
Simulation-Based Pretraining and Domain Adaptation for Astronomical Time Series Tasks with Minimal Labeled Data

Rithwik Gupta^{1 2} Daniel Muthukrishna¹ Jeroen Audenaert¹

Abstract

Astronomical time-series analysis faces a critical limitation: the scarcity of labeled observational data. We present a pre-training approach that leverages physics-informed simulations, significantly reducing the need for labeled examples from real observations. Using classifier-based architectures enhanced with contrastive and adversarial objectives, we create domain-agnostic models that recognize similar astronomical phenomena across different instrumental contexts and learn generalizable representations that transfer effectively to downstream tasks. Our models are trained on simulated astronomical transients from multiple telescope surveys (ZTF and LSST), and demonstrate substantial performance improvements over previous methods in classification, redshift estimation, and anomaly detection tasks when fine-tuned with minimal real data. Remarkably, our models exhibit effective zero-shot transfer capabilities, achieving comparable performance on future telescope (LSST) simulations when trained solely on existing telescope (ZTF) data. Furthermore, they generalize to entirely different astronomical phenomena (namely variable stars from NASA’s *Kepler* telescope) despite being trained on transient events, demonstrating cross-domain capabilities.

1. Introduction

Time-series analysis in astronomy often requires substantial labeled data for supervised learning approaches. Models have been developed to classify variable stars and transient events (e.g. Narayan et al., 2018; Muthukrishna et al., 2019; Rehemtulla et al., 2024), detect anomalies (Perez-Carrasco et al., 2023; Muthukrishna et al., 2022; Villar et al., 2021),

and estimate physical parameters like redshift (e.g. Qu & Sako, 2023; Zhang et al., 2024). However, these models are typically trained for specific instruments or tasks, missing opportunities to leverage commonalities across astronomical datasets. Furthermore, they struggle with new telescopes where labeled data is initially scarce.

Foundation models (FMs) (Bommasani et al., 2022) have transformed natural language processing and computer vision by learning generalizable representations from large quantities of data. However, astronomical data presents unique challenges: (1) limited publicly available labeled data compared to other domains; (2) instrument-specific characteristics that hinder cross-survey generalization; and (3) complex physical phenomena that require domain expertise to model effectively. Recent works have begun to develop foundation models for astrophysical data (e.g. Parker et al., 2024; Donoso-Oliva et al., 2023; Smith et al., 2024; Zhang et al., 2024; Audenaert et al., 2025), but these do not generalize to zero-shot transfer between different time-domain surveys. Moreover, existing FMs typically use self-supervised methods that do not leverage the astronomical class structure that is fundamental to astronomical understanding.

Astronomy has a significant advantage over many domains: decades of physical understanding encoded in simulations. Astronomers have developed detailed models of astrophysical phenomena that generate synthetic light curves (e.g. PLAsTiCC Modelers, 2019). While these simulations do not perfectly match real observations (Gupta & Muthukrishna, 2025), they effectively encode some domain knowledge that can bootstrap learning.

Motivated by the success of supervised classifier-based methods for anomaly detection (Gupta et al., 2024), we propose a novel approach to learning generalizable representations for astronomical time series that:

1. Leverages the latent space of classifiers pretrained on physics-informed simulations
2. Develops domain-agnostic representations through adversarial and contrastive learning objectives
3. Enables effective downstream task performance with minimal labeled real data
4. Facilitates zero-shot transfer to new telescopes

¹Massachusetts Institute of Technology, Cambridge, MA 02139, USA ²Irvington High School, Fremont, CA 94538, USA. Correspondence to: Daniel Muthukrishna <danmuth@mit.edu>.

