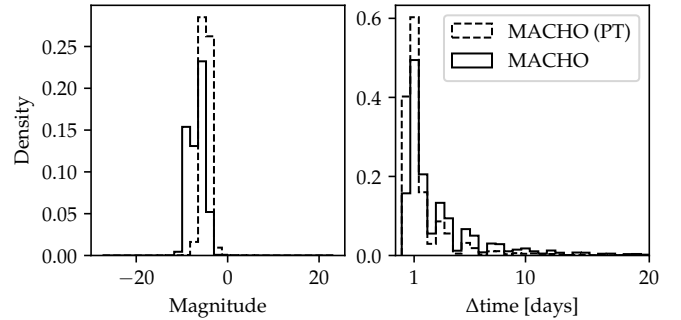


Fig. 6. Comparison between the labeled and unlabeled pre-training MACHO datasets. The continuous lines represent the labeled MACHO dataset containing variable stars from Alcock et al. (2003). Left: Distribution of light curve magnitudes. Right: Distribution of the differences between observation times.

data can be helpful to regularize weights (Bishop 1995), we discard some of the light curves that show white noise behavior (i.e., $|\text{Kurtosis}| > 10$, $|\text{Skewness}| > 1$, and $\text{Std} > 0.1$). Removing white noise samples avoids training on a dataset dominated by irrelevant information while keeping variable objects that might contain informative variability. After this filtering process, 1 529 386 light curves are left, with a cadence mean of 2.9 ± 17.3 days, which fits the lower bound of the PE wavelengths mentioned at the end of Sect. 3.2.

5.2. Labeled datasets

We employ 20 894 labeled variable stars from the MACHO survey (Alcock et al. 2003) to fine-tune and evaluate the downstream task. Table 1 shows the class distribution.

Similar to the unlabeled dataset, this catalog contains objects from the LMC. Consequently, it contains light curves quite similar to the pre-training ones. However, labeled objects come from all the survey’s fields, while the pre-training dataset is a subset (fields 1, 10, and 101–104). Figure 6 shows the distribution of magnitudes and the observational time difference (Δtime) from the labeled and unlabeled (pre-training) datasets. In the left histogram, we can see the pre-training dataset overlaps one of the modes from the labeled dataset distribution. Similarly, distributions of delta times follow the same trend along the x -axis. However, they are slightly different for values of delta time greater than one day, where a greater density from the labeled dataset is observed. The labels were updated to conform to modern classifications. In particular, we group the LPV classes into one category and remove the categories “RR Lyrae and GB blends” and “RRe”. The former are excluded because of their small sample size, and the latter because RRe is not considered a subclass on its own, based on later studies (Catelan 2004).

Table 2. OGLE-III catalog distribution.

Tag	Class name	# of sources
EC	Eclipsing binary	6862
ED	Detached binary	21 503
ESD	Semidetached binary	9475
Mira	Mira	6090
OSARG	Small-amplitude red giant	234 932
RRab	RRLyra type ab	25 943
RRc	RRLyra type c	7990
SRV	Semi-regular variable	34 835
cep	Cepheid	7836
dsct	Delta scuti	2822
Total		358 288

Table 3. Surveys and the filters used in this work.

Survey	Filter name	λ_{\min} $\circ A$	λ_{\max} $\circ A$
MACHO ⁽¹⁾	R	6300	7600
OGLE ⁽²⁾	I	7270	8750
ATLAS ⁽³⁾	Orange	5620	8200

Notes. In [Szymański et al. \(2011\)](#), it is noted that OGLE-III and OGLE-IV do not have the same I filters, exhibiting a wide transparency “wing” toward the infrared part of the spectrum. Although different, we use [Udalski et al. \(2015\)](#) numbers as a reference, as they are intended to show the approximate ranges of the filters.

References. (1) [Alcock et al. \(1999\)](#), (2) [Udalski et al. \(2015\)](#), (3) [Tonry et al. \(2018\)](#).

Table 4. ATLAS catalog distribution.

Tag	Class name	# of sources
CB	Close binaries	80 218
DB	Detached binary	28 767
Mira	Mira	7370
Pulse	RR Lyrae, δ -Scuti, Cepheids	25 021
Total		141 376

We also tested ASTROMER on the Optical Gravitational Lensing Experiment (OGLE-III; [Udalski 2003](#)) and ATLAS [Heinze et al. \(2018\)](#) catalogs of variable stars. The OGLE-III data used were selected by [Becker et al. \(2020\)](#) and contain 358 288 labeled variable stars, and correspond to I-band light curves with a mean cadence of 3.8 ± 14.6 days. Table 2 shows their class distribution. The OGLE-III catalog is chosen because its observations are not captured in the same filter as MACHO, but close to its wavelength range, as shown in Table 3.

The ATLAS dataset, published by [Heinze et al. \(2018\)](#), contains labeled and unclassified objects, as well as a dubious class, which amounts to roughly 10% actual variable stars and 90% instrumental noise, according to their estimates. From this dataset, 422 630 light curves are used, measured in the orange passband, as seen in Table 4, with a mean cadence of 4.7 ± 19.1 days. The reported classes are grouped to obtain labels similar to the other datasets. In particular, we group detached eclipsing binaries with either a full or half period identified into close binaries (CB) and the same for detached binaries (DB).

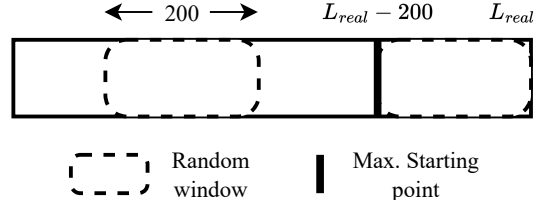


Fig. 7. Windows sampling diagram. The rectangle represents the entire light curve whose length is L_{real} . The dashed line denotes the sampled windows of size 200 observations. Windows are randomly generated along the light curve. However, we constrain the starting point to be smaller than $L_{real} - 200$.

We do not use the remaining objects, as their labels are based on Fourier analysis and do not correspond exactly to astrophysical categories.

