

Introduction

Deep-sky surveys continue to exponentially increase our discovery rates of variable cosmic phenomena. Researchers are turning towards automated methods of anomaly detection. Here, we provide a convolutional autoencoder based anomaly detection pipeline for

1. *Anomaly detection*
2. *Classifications*

of periodic variable stars (PVSs) detected with Zwicky Transient Facility (ZTF), a pathfinder survey for Large Synoptic Survey Telescope (LSST).

Methodology

We extracted light curves of PVSs from the ZTF catalog of periodic variable stars (ZTF CPVS). The data are given in two filters (g-band and r-band) of magnitudes. We phase-folded the light curves by their joint period and chose to interpolate them using the method of **multivariate Gaussian process (MGPR)**. After that, we generated 160 evenly spaced data points along with the phase direction for both bands. We stacked them horizontally to form an “image” of size 2×160 . They will be encoded through a **convolutional variational autoencoder** to generate latent features. We ran an isolation forest in the latent space to rank each variable by a calculated anomaly score.

In addition, we cross-matched PVSs in the ZTF catalog with those presented in the **SIMBAD** catalog to obtain more reliable and robust class labels. We found 31,541 successfully cross-matched objects. We extracted their latent features together with their class labels, and they were input into a **hierarchical random forest** to train our classifier.