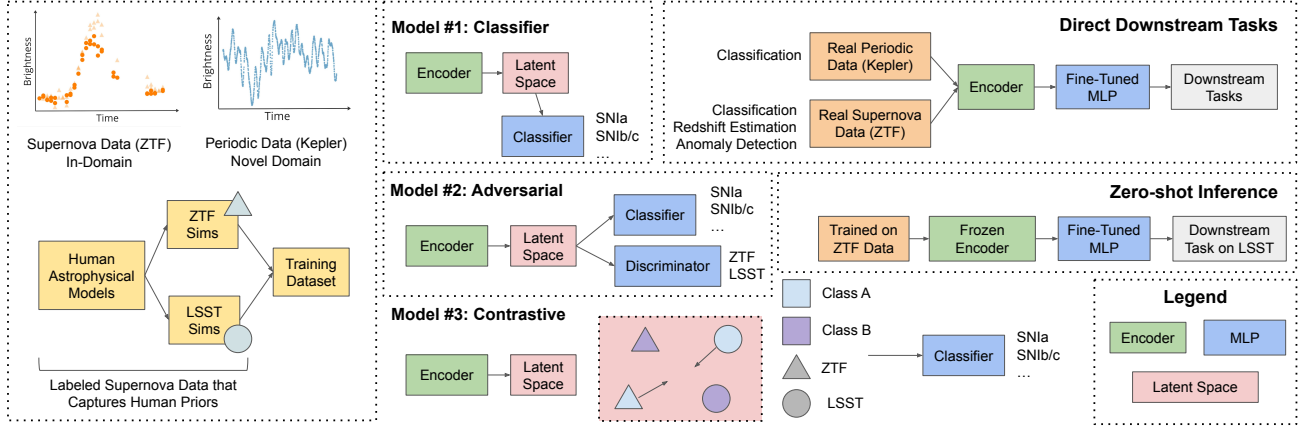


Figure 1. Overview of our simulation-based pre-training methodology. We first pre-train various classifiers and domain-agnostic models using simulated data. We then evaluate these models on various downstream tasks, including zero-shot estimation for new telescopes.

This approach is particularly valuable for upcoming surveys like the Vera C. Rubin Observatory’s Legacy Survey of Space and Time (LSST), which will produce millions of time-series alerts nightly (Ivezić et al., 2019). Having models ready to analyze LSST data from day one—without requiring extensive new labeled datasets—would dramatically accelerate scientific discovery.¹

2. Datasets and Benchmarks

We pretrain our models using 151,468 simulated astronomical transients: 87,080 from the Zwicky Transient Facility (ZTF; Bellm et al., 2018) and 64,388 from LSST, covering eight astronomical transient classes (see Table 3) generated using astrophysical models reflecting each telescope’s observational characteristics (Kessler et al., 2019; Muthukrishna et al., 2019; Narayan & ELAsTiCC Team, 2023). Each object is a multi-channel time-series known as a light curve represented as $[\lambda_p, t_i, f_i, \epsilon_i]$, indicating the passband wavelength, time since first observation, normalized flux, and flux error.

We fine tune and evaluate our pre-trained models on the following datasets and downstream tasks.

ZTF Real Data (In-Domain): 3,747 labeled transients for classification (AUROC), redshift estimation (MSE), and anomaly detection (AUROC, 38 rare objects). **Kepler Real Data (Cross-Domain):** 9,501 variable stars for stellar classification, demonstrating generalization beyond transients. **LSST Simulations (Zero-Shot):** 2,596 objects within ZTF’s redshift range for zero-shot transfer evalu-

ation.

We evaluate under two scenarios: *Limited* (512 labeled objects, simulating early survey deployment) and *Full* (all available labels). This design addresses the practical challenge of having effective models ready for new surveys before extensive expert annotation becomes available.

3. Methods

The first foundation model we propose is a classifier built using simulated data from both ZTF and LSST, and we fine-tune this model for downstream tasks on real data. We further enhance the classifier’s architecture domain-agnostic models by incorporating additional adversarial and contrastive loss components. Fine-tuning a domain-agnostic model on one domain directly translates to performance on other domains. More specific training and model information can be found in Appendix B.

3.1. Classifier

Neural network classifiers have demonstrated the ability to capture the underlying structure of astronomical transients and have been applied to various tasks beyond classification (Etsebeth et al., 2023; Walmsley et al., 2022; Gupta & Muthukrishna, 2025). Motivated by these prior successes, we propose building a foundation model that is a classifier trained on data from both ZTF and LSST. Once the classifier has been trained, we discard the output classification layer, leaving the penultimate layer as the output of our FM. This penultimate layer, henceforth referred to as the latent space, provides meaningful features that exhibit coherent clustering patterns. Although the classifier successfully learns a structured representation of transients from a single telescope (Fig. 4 of Gupta et al., 2024), it does not produce

¹The code used in this work is publicly available: <https://github.com/Rithwik-G/astrofm2.0> and <https://github.com/Rithwik-G/Kepler-FM>

a domain-agnostic latent space. Appendix C qualitatively discusses this issue with visualizations of the latent space (Figures 5 and 6).