5.3. Preprocessing

Neural networks must work with tensors holding equal-length samples. It is not the case for light curves that differ on the number of measurements. As such, it is necessary to set a maximum number of observations for all sequences, padding with zero values if necessary. Since we do not want to constrain the encoder to learn from long series exclusively, and 99.52% of the unlabeled dataset is longer than 200 (see Appendix B), we set 200 as the maximum light curve length to be fed to the model. If the light curve is longer than 200, we sample temporal windows starting from a random position, at least 200 observations behind the last measurement. Figure 7 shows a diagram that exemplifies the process. On the other hand, if the light curve has fewer than 200 observations, we pad it with zero values after the last point in order to obtain a sequence of length 200. Those filler observations are masked during training to be excluded from both the attention vector and the loss function.

After generating windows, we independently subtract the mean from each sample. It implies both magnitude and time vectors have zero mean. We scale neither the magnitude nor time windows by their standard deviation. Scaling by time dispersion would lose some of the interpretability for the positional encoder. Standardizing magnitudes may lose amplitude-related information that is important for discriminating between some classes.

6. Training strategy

The following sections describe the training strategy consisting of two steps: pre-training and fine-tuning. In the pre-training, 60% of the unlabeled samples are used for training, 20% for validation, and 20% for testing. Alternatively, we set 100 samples per class of the labeled dataset to create the test set for evaluating the fine-tuning step. The rest of the labeled dataset is divided into 80% and 20% for training and validation, respectively. We recall that the pre-training dataset is different from the labeled one, so the testing samples used to evaluate the fine-tuning are never seen by ASTROMER.

Additionally, we study the scientific case of having small target datasets to fine-tune and perform downstream tasks. Therefore, we do not change our 100 samples per class test set, but only the training and validation subsets. First, we sample 20, 50, 100, and 500 objects per class from each labeled dataset. Then we divide them by 80/20 for training and validation. We aim to see the impact of the number of samples when adjusting pre-trained weights on a different survey.

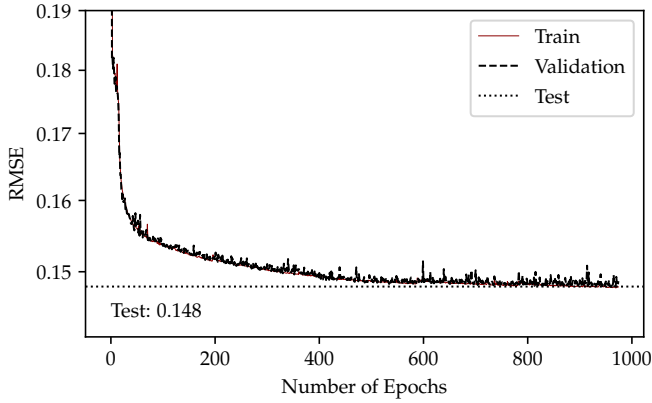


Fig. 8. Pre-training learning curves using the MACHO unlabeled dataset.

6.1. Pre-training

We implemented ASTROMER using Tensorflow 2 (Abadi et al. 2015). Both pre-training and the other experiments can be reproduced by following the official implementation on Github³.

The pre-training constitutes the first stage of representation learning. It defines the preliminary data requirements, such as cadence and filter, which condition other target domains. In this case, we use the massive unlabeled dataset from the MACHO survey presented in Sect. 5.1.

ASTROMER weights are initialized using the Xavier uniform initializer (Glorot & Bengio 2010). Training big models such as ASTROMER uses a large amount of hardware resources and computational time. Despite this, once the weights are adjusted, they can be shared, avoiding having to re-train ASTROMER from scratch.

As mentioned in Sect. 5.3, the light curves are split into windows of 200 consecutive observations, increasing the effective number of training samples to 6 201 030. An epoch is completed after training all samples once, in batches of 5000. We used Adam (Kingma & Ba 2014) with a learning rate of 10^{-3} .

The training performance is evaluated on the validation dataset and the results on the test set. Early stopping is used to stop the training after 40 epochs without improving the validation loss. Once the training is finished, the weights corresponding to the lowest RMSE are saved.

Figure 8 shows the learning curves of the training process associated with the architecture in Fig. 3. We also trained on smaller models of 64 and 128 attention sizes. However, we chose the 256-dimensional setting as it reached the minimum validation loss on the unlabeled dataset. The model achieved the lowest RMSE close to epoch 1000⁴, obtaining an RMSE of 0.148 on the testing samples (dotted line).

6.2. Fine-tuning

After pre-training on the MACHO unlabeled dataset, the model has learned much of the variability patterns of the light curves. However, it is still necessary to incorporate downstream task-related information. In this work, we use catalogs of variable stars described in Sect. 5.2 to fine-tune the attention-based embeddings.

The loss function and self-supervised strategy remain the same (i.e., RMSE and masked self-attention technique).

³ <https://github.com/astromer-science>

⁴ ~10 days training on a Nvidia A100 GPU.

Table 5. Fine-tuning results.

Dataset	# Epochs	# Time	RMSE
MACHO (PT)	970	9 days 17 h	-/0.15
MACHO	115	23 min	0.09/0.10
ATLAS	147	9 h 58 min	0.07/0.22
OGLE-III	244	1 day 8 h	0.06/0.08

Notes. As a reference, the first row shows the performance of the pre-trained model used to initialize weights in the fine-tuning. From the second row forward, the first column indicates the name of the labeled dataset used to fine-tune ASTROMER. The second and third columns show the number of epochs and training time the models spend to converge. The last column is the testing RMSE evaluated on the fine-tuned (left) and pre-trained (right) models. For the fine-tuning metrics, we employ balanced testing with 100 objects per class from each labeled dataset.