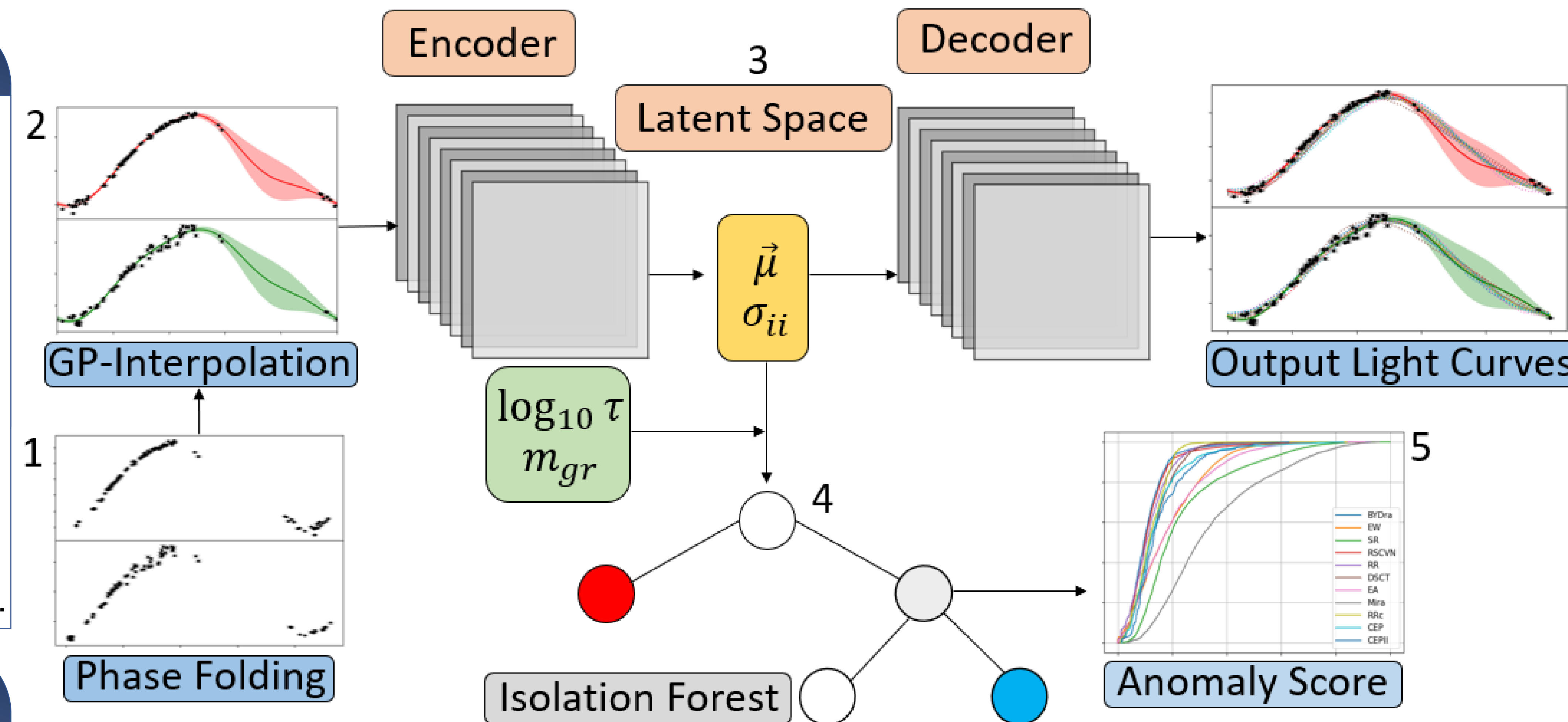


Figure 1. The anomaly detection pipeline: 1. Phase-folding the raw detection data. 2. Interpolation using the Multivariate Gaussian Process. 3. Encoding to generate latent features. 4. Append additional hand-engineered features. 5. Isolation forest and ranking the anomalies.