3.2. Adversarial Training

To encourage domain-agnostic representations across observatories, we adopt an adversarial training framework, inspired by Generative Adversarial Networks (GANs; Goodfellow et al., 2014). This method promotes the unification of feature representations for similar transients originating from different surveys, such as ZTF and LSST. Specifically, we jointly train a classifier $C(X_i)$, which outputs both class predictions and provides the intermediate latent representations, and a discriminator $D(L_i)$, which attempts to predict the observatory domain. The classifier is trained to both classify the transient correctly and to confuse the discriminator, thereby learning a latent space that is useful for classification yet agnostic to the survey domain. The discriminator, on the other hand, is trained to distinguish between observatories using the latent representations. The adversarial training process is detailed further in Appendix B.2.

3.3. Supervised Contrastive Training

While adversarial training encourages ZTF and LSST transients to be embedded more closely in the latent space, it does not explicitly enforce that objects of the same class should be clustered together. Since we know that class identity is a strong indicator of similarity, we propose using a supervised contrastive loss (Khosla et al., 2020). Rather than relying on data augmentations like traditional contrastive losses, we treat all samples belonging to the same class as positive pairs and apply the contrastive objective accordingly. This modification provides a more explicit signal for class-based alignment and naturally encourages unification across domains, as long as examples from the same class are drawn from both surveys. We use the contrastive loss proposed in Chen et al. (2020) for model pretraining. More details can be found in Appendix B.3.

Our adversarial and contrastive models are trained to learn an explicit relationship between ZTF and LSST. This is useful because once these models are fine-tuned on for a task on one domain, they can be directly applied to another domain. However, for downstream tasks on a single survey, classifiers perform as well as these domain-agnostic models.

3.4. Downstream Tasks

Fine Tuning Foundation Models To fine-tune our classifier-based foundation models for regression and classification tasks, we attach a multi-layer perceptron (MLP) to the penultimate neural network layers of our models. Normally, we would then freeze the initial model and fine-tune

just the MLP, but this hyperparameter setup does not work effectively for all tasks. Further details on how we fine-tune can be found in Appendix B.4 and in the released code.

Baselines and Our Models When evaluating our models for downstream tasks, we compare them with models trained directly on a downstream task (No Pretraining). This is largely how models for real data have been trained in past research. We propose three foundation models for evaluation: a Classifier, an Adversarial model, and a Contrastive model, all trained on both ZTF and LSST data. We evaluate these three methods of pretraining on various downstream and zero-shot tasks. While our FMs are trained on ZTF and LSST, we also evaluate pretraining Classifiers on telescopes individually. In Table 1, we clearly disambiguate which dataset the model used for pretraining. In Table 2, Classifier refers to our proposed foundation model trained on both ZTF and LSST data.

4. Results

We evaluate our foundation models on multiple downstream tasks using both real observational data and simulated data from different telescopes. Our experiments demonstrate that pretraining on physics-informed simulations provides substantial improvements over training from scratch, with particularly strong results for cross-survey generalization. Further analysis can be found in Appendix D.

4.1. Downstream Performance on Real Astronomical Data

As seen in Table 1, our classifier-based foundation models outperform previous baseline methods trained directly on numerous tasks. Overall, these models achieve better performance than state-of-the-art (SoTA) no pretraining methods for astronomical tasks on real data. We also see performance improvements for tasks from the Limited to Full testing scenarios, which makes sense as more data for fine-tuning should result in better performance. This also indicates that the simulations used for pretraining do not perfectly reflect real data, which we further discuss in Appendix D.2.

Surprisingly, our models also improve performance on tasks on Kepler Data, even though Kepler data is not given to the model during pretraining. Kepler specifically looks for periodic transient events, in comparison to the supernova transients from ZTF and LSST. We find this to be similar to how LLMs perform well on tasks they were not trained on.

4.2. Cross-Survey Generalization: Zero-Shot Inference for LSST

Our contrastive and adversarially trained models are designed to be domain-agnostic and understand relationships