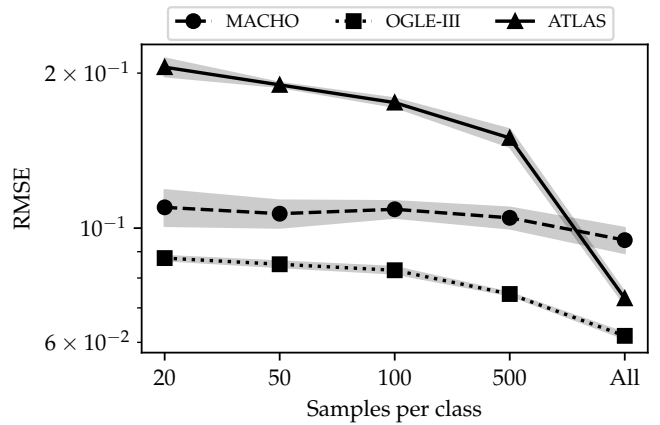


Fig. 9. Test RMSE from models fine-tuned on subsets of labeled datasets.

Similarly, the same hyper-parameters from the pre-training stage are used.

Table 5 shows the number of epochs and total time to fine-tune ASTROMER on each labeled dataset using all samples, comparing it to the pre-training stage. A significant difference to the pre-training time is evident. The fine-tuned models converge faster depending on their similarities and the number of examples. For instance, the model fine-tuned on ATLAS light curves takes more time than the one trained on the MACHO labeled dataset, principally due to the survey differences and the number of samples that allow the model to get a much lower RMSE than the MACHO labeled dataset. At the same time, similar behavior is observed on the OGLE fine-tuning. However, the minor improvements in OGLE’s RMSE suggest the training time is highly dominated by the significant number of samples in the catalog (see Table 2).

On the other hand, Fig. 9 shows the RMSE for the models fine-tuned on smaller datasets of 20, 50, 100, and 500 labels per class. The testing set contains the same 100 objects per class, independent of the training subsets. We can see a descending trend in the RMSE while increasing the number of samples in the fine-tuning, which is expected. However, there is not much improvement along the training subsets in the MACHO labeled dataset since it comes from the same survey as the pre-training light curves. Validation learning curves associated with the training process in all datasets can be found in Appendix C.

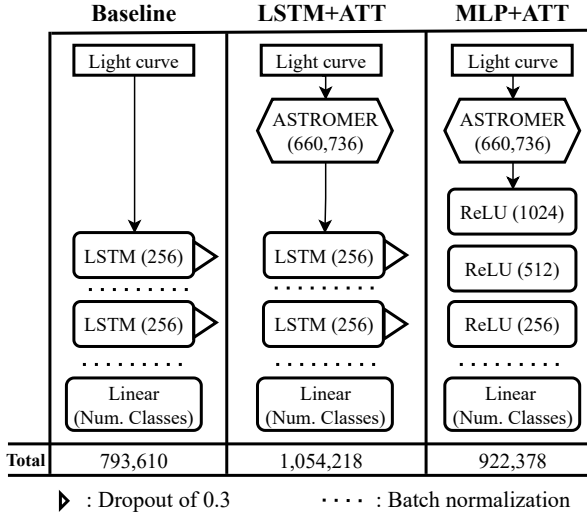


Fig. 10. Classifiers architectures. Architecture names with the +ATT tag refer to models trained on ASTROMER embeddings. In contrast, the baseline uses only the light curves for training. The number in parenthesis is the number of units from the corresponding layer. We notice that the 256 in the LSTM corresponds to the states' size of the RNN. ReLU and Linear correspond to the activations for each feed-forward layer. The dimensionality of the last linear layer depends on the number of classes from each labeled dataset. The total number of parameters at the bottom does not include the ASTROMER size.

7. Specific task: Classification

ASTROMER is designed as a general embedding extractor that can be fine-tuned to solve downstream tasks such as classification or regression. In this work, we cover the problem of classifying variable stars from different surveys. As explained in Sect. 6, the model is fine-tuned in subsets of 20, 50, 100, and 500 samples per class, as well as using the full catalog (see Sect. 5). Fine-tuned models are then used to transform light curves into attention-based embeddings. In practice, we use ASTROMER at the beginning of the classifiers as an additional layer (see Fig. 10).

We selected a baseline model from the literature (Donoso-Oliva et al. 2021) to compare the benefits of using the embeddings as opposed to other deep learning approaches. It consists of two layers of LSTM with a memory cell and hidden state of 256 units each. Following the original architecture, we normalize each recurrent layer's output (Ba et al. 2016) and then apply a dropout (Semeniuta et al. 2016) of 0.3 over them. We do not validate this architecture since it was already tested in the previous work. However, we looked for the best learning rate on the MACHO light curves since we changed conditions, such as the batch size and normalization techniques. As shown in Appendix D, we confirm that a 10^{-3} learning rate achieved the best performance in terms of the validation error. In the results shown in Fig. 10, we employed the same LSTM architecture but training with ASTROMER embeddings.

In addition to the LSTM classifier, we tested a multilayer perceptron (MLP+ATT in Fig. 10) with three feed-forward hidden layers of 1024, 512, and 256 neurons, and used a rectified linear activation function (ReLU) in each of them. Since MLPs cannot process time, we use the average along the step dimension ($L = 200$) of the embedding matrix⁵. By doing this, we

⁵ We also tried taking the first, last, and arbitrary attention vectors to feed the input of the MLP. However, the results were not better than using the average.

collapse all the encoded observations in a single vector of size 256 that fits the input dimensionality of the first hidden layer of the MLP+ATT classifier (feed-forward layer with 1024 units from Fig. 10).

For small datasets, freezing the weights of ASTROMER allows us to optimize a classifier without training a huge number of weights. The ASTROMER encoder can also be trained in tandem with the classification network, to better tune weights to solve a specific task, which can be seen as another fine-tuning step. The following experiments evaluate these two approaches, that is, training and non-training of the encoder layer of ASTROMER.

To measure the classification performance, we use the $F1$ score metric,

$$F1 = \frac{1}{K} \sum_{k=0}^{K-1} 2 \times \frac{\text{Precision}_k \times \text{Recall}_k}{\text{Precision}_k + \text{Recall}_k} \quad (7)$$

where K is the number of classes, and

$$\text{Recall}_k = \frac{\text{TP}_k}{\text{TP}_k + \text{FN}_k} \quad \text{Precision}_k = \frac{\text{TP}_k}{\text{TP}_k + \text{FP}_k} \quad (8)$$

In Eq. (8), TP, FN, and FP are true positive, false negative, and false positive cases, respectively. Intuitively, the precision score identifies how many predicted classes are actually valid, and the recall indicates the number of objects in the testing set the model could identify correctly.