Top Anomalies

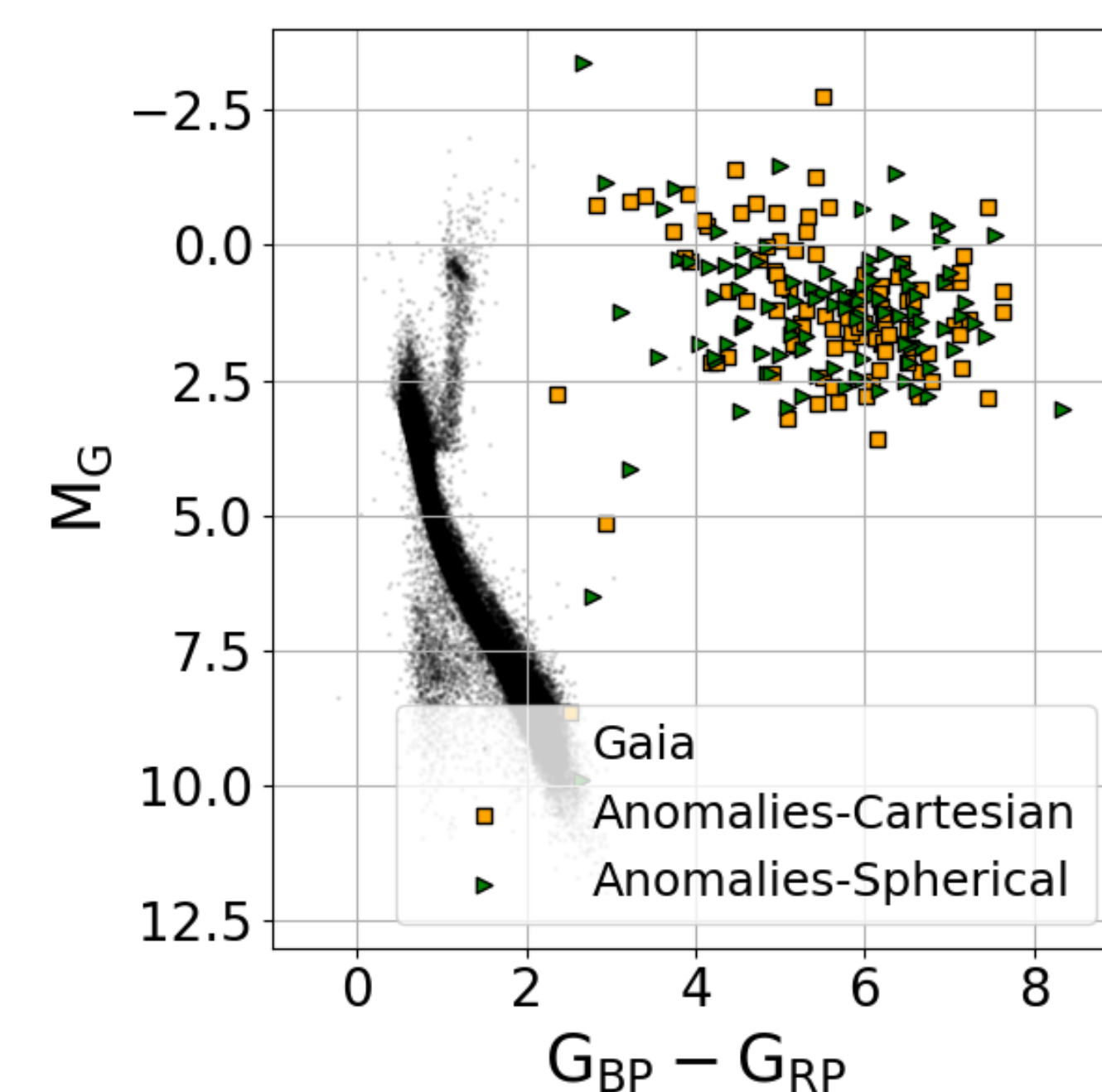


Figure 2. Distribution of the anomalies in the Gaia HR-Diagram against main-sequence stars in black dots