Model (Pretraining Data)	Classification		ZTF Real Data Redshift ($\times 10^2$)				AD		Kepler Real Data Classification Full	Simulations (Redshifting)	
	Limited	Full	Limited	Full	Limited	Full	Limited	Full		ZTF Full	LSST Full
Previous Work	0.637	0.853	0.602	0.385	0.498	0.527			0.901	0.079	0.289
No Pretraining	± 0.005	± 0.013	± 0.005	± 0.031	± 0.013	± 0.027			± 0.003	± 0.022	± 0.018
Classifier	0.875	0.904	0.491	0.387	0.605	0.596				0.028	0.252
ZTF	± 0.020	± 0.012	± 0.010	± 0.036	± 0.025	± 0.008				± 0.003	± 0.004
Classifier	0.879	0.910	0.479	0.382	0.622	0.616			0.968	0.026	0.177
ZTF and LSST	± 0.011	± 0.013	± 0.006	± 0.039	± 0.018	± 0.036			± 0.006	± 0.002	± 0.008
Contrastive	0.886	0.914	0.487	0.373	0.584	0.576			0.946	0.028	0.191
	± 0.026	± 0.005	± 0.003	± 0.017	± 0.018	± 0.028			± 0.021	± 0.002	± 0.013
Adversarial	0.844	0.853	0.520	0.419	0.559	0.546			0.925	0.030	0.197
	± 0.016	± 0.012	± 0.004	± 0.042	± 0.024	± 0.077			± 0.015	± 0.003	± 0.014
Performance Metric	AUROC	AUROC	MSE	MSE	AUROC	AUROC			AUROC	MSE	MSE

Table 1. Foundation model performance on various tasks. We train five different foundation models and fine-tune each of them five times. We report the mean and standard deviation of these recorded results. The final three model rows are the main contributions of this work and the first row is the current baseline and often the SoTA. The best performance is bolded for each task.

	Pretraining	Redshifting Data		LSST MSE
		ZTF	LSST	
Direct	No Pretraining	No	Yes	0.0750 \pm 0.0252
	No Pretraining	Yes	Yes	0.0614 \pm 0.0033
	Classifier	Yes	Yes	0.0579 \pm 0.0061
Zero-shot	No Pretraining	Yes	No	0.1869 \pm 0.0127
	Classifier	Yes	No	0.1035 \pm 0.0081
	Contrastive	Yes	No	0.0727 \pm 0.0056
	Adversarial	Yes	No	0.0744 \pm 0.0063
	Contrastive kNN	Yes	No	0.0854 \pm 0.0040

Table 2. Performance of various models for LSST redshift estimation. Performance is reported as the mean and standard deviation of training five different FMs and five iterations of fine-tuning each of them. Our zero-shot methods achieve similar performance to previous methods directly trained on redshifting LSST. We find no benefit to training with both the contrastive and adversarial objective.

between different telescopes in their latent space. To leverage this learned relationship, we first fine-tune our FM on a downstream task for ZTF and then evaluate this model’s zero-shot performance on LSST. Importantly, we freeze the entire FM to preserve the learned relationship between the surveys during pretraining. If our FM is indeed encoding a unified latent space, the zero-shot performance on LSST should improve as we train on ZTF. When evaluating LSST redshifting in the zero-shot setting, we restrict the evaluation to LSST light-curves in the same redshift range as ZTF. This decision is discussed further in Appendix A.1.

As seen in Table 2, domain-agnostic FMs fine-tuned only ZTF redshifting data (simulated time-series with corresponding redshift measurements) work exceptionally well when repurposed to redshift simulated LSST transients. We compare these zero-shot methods to previous baseline methods that do not involve pretraining and our proposed method involving pretraining (similar to the fine-tuning done for Table 1). Models trained on both ZTF and LSST redshifting data (Rows 2 and 3 of Table 2) are first fine-tuned to redshift ZTF and then to redshift LSST. We describe the k Nearest Neighbors (kNN) zero-shot estimation method in Appendix E. Overall, our domain-agnostic models achieve

the performance of baseline methods trained directly on LSST without any LSST data and significantly improve the performance of previous zero-shot methods.

Table 2 further reiterates that pretraining on adjacent domains and tasks produces SoTA models. The best model for this redshift estimation task outside the zero-shot scenario is a pretrained classifier fine-tuned on both ZTF and LSST redshifting data, essentially incorporating two tasks across two domains.

5. Conclusion

The shift from manual discovery to data driven discovery has motivated the development for machine learning in many scientific domains. Effective foundation models can expedite this process, and to build such models for astronomy, we propose leveraging existing physics-informed simulations. Training specialized classifiers on human-generated simulated data proves to be an effective way to incorporate domain expertise into these models. Fine-tuning our models for tasks on real data achieves SoTA performance on numerous downstream tasks and has excellent zero-shot task performance. We believe that the development of foundation models for astronomy is the next major step in expediting discovery and we hope that this work facilitates future research in the development of FMs for science through supervised training.

We see numerous promising research directions for future work. Incorporating unlabeled data in model fine-tuning could yield better results by better exposing models to the structure of real data after being pretrained. Further, different methods of supervision could help models extract more meaningful information from physics-informed simulations. In conclusion, our work aims to bridge the gap between past research for machine learning for astronomy, with the current era of discovery necessitating the development of models that leverage all that we know for novel tasks.