Figure 11 shows the scores obtained by the classifiers by (a) freezing the encoder and (b) training the encoder. In addition, a third case (c) is included where all the light curves in the dataset were used to perform fine-tuning, but still doing the classification on the smaller subsets. As in experiment (b), the third column allows gradients to flow into the ASTROMER encoder while training the classifiers. These three scenarios explore often-used strategies when implementing pre-trained models.

The results show that models trained on ASTROMER embeddings perform better than the baseline in MACHO and OGLE-III across all the experiments. For ATLAS, the difference is smaller but not worse than the baseline. In general, the LSTM trained on attention vectors performs better than the MLP+ATT. However, the MLP+ATT outperforms the baseline and approaches the LSTM performance as the number of samples per class is increased, in all datasets.

It should be noted that in (c), even fine-tuning with all the light curves, the score improvement is marginal compared to in (b). Similarly, minor improvements can be seen in the classification metrics when optimizing the encoder layer instead of freezing it (i.e., Cols. (b) and (c) in Fig. 11). The MLP+ATT classifier shows the most notorious gain when training with more than 100 samples.

We evaluate the improved learning speed when using ASTROMER. Figure 12 shows the validation loss associated with the best classifier among the three folds for each dataset and the training scenarios (i.e., freezing (a) and training (b) the ASTROMER encoder). We evaluate learning speed on the most significant subset, which contains 500 objects per class.

Attention-based models generally take fewer epochs to achieve better results than the baseline. When the encoder stays frozen, the validation curve of the LSTM+ATT is significantly shorter than the other methods, achieving a high $F1$ score, as shown in Fig. 11. In contrast, when training the encoder, the LSTM+ATT takes almost the same number of epochs as the baseline. However, according to Fig. 11, training the encoder

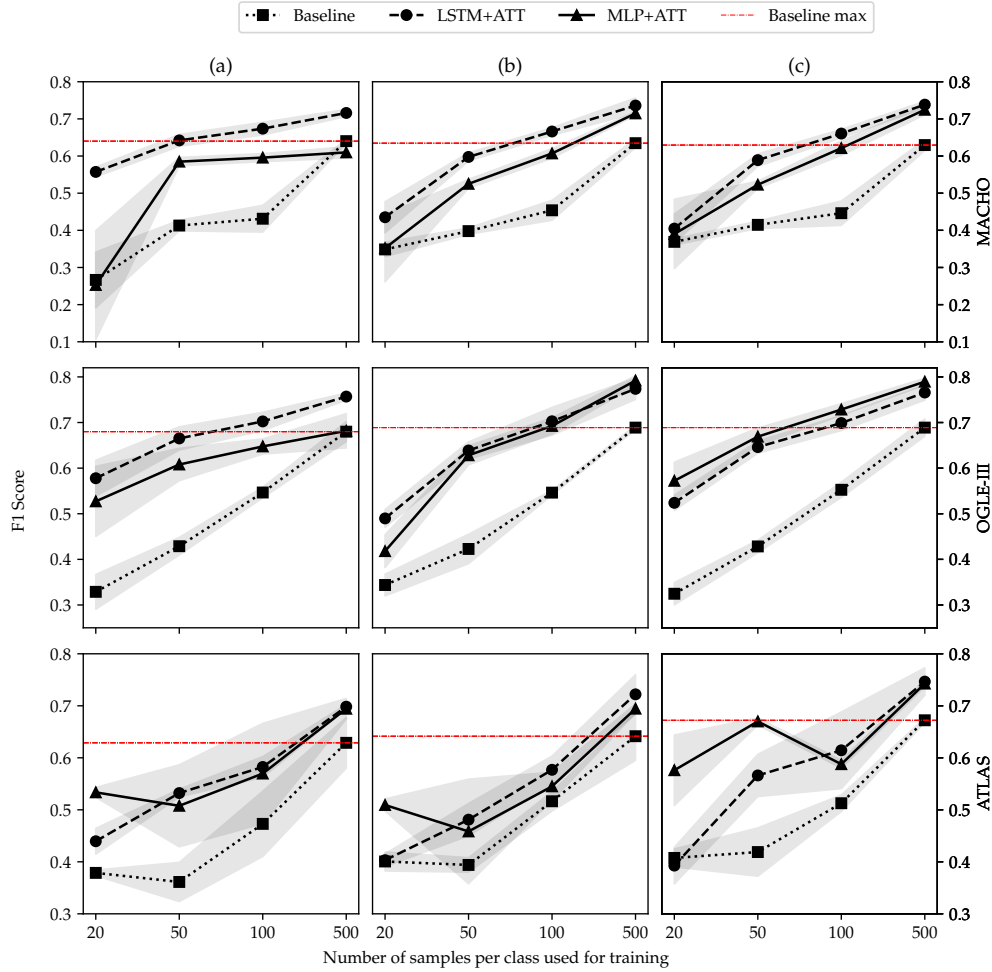


Fig. 11. Testing $F1$ scores for an LSTM trained on light curves directly (baseline) and models trained on ASTROMER embeddings (LSTM+ATT and MLP+ATT). The gray shading represents the standard deviation of the three cross-validation split. Each row corresponds to the experiments on each survey, MACHO, OGLE, and ATLAS, respectively. In (a), we fine-tune ASTROMER and optimize classifiers on smaller subsets of 20, 50, 100, and 500 samples per class. The weights of ASTROMER are kept frozen when classifying. However, in (b), we allow gradients to flow into ASTROMER. The third case, (c), shows the results of fine-tuning with the entire set of light curves and classifying on smaller subsets, training ASTROMER simultaneously.

using the LSTM+ATT does not improve the $F1$ score. The findings are reversed when using the MLP+ATT, in this case, training the encoder decreases the number of epochs while improving the $F1$ score in larger subsets.