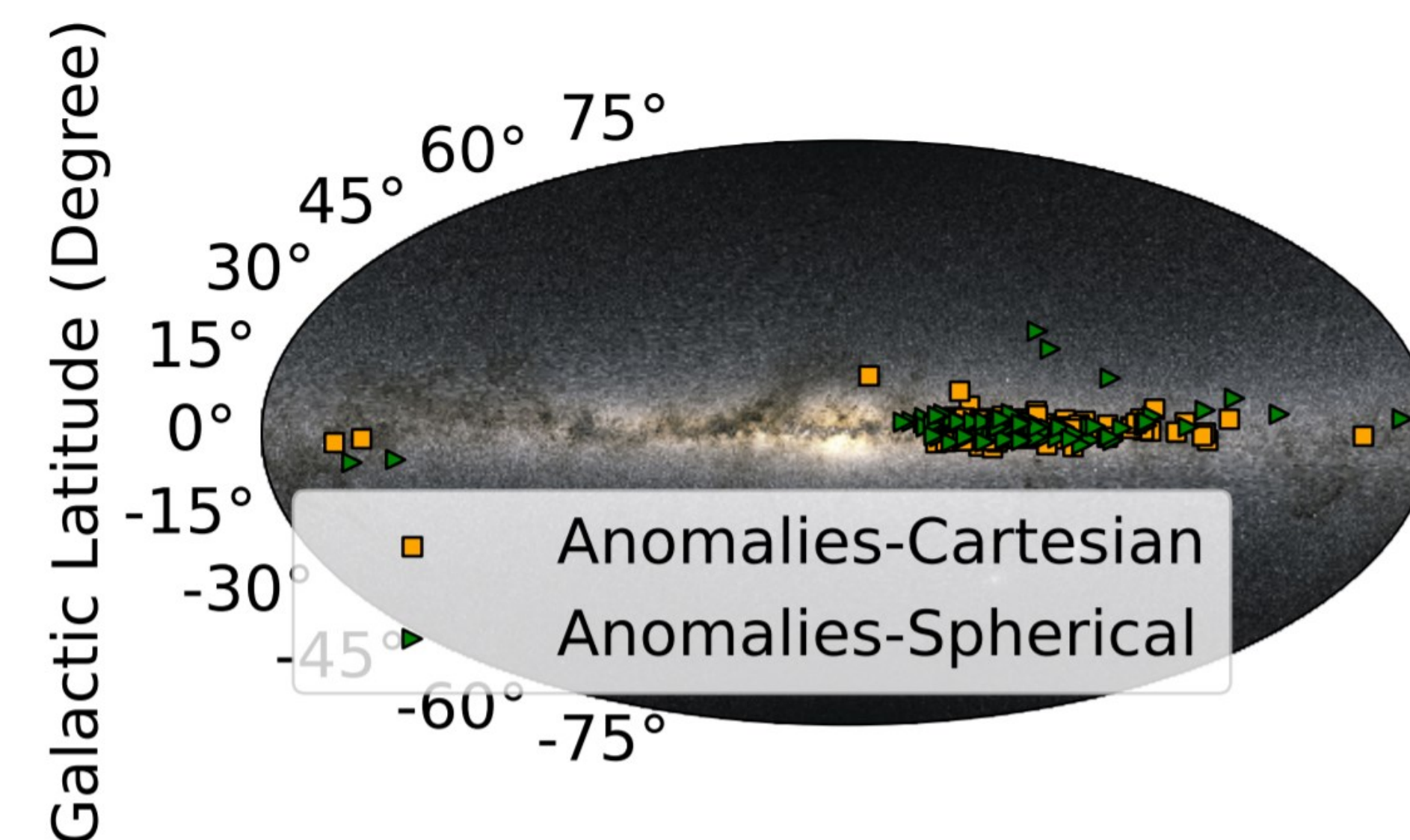


Figure 3. Distribution of the anomalies in the Milky Way galactic coordinate. Galactic longitude is omitted. Image credit: ESA/Gaia/DPAC.

We find that our learned latent space exhibits an annular structure, which inspires us to transform the latent space using **spherical coordinates**. We search for the top 100 anomalies in both latent spaces. We show their distribution in the Gaia HR-Diagram (M_G vs $G_{BP} - G_{RP}$) in Figure 2. We found that they are **bright** and **cool**, which corresponds to evolved stars. Their Galactic positions as shown in Figure 3 indicate that they are tightly located near the Galactic plane. Taken together, these stars are consistent with **young** and **massive Red Giant/AGB** stars.

The Classification Model

True label \ Predicted label	C-Type	HB	RGB	S-Type	Voth	YSOL
C-Type	55.6%	9.6%	1.9%	5.0%	11.9%	16.1%
HB	0.0%	74.1%	0.0%	3.7%	11.1%	11.1%
RGB	6.7%	0.0%	80.0%	0.0%	6.7%	6.7%
S-Type	35.5%	3.2%	6.5%	38.7%	9.7%	6.5%
Voth	11.4%	5.9%	3.0%	3.0%	59.9%	16.9%
YSOL	11.4%	10.3%	1.6%	2.7%	37.5%	36.4%

True label \ Predicted label	AGNL	CEP	EB	LPV	Mira	Pec	Puloth	RR
AGNL	92.6%	0.0%	0.7%	0.7%	0.0%	5.1%	0.0%	0.7%
CEP	0.0%	48.5%	6.1%	1.5%	0.0%	25.8%	1.5%	16.7%
EB	0.1%	0.2%	92.0%	0.4%	0.1%	4.7%	0.4%	2.1%
LPV	1.1%	0.0%	16.7%	51.7%	5.1%	21.2%	0.6%	3.7%
Mira	0.5%	0.0%	19.9%	5.0%	46.3%	24.9%	0.5%	3.0%
Pec	1.2%	1.3%	16.7%	8.8%	2.2%	62.4%	1.2%	6.2%
Puloth	0.0%	2.0%	22.4%	1.0%	1.0%	10.2%	51.0%	12.2%
RR	0.0%	0.0%	4.4%	0.3%	0.2%	2.7%	0.1%	92.3%

Figure 3. Confusion matrix of our new classification model. We show the completeness of each class. The description of class labels are temporarily omitted.

Finally, we highlight the classification results from our hierarchical classifier in Figure 3, where we show the completeness of each class. We find that our new classification model for the ZTF CPVS is reasonably accurate.

Conclusion

We present a convolutional autoencoder-based pipeline for anomaly detection and classification. Techniques similar to those presented here can be used in broader applications to identify anomalies in periodic, multi-variate time series.