References

- Audenaert, J. and Muthukrishna, D., Gregory, P., Hogg, D., and Villar, V. A. Causal Foundation Models: Disentangling Physics from Instrument Properties. *ICML 2025 Workshop on Foundation Models for Structured Data*, 2025.
- Audenaert, J., Kuszlewicz, J. S., Handberg, R., Tkachenko, A., Armstrong, D. J., Hon, M., Kgoadi, R., Lund, M. N., Bell, K. J., Bugnet, L., Bowman, D. M., Johnston, C., García, R. A., Stello, D., Molnár, L., Plachy, E., Buzasi, D., and Aerts, C. Tess data for astero-seismology (t'da) stellar variability classification pipeline: Setup and application to the kepler q9 data. *The Astronomical Journal*, 162(5):209, October 2021. ISSN 1538-3881. doi: 10.3847/1538-3881/ac166a. URL <http://dx.doi.org/10.3847/1538-3881/ac166a>.
- Bellm, E. C., Kulkarni, S. R., Graham, M. J., Dekany, R., Smith, R. M., Riddle, R., Masci, F. J., Helou, G., Prince, T. A., Adams, S. M., Barbarino, C., Barlow, T., Bauer, J., Beck, R., Belicki, J., Biswas, R., Blagorodnova, N., Bodewits, D., Bolin, B., Brinnel, V., Brooke, T., Bue, B., Bulla, M., Burruss, R., Cenko, S. B., Chang, C.-K., Connolly, A., Coughlin, M., Cromer, J., Cunningham, V., De, K., Delacroix, A., Desai, V., Duev, D. A., Eadie, G., Farnham, T. L., Feeney, M., Feindt, U., Flynn, D., Franckowiak, A., Frederick, S., Fremling, C., Gal-Yam, A., Gezari, S., Giomi, G., Goldstein, D. A., Golkhou, V. Z., Goobar, A., Groom, S., Hacopians, E., Hale, D., Henning, J., Ho, A. Y. Q., Hover, D., Howell, J., Hung, T., Huppenkothen, D., Imel, D., Ip, W.-H., Ivezić, Ž., Jackson, E., Jones, L., Juric, M., Kasliwal, M. M., Kaspi, S., Kaye, S., Kelley, M. S. P., Kowalski, M., Kramer, E., Kupfer, T., Landry, W., Laher, R. R., Lee, C.-D., Lin, H. W., Lin, Z.-Y., Lunnan, R., Giomi, M., Mahabal, A., Mao, P., Miller, A. A., Monkewitz, S., Murphy, P., Ngeow, C.-C., Nordin, J., Nugent, P., Ofek, E., Patterson, M. T., Penprase, B., Porter, M., Rauch, L., Rebbapragada, U., Reiley, D., Rigault, M., Rodriguez, H., van Roestel, J., Rusholme, B., van Santen, J., Schulze, S., Shupe, D. L., Singer, L. P., Soumagnac, M. T., Stein, R., Surace, J., Sollerman, J., Szkody, P., Taddia, F., Terek, S., Sistine, A. V., van Velzen, S., Vestrand, W. T., Walters, R., Ward, C., Ye, Q.-Z., Yu, P.-C., Yan, L., and Zolkower, J. The zwicky transient facility: System overview, performance, and first results. *Publications of the Astronomical Society of the Pacific*, 131(995):018002, dec 2018. doi: 10.1088/1538-3873/aaecbe. URL <https://doi.org/10.1088/1538-3873/aaecbe>.