8. Discussion

Our results have shown the benefits of using pre-trained models to learn representations of light curves. We successfully adapted and applied an NLP self-supervised technique (Devlin et al. 2018) to the domain of light curves. It offers a solution for learning representations when labels are in limited quantity. Furthermore, it is an alternative to fully unsupervised models that focus more on clustering than making predictions from the data.

Our representation can be extended to light curves coming from different surveys. Using ASTROMER on new data requires a short training process called fine-tuning to adjust the representation and minimize the RMSE. A lower RMSE implies a better representation and, subsequently, a better result in downstream tasks.

The size of the fine-tuning dataset depends on the similarity between the pre-training and fine-tuning magnitudes and cadence distributions. For illustration, Fig. 11c shows no

difference in the $F1$ scores when all the light curves in the dataset were used to fine-tune the encoder. Indeed it is more important to bring the classifier more labels than more light curves in the fine-tuning. As shown in Table 5, minor improvements in the RMSE were obtained when fine-tuning ASTROMER on MACHO and OGLE-III, both of which are more similar to the pre-training dataset than ATLAS. In this scenario, the fine-tuning stage can even be omitted without affecting the results.

Regarding the downstream task, we demonstrated the power of using embeddings to train classifiers, instead of directly using light curves to predict classes. We evaluated the classification performances in three commonly used scenarios using limited labeled data. As seen in Fig. 11, we outperform a recurrent neural network classifier trained on light curves (baseline) in all the experiments. The improvements in the $F1$ score became more significant when fewer than 500 labels per class were provided.

ASTROMER can learn discriminative features from the pre-training and fine-tuning steps using only a self-supervised task. As shown in Fig. 11, the LSTM+ATT can predict classes with high $F1$ scores without the need to train the encoder along the classifier. Furthermore, its performance remains the same if the encoder is trained in tandem, independent of the number of samples per class.

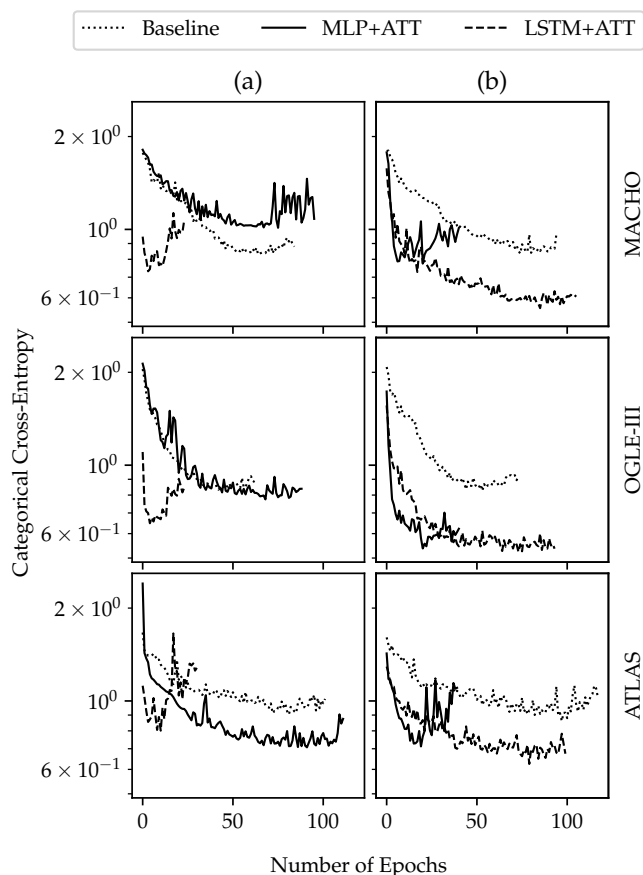


Fig. 12. Best model validation learning curves for classifiers trained on 500 samples per class on each catalog. The columns show the science cases, i.e., freezing (a) and training (b) ASTROMER when classifying.

While the LSTM+ATT uses its hidden state to capture long- and short-term dependencies on the attention vectors over time, the MLP+ATT deals with their average, thus losing global context information. However, when training the encoder along with MLP+ATT using more than 100 samples per class, the encoder is able to counteract this effect, as shown in Fig. 11.

In both models, the discriminative features learned by the ASTROMER encoder allow the classifiers to converge faster⁶. From the best performing scenarios, the LSTM+ATT converges faster than the baseline when freezing the encoder weights, as it only needs to recognize discriminative patterns in the embeddings. On the other hand, the MLP+ATT only needs to mitigate the effect of averaging the attention vectors. This task is simpler than learning a representation from scratch, which is the case of the baseline classifier. Specifically, the baseline has to learn how to extract an informative representation and learn discriminative patterns simultaneously. In the training process, any change in the representation will impact the subsequent layers, increasing the number of epochs to converge. This behavior is related to the internal covariance shift (Ioffe & Szegedy 2015).

Finally, even though this work’s main contribution is to help train downstream models on small datasets, we achieved competitive results against Donoso-Oliva et al. (2021) using OGLE-III unfolded light curves (88.1% vs. 88.0% ours). For the MACHO

⁶ It is important to note that shorter training times using shared pre-trained representations imply the use of less powerful hardware and fewer resources, decreasing the time and energy consumption when training deep learning models (Dhar 2020).

catalog, we compare ASTROMER-based classifiers against the best (78.1% accuracy) single-band non-period informed model from Jamal & Bloom (2020). We achieved similar results, 78.2% and 76.7% classification accuracy for the LSTM+ATT and MLP+ATT, respectively. It is important to mention that all models from the literature use random smaller partitions, while we use a balanced testing set of 100 objects per class, which is more representative.

9. Python package

We provide a Python package, ASTROMER, which includes the pre-trained weights obtained in Sect. 6.1. More pre-trained models will be uploaded in the future, either from the ASTROMER team or the community.

The stable version `0.x.y` matches the code of this paper. Variables `x` and `y` refer to the minor changes and patches of the current version `0`. Major changes will not be related to this work but must be duly evaluated to guarantee at least the same performance. ASTROMER `v0` can be easily installed from the Python Package Index (PyPI) repository as follows:

```
pip install ASTROMER
```

The principal module `ASTROMER.models` includes the models associated with the different encoders. So far, we only have the `SingleBandEncoder`, corresponding to the ASTROMER model trained on the MACHO unlabeled dataset. To import and use ASTROMER, we can type:

```
from ASTROMER.models import SingleBandEncoder
model = SingleBandEncoder()
```

where `model` is a new instance of `ASTROMER`. To load pre-trained weights, we must use the `from_pretraining()` method from the model object,

```
model = model.from_pretraining("macho")
```

The name in parenthesis matches the zip file name in the public repository⁷.

The simplest way to obtain embeddings is via a collection of NumPy arrays, including light curve information in the form:

```
data = [ np.array([[5200, 0.3, 0.2],
                  [5300, 0.5, 0.1],
                  [5400, 0.2, 0.3]]),
         np.array([[4200, 0.3, 0.1],
                  [4300, 0.6, 0.3]]) ]
```

where the axes (from left to right) of the NumPy array are the times, magnitudes, and standard deviation of the magnitudes. Then, to get the embeddings, we use the `encode` method:

```
att_vectors = model.encode(data)
```

where `att_vectors` is a list containing the embeddings for each sample in data.

ASTROMER can be easily trained using the `fit()` method from the `SingleBandEncoder` instance:

```
model.fit(training_data,
          validation_data,
          epochs=10)
```

⁷ <https://github.com/astromer-science/weights>

Within the fit method, `training_data` and `validation_data` are TensorFlow datasets. We provide a function in the `ASTROMER.preprocessing` module to format datasets,

```
from ASTROMER.preprocessing import load_numpy
```

```
training_data = load_numpy(data,
                           batch_size=2,
                           msk_frac=.5,
                           rnd_frac=.2,
                           same_frac=.2,
                           max_obs=200)
```

It should be noted that `data` is the same collection of NumPy arrays we used in the previous examples. In the `load_numpy()`:

- `batch_size`: Number of samples to process in one forward pass of the training;
- `msk_frac`: Fraction of the sequence to be masked;
- `rnd_frac`: Fraction of the mask to be replaced with random values (see Sect. 4.3);
- `rnd_frac`: Fraction of the mask to be replaced with actual values (see Sect. 4.3);
- `max_obs`: Maximum windows length (see Sect. 5.3).

For more information we open project repositories⁸ where data, tutorials, documentation, and contributions via issuing pull requests can be found. These contributions can be in the form of functionalities or pre-trained models.

10. Conclusion

We present ASTROMER, a single-band embedding for light curve representation. It is based on the BERT NLP model, which codifies observations into attention vectors via self-supervised learning, taking advantage of the massive unlabeled volume of data. In this work, we pre-train ASTROMER on millions of R-band light curves from the MACHO survey.

ASTROMER can be fine-tuned on specific domain datasets to solve downstream tasks, such as classification or regression. Here, we use labeled catalogs from the MACHO, OGLE-III, and ATLAS surveys to evaluate the effect of using embeddings on the classification of variable stars. By training MLP and LSTM classifiers on embeddings, we outperform a baseline LSTM network trained on light curves. We evaluated classification performances of fine-tuned models using 20, 50, 100, and 500 samples per class, obtaining better F1 scores in all experiments. Additionally, we showed that the embeddings-based classifiers achieved competitive scores against state-of-the-art solutions. In terms of training times, models trained on ASTROMER representations took fewer epochs than the baseline LSTM classifier to reach the lowest validation loss.

Our self-supervised approach only considers the reconstruction of magnitudes, not including other tasks, such as the next sentence prediction applied in BERT. As a result, the final representation may not satisfy other downstream tasks correctly. For instance, the consequences of losing global context information were evidenced in the difference between the MLP and the LSTM classifiers, both using attention vectors. In order to make embeddings more informative, we will explore adding new tasks during the pre-training and fine-tuning steps in a future work.

Finally, we provided a python library including MACHO pre-training weights and fine-tuned models used in this work. Moreover, our library allows users to pre-train, fine-tune, and get

embeddings on new single-band datasets. We aim to create a collaborative research environment where pre-trained ASTROMER weights can be shared, saving computational resources and improving state-of-the-art models.

Acknowledgements. This research was supported by the ANID Millennium Science Initiative ICN12 009, awarded to the Millennium Institute of Astrophysics; FONDECYT Initiation No. 11191130 (G.C.V.); the Patagón supercomputer of Universidad Austral de Chile (FONDEQUIP EQM180042) and CONICYT-PFCHA/Doctorado Nacional/2018-21181990.