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N., Chen, A., Creel, K., Davis, J. Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K., Goodman, N., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D. E., Hong, J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P. W., Krass, M., Krishna, R., Kuditipudi, R., Kumar, A., Ladhak, F., Lee, M., Lee, T., Leskovec, J., Levent, I., Li, X. L., Li, X., Ma, T., Malik, A., Manning, C. D., Mirchandani, S., Mitchell, E., Munyikwa, Z., Nair, S., Narayan, A., Narayanan, D., Newman, B., Nie, A., Niebles, J. C., Nilforoshan, H., Nyarko, J., Ogut, G., Orr, L., Papadimitriou, I., Park, J. S., Piech, C., Portelance, E., Potts, C., Raghunathan, A., Reich, R., Ren, H., Rong, F., Roohani, Y., Ruiz, C., Ryan, J., Ré, C., Sadigh, D., Sagawa, S., Santhanam, K., Shih, A., Srinivasan, K., Tamkin, A., Taori, R., Thomas, A. W., Tramèr, F., Wang, R. E., Wang, W., Wu, B., Wu, J., Wu, Y., Xie, S. M., Yasunaga, M., You, J., Zaharia, M., Zhang, M., Zhang, T., Zhang, X., Zhang, Y., Zheng, L., Zhou, K., and Liang, P. On the opportunities and risks of foundation models, 2022. URL <https://arxiv.org/abs/2108.07258>.
- Boone, K. Avocado: Photometric classification of astronomical transients with gaussian process augmentation. *The Astronomical Journal*, 158(6):257, dec 2019. doi: 10.3847/1538-3881/ab5182. URL <https://doi.org/10.3847/1538-3881/ab5182>.
- Borucki, W. J., Koch, D., Basri, G., Batalha, N., Brown, T., Caldwell, D., Caldwell, J., Christensen-Dalsgaard, J., Cochran, W. D., DeVore, E., Dunham, E. W., Dupree, A. K., Gautier, T. N., Geary, J. C., Gilliland, R., Gould, A., Howell, S. B., Jenkins, J. M., Kondo, Y., Latham, D. W., Marcy, G. W., Meibom, S., Kjeldsen, H., Lissauer, J. J., Monet, D. G., Morrison, D., Sasselov, D., Tarter, J., Boss, A., Brownlee, D., Owen, T., Buzasi, D., Charbonneau, D., Doyle, L., Fortney, J., Ford, E. B., Holman, M. J., Seager, S., Steffen, J. H., Welsh, W. F., Rowe, J., Anderson, H., Buchhave, L., Ciardi, D., Walkowicz, L., Sherry, W., Horch, E., Isaacson, H., Everett, M. E., Fischer, D., Torres, G., Johnson, J. A., Endl, M., MacQueen, P., Bryson, S. T., Dotson, J., Haas, M., Kolodziejczak, J., Van Cleve, J., Chandrasekaran, H., Twicken, J. D., Quintana, E. V., Clarke, B. D., Allen, C., Li, J., Wu, H., Tenenbaum, P., Verner, E., Bruhweiler, F., Barnes, J., and Prsa, A. Kepler Planet-Detection Mission: Introduction and First Results. *Science*, 327(5968):977, February 2010. doi: 10.1126/science.1185402.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations, 2020. URL <https://arxiv.org/abs/2002.05709>.

- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1724–1734. Association for Computational Linguistics, 2014. doi: 10.3115/v1/D14-1179.
- Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NIPS 2014 Workshop on Deep Learning, December 2014*, 2014.
- Donoso-Oliva, C., Becker, I., Protopapas, P., Cabrera-Vives, G., Vishnu, M., and Vardhan, H. Astromer: A transformer-based embedding for the representation of light curves. *Astronomy and Astrophysics*, 670: A54, February 2023. ISSN 1432-0746. doi: 10.1051/0004-6361/202243928. URL <http://dx.doi.org/10.1051/0004-6361/202243928>.
- Etsebeth, V., Lochner, M., Walmsley, M., and Grespan, M. Astronomy at scale: Searching for anomalies amongst 4 million galaxies, 2023.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial networks, 2014. URL <https://arxiv.org/abs/1406.2661>.
- Gupta, R. and Muthukrishna, D. Transfer learning for transient classification: From simulations to real data and ztf to lsst, 2025. URL <https://arxiv.org/abs/2502.18558>.
- Gupta, R., Muthukrishna, D., and Lochner, M. A classifier-based approach to multi-class anomaly detection for astronomical transients, 2024. URL <https://arxiv.org/abs/2403.14742>.
- Huang, H., Muthukrishna, D., Nair, P., Zhang, Z., Fausnaugh, M., Majumder, T., Foley, R. J., and Ricker, G. R. Predicting the age of astronomical transients from real-time multivariate time series. *arXiv preprint arXiv:2311.17143*, 2023.
- Ivezić, Ž., Kahn, S. M., Tyson, J. A., Abel, B., Acosta, E., Allsman, R., Alonso, D., AlSayyad, Y., Anderson, S. F., Andrew, J., and et al. LSST: From Science Drivers to Reference Design and Anticipated Data Products. *The Astrophysical Journal*, 873:111, March 2019. doi: 10.3847/1538-4357/ab042c.
- Kessler, R., Narayan, G., Avelino, A., Bachelet, E., Biswas, R., Brown, P. J., Chernoff, D. F., Connolly, A. J., Dai, M., Daniel, S., Stefano, R. D., Drout, M. R., Galbany, L., González-Gaitán, S., Graham, M. L., Hložek, R., Ishida, E. E. O., Guillochon, J., Jha, S. W., Jones, D. O., Mandel, K. S., Muthukrishna, D., O’Grady, A., Peters, C. M., Pierel, J. R., Ponder, K. A., Prša, A., Rodney, S., and and, V. A. V. Models and simulations for the photometric LSST astronomical time series classification challenge (PLAsTiCC). *Publications of the Astronomical Society of the Pacific*, 131(1003):094501, jul 2019. doi: 10.1088/1538-3873/ab26f1. URL <https://doi.org/10.1088%2F1538-3873%2Fab26f1>.
- Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., and Krishnan, D. Supervised contrastive learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 18661–18673. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/d89a66c7c80a29b1bdbab0f2a1a94af8-Paper.pdf.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization, 2017.
- Koch, D. G., Borucki, W. J., Basri, G., Batalha, N. M., Brown, T. M., Caldwell, D., Christensen-Dalsgaard, J., Cochran, W. D., DeVore, E., Dunham, E. W., Gautier, Thomas N., I., Geary, J. C., Gilliland, R. L., Gould, A., Jenkins, J., Kondo, Y., Latham, D. W., Lissauer, J. J., Marcy, G., Monet, D., Sasselov, D., Boss, A., Brownlee, D., Caldwell, J., Dupree, A. K., Howell, S. B., Kjeldsen, H., Meibom, S., Morrison, D., Owen, T., Reitsema, H., Tarter, J., Bryson, S. T., Dotson, J. L., Gazis, P., Haas, M. R., Kolodziejczak, J., Rowe, J. F., Van Cleve, J. E., Allen, C., Chandrasekaran, H., Clarke, B. D., Li, J., Quintana, E. V., Tenenbaum, P., Twicken, J. D., and Wu, H. Kepler Mission Design, Realized Photometric Performance, and Early Science. *The Astrophysical Journal Letters*, 713(2):L79–L86, April 2010. doi: 10.1088/2041-8205/713/2/L79.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- Liu, F. T., Ting, K. M., and Zhou, Z.-H. Isolation forest. In *2008 eighth ieee international conference on data mining*, pp. 413–422. IEEE, 2008.
- McInnes, L., Healy, J., and Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction, 2020.
- Muthukrishna, D., Narayan, G., Mandel, K. S., Biswas, R., and Hložek, R. RAPID: Early classification of explosive