References

- Abadi, M., Agarwal, A., Barham, P., et al. 2015, *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*, software available from [tensorflow.org](https://www.tensorflow.org)
- Alcock, C., Allsman, R., Alves, D. R., et al. 1999, *PASP*, **111**, 1539
- Alcock, C., Allsman, R., Alves, D. R., et al. 2000, *ApJ*, **542**, 281
- Alcock, C., Allsman, R., Alves, D., et al. 2003, *VizieR Online Data Catalog*: **II/247**
- Allam Jr, T., & McEwen, J. D. 2021, ArXiv e-prints [arXiv:2105.06178]
- Ba, J. L., Kiros, J. R., & Hinton, G. E. 2016, ArXiv e-prints [arXiv:1607.06450]
- Becker, I., Pichara, K., Catelan, M., et al. 2020, *MNRAS*, **493**, 2981
- Bishop, C. M. 1995, *Neural Comput.*, **7**, 108
- Catelan, M. 2004, in *International Astronomical Union Colloquium*, 193 (Cambridge University Press), 113
- Catelan, M., & Smith, H. A. 2015, *Pulsating Stars* (John Wiley & Sons)
- Charnock, T., & Moss, A. 2017, *ApJ*, **837**, L28
- de Vries, W., van Cranenburgh, A., Bisazza, A., et al. 2019, ArXiv e-prints [arXiv:1912.09582]
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. 2018, ArXiv e-prints [arXiv:1810.04805]
- Dhar, P. 2020, *Nat. Mach. Intell.*, **2**, 423
- Donoso-Oliva, C., Cabrera-Vives, G., Protopapas, P., Carrasco-Davis, R., & Estevez, P. 2021, *MNRAS*, **505**, 6069
- Glorot, X., & Bengio, Y. 2010, in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, JMLR Workshop and Conference Proceedings, 249
- Heinze, A., Tonry, J. L., Denneau, L., et al. 2018, *AJ*, **156**, 241
- Ioffe, S., & Szegedy, C. 2015, in *International conference on machine learning*, *PMLR*, 448
- Ivezić, Ž., Kahn, S. M., Tyson, J. A., et al. 2019, *ApJ*, **873**, 111
- Jamal, S., & Bloom, J. S. 2020, *ApJS*, **250**, 30
- Kingma, D. P., & Ba, J. 2014, ArXiv e-prints [arXiv:1412.6980]
- Kremer, J., Stensbo-Smidt, K., Gieseke, F., Pedersen, K. S., & Igel, C. 2017, *IEEE Intell. Syst.*, **32**, 16
- LeCun, Y., Bengio, Y., & Hinton, G. 2015, *Nature*, **521**, 436
- Liu, Y., Ott, M., Goyal, N., et al. 2019, ArXiv e-prints [arXiv:1907.11692]
- Liu, X., Zhang, F., Hou, Z., et al. 2021, *IEEE Trans. Knowl. Data Eng.*, **1**
- Masala, M., Ruseti, S., & Dascalu, M. 2020, in *Proceedings of the 28th International Conference on Computational Linguistics*, 6626
- Moradshahi, M., Palangi, H., Lam, M. S., Smolensky, P., & Gao, J. 2019, ArXiv e-prints [arXiv:1910.12647]
- Naul, B., Bloom, J. S., Pérez, F., & van der Walt, S. 2018, *Nat. Astron.*, **2**, 151
- Pan, J., Ting, Y.-S., & Yu, J. 2022, ArXiv e-prints [arXiv:2207.02787]
- Patel, S. 2020, PhD thesis, The Cooper Union for the Advancement of Science and Art
- Pimentel, Ó., Estévez, P. A., & Förster, F. 2023, *AJ*, **165**, 18
- Polignano, M., Basile, P., De Gemmis, M., Semeraro, G., & Basile, V. 2019, in *6th Italian Conference on Computational Linguistics, CLiC-it 2019*, 2481, CEUR, 1
- Sánchez-Sáez, P., Reyes, I., Valenzuela, C., et al. 2021, *AJ*, **161**, 141
- Semeniuta, S., Severyn, A., & Barth, E. 2016, ArXiv e-prints [arXiv:1603.05118]
- Szymański, M., Udalski, A., Soszyński, I., et al. 2011, *Acta Astron.*, **61**, 83
- Tensorflow 2022, Positional encoding Transformer model for language understanding
- Tonry, J., Denneau, L., Heinze, A., et al. 2018, *PASP*, **130**, 064505
- Tsang, B. T.-H., & Schultz, W. C. 2019, *ApJ*, **877**, L14
- Udalski, A. 2003, *Acta Astron.*, **53**, 291
- Udalski, A., Szymański, M., & Szymański, G. 2015, *Acta Astron.*, **65**, 1
- Vaswani, A., Shazeer, N., Parmar, N., et al. 2017, in *Advances in Neural Information Processing Systems*, 5998
- Vunikili, R., Supriya, H., Marica, V. G., & Farri, O. 2020, in *IberLEF@ SEPLN*, 505

⁸ <https://github.com/astromer-science>

Appendix A: Summing magnitudes to positional encoding

According to the classical approach (Vaswani et al. 2017, Devlin et al. 2018), magnitudes and times should be added to build the input of self-attention blocks. However, summing magnitudes to the PE information is not straightforward as we can interfere with the encoded variability of times. In principle, most of the values in the last dimensions of the PE are constant so that we can sum magnitudes on the unused part (see Fig. 2). Instead of imposing the above assumption, we let the model learn the way to sum magnitudes from the data. We use a single FFN, transforming magnitude scalars to 256-dimensional vectors that match the PE dimensionality. After training, we confirmed that most of the projected dimensions are zero, except for one close to 200th dimension (see Fig A.1). We note that the resulting transformation is almost a one-hot encoding. Thus, we showed that the magnitudes do not interfere with the PE information when joining both vectors.

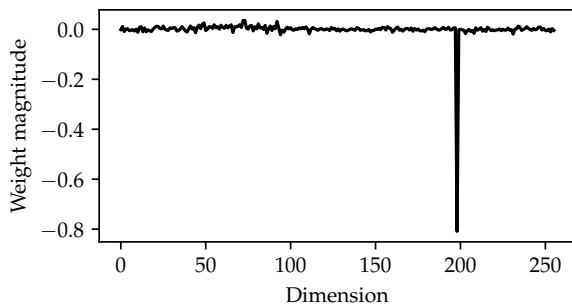


Fig. A.1. Feed-forward network weights that project magnitude scalars to 256-dimensional vectors. The new magnitude vector is then summed to the PE information of equal size.

Appendix B: Looking for a window size

Our preprocessing pipeline consists of windows that move along the light curve, sampling a subset of observations. Using this technique, we deal with the variable length problem while generating more samples for training. Figure B.1 shows the distribution of the lengths of the light curves. Considering that most samples are longer than 200 observations, we defined 200 as the window size.

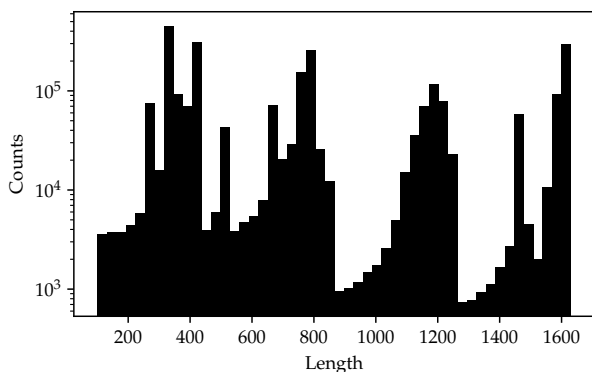


Fig. B.1. Distributions of light curves lengths from the MACHO pre-training dataset.

Appendix C: Fine-tuning validation light curves

Figure C.1 shows the validation learning curves associated with ASTROMER fine-tuned on different training sets. The figure rows are associated with subsets of 20, 50, 100, and 500 samples per class. We also show the learning curves for the fine-tuning using all light curves in the respective catalogs in the last row. The different lines within each subplot show the RMSE of the cross-validation splits.

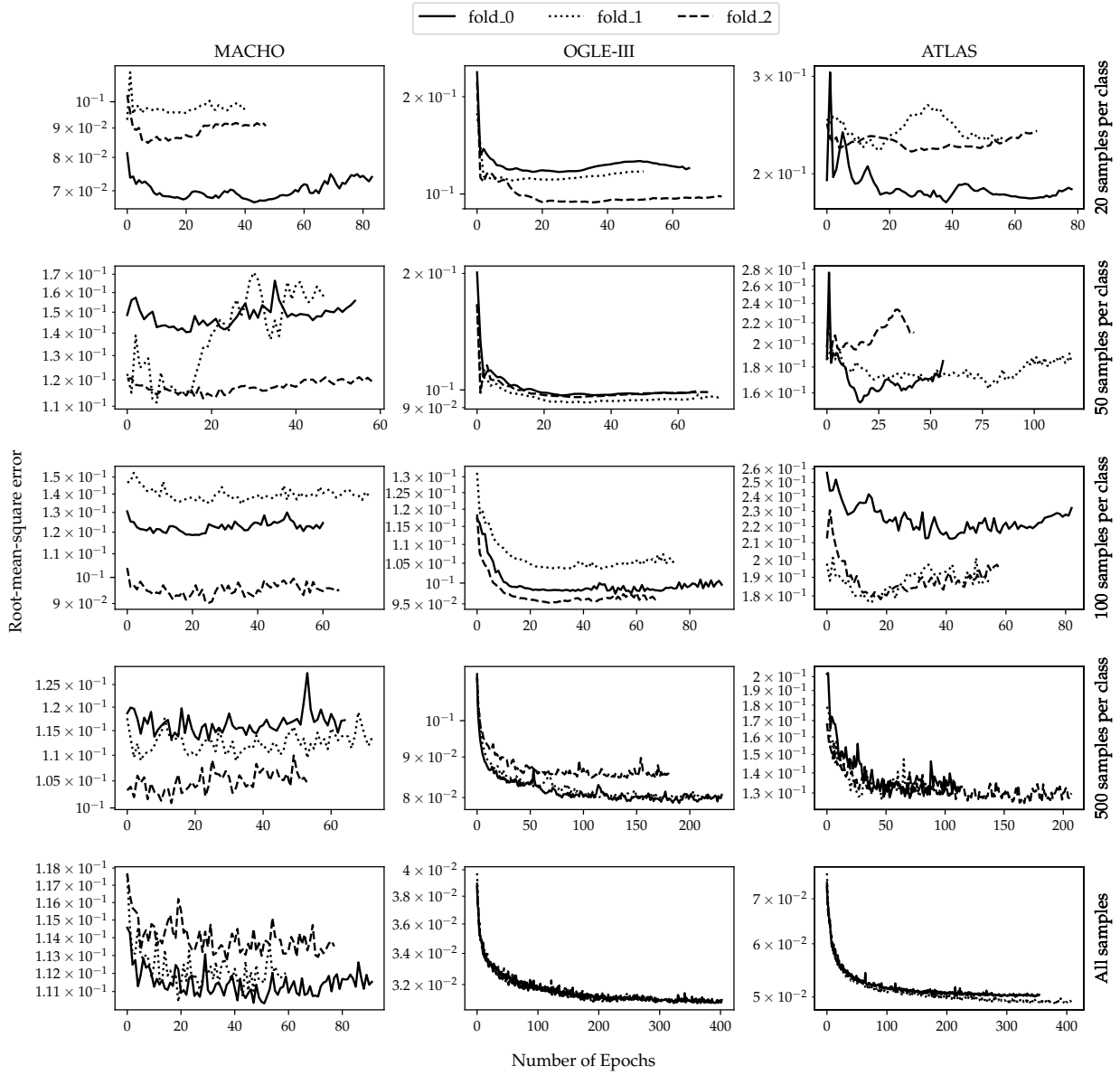


Fig. C.1. Validation learning curves from all datasets during fine-tuning.

Appendix D: Looking for the LSTM optimum learning rate

In order to evaluate the effectiveness of ASTROMER embeddings, we selected a state-of-the-art light curve classifier from the literature (Donoso-Oliva et al. 2021). The classifier consists of two layers of LSTM with 256 units each. The idea is to compare how the model performs when changing the input from light curves to light curve embeddings. We kept hyper-parameters as they were defined in the previous work. However, since we changed the batch size and normalization technique, we looked for the best learning rate on MACHO light curves. Figure D.1 summarizes the cross-validated results showing that a rate of 0.001 performs better in terms of minimal validation loss.

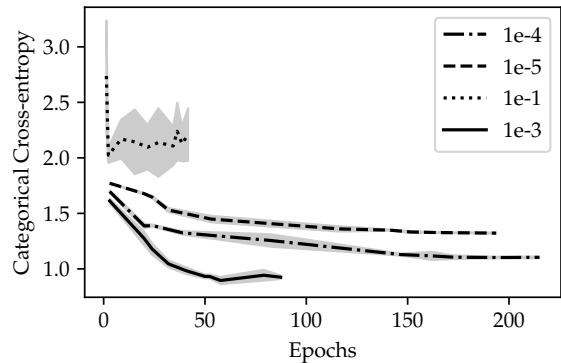


Fig. D.1. Validation learning curves from the baseline classifier (LSTM) using different learning rates.