- transients using deep learning. *Publications of the Astronomical Society of the Pacific*, 131(1005):118002, sep 2019. doi: 10.1088/1538-3873/ab1609. URL <https://doi.org/10.1088%2F1538-3873%2Fab1609>.
- Muthukrishna, D., Mandel, K. S., Lochner, M., Webb, S., and Narayan, G. Real-time detection of anomalies in large-scale transient surveys. *Monthly Notices of the Royal Astronomical Society*, 517(1):393–419, November 2022. doi: 10.1093/mnras/stac2582.
- Muthukrishna, D., Mandel, K. S., Lochner, M., Webb, S., and Narayan, G. Real-time detection of anomalies in large-scale transient surveys. *Monthly Notices of the Royal Astronomical Society*, 517(1):393–419, sep 2022. doi: 10.1093/mnras/stac2582. URL <https://doi.org/10.1093%2Fmnras%2Fstac2582>.
- Narayan, G. and ELAsTiCC Team. The Extended LSST Astronomical Time-series Classification Challenge (ELAsTiCC). In *American Astronomical Society Meeting Abstracts*, volume 241 of *American Astronomical Society Meeting Abstracts*, pp. 117.01, January 2023.
- Narayan, G., Zaidi, T., Soraisam, M. D., Wang, Z., Lochner, M., Matheson, T., Saha, A., Yang, S., Zhao, Z., Kececioglu, J., Scheidegger, C., Snodgrass, R. T., Axelrod, T., Jenness, T., Maier, R. S., Ridgway, S. T., Seaman, R. L., Evans, E. M., Singh, N., Taylor, C., Toeniskoetter, J., Welch, E., Zhu, S., and ANTARES Collaboration. Machine-learning-based Brokers for Real-time Classification of the LSST Alert Stream. *The Astrophysical Journal*, 236:9, May 2018. doi: 10.3847/1538-4365/aab781.
- Parker, L., Lanusse, F., Golkar, S., Sarra, L., Cranmer, M., Bietti, A., Eickenberg, M., Krawezik, G., McCabe, M., Morel, R., Ohana, R., Pettee, M., Régalo-Saint Blancard, B., Cho, K., and Ho, S. Astroclip: a cross-modal foundation model for galaxies. *Monthly Notices of the Royal Astronomical Society*, 531(4):4990–5011, June 2024. ISSN 1365-2966. doi: 10.1093/mnras/stae1450. URL <http://dx.doi.org/10.1093/mnras/stae1450>.
- Perez-Carrasco, M., Cabrera-Vives, G., Hernandez-García, L., Förster, F., Sanchez-Saez, P., Arancibia, A. M. M., Arredondo, J., Astorga, N., Bauer, F. E., Bayo, A., Cateilan, M., Dastidar, R., Estévez, P. A., Lira, P., and Pignata, G. Alert classification for the alerce broker system: The anomaly detector. *The Astronomical Journal*, 166(4):151, sep 2023. doi: 10.3847/1538-3881/ace0c1. URL <https://dx.doi.org/10.3847/1538-3881/ace0c1>.
- PLAsTiCC Modelers. Libraries & Recommended Citations for using PLAsTiCC Models, March 2019. URL <https://doi.org/10.5281/zenodo.2612896>.
- Qu, H. and Sako, M. Photo-zsnthesis: Converting type ia supernova lightcurves to redshift estimates via deep learning, 2023. URL <https://arxiv.org/abs/2305.11869>.
- Rehemitulla, N., Miller, A. A., Laz, T. J. D., Coughlin, M. W., Fremling, C., Perley, D. A., Qin, Y.-J., Sollerman, J., Mahabal, A. A., Laher, R. R., Riddle, R., Rusholme, B., and Kulkarni, S. R. The zwicky transient facility bright transient survey. iii. btsbot: Automated identification and follow-up of bright transients with deep learning. *The Astrophysical Journal*, 972(1):7, aug 2024. doi: 10.3847/1538-4357/ad5666. URL <https://dx.doi.org/10.3847/1538-4357/ad5666>.
- Smith, M. J., Roberts, R. J., Angeloudi, E., and Huertas-Company, M. Astropt: Scaling large observation models for astronomy, 2024. URL <https://arxiv.org/abs/2405.14930>.
- Villar, V. A., Cranmer, M., Berger, E., Contardo, G., Ho, S., Hosseinzadeh, G., and Lin, J. Y.-Y. A deep-learning approach for live anomaly detection of extragalactic transients. *The Astrophysical Journal Supplement Series*, 255(2):24, 2021.
- Walmsley, M., Scaife, A. M. M., Lintott, C., Lochner, M., Etsebeth, V., Géron, T., Dickinson, H., Fortson, L., Kruk, S., Masters, K. L., Mantha, K. B., and Simmons, B. D. Practical galaxy morphology tools from deep supervised representation learning. *Monthly Notices of the Royal Astronomical Society*, 513(2):1581–1599, 02 2022. ISSN 0035-8711. doi: 10.1093/mnras/stac525. URL <https://doi.org/10.1093/mnras/stac525>.
- Zhang, G., Helfer, T., Gagliano, A. T., Mishra-Sharma, S., and Villar, V. A. Maven: A multimodal foundation model for supernova science, 2024. URL <https://arxiv.org/abs/2408.16829>.

A. Task Descriptions

In this section, we describe each downstream task in detail.

1. Training Data (Supernova data)
 - (a) Classification (Sims): **Simulated time series** with corresponding class labels from ZTF (Muthukrishna et al., 2019; PLAsTiCC Modelers, 2019) and LSST (Narayan & ELAsTiCC Team, 2023). This task is only used for training our FMs.
2. Supernova Data Downstream Tasks (in domain)
 - (a) Classification (Real Data): **Real time series from ZTF** data with corresponding class labels (Rehemtulla et al., 2024). This task is evaluated using the macro-averaged Area under the ROC Curve (AUROC).
 - (b) Anomaly Detection: **Real time series from ZTF** with anomalous objects labeled (Rehemtulla et al., 2024). This task is also evaluated using the macro-averaged AUROC metric, treating anomaly detection (AD) as a binary classification task. Models are not provided any anomalous data during training to emulate real-world AD.
 - (c) Redshift Estimation (Real Data): **Real time series from ZTF** with corresponding spectroscopic redshifts (Rehemtulla et al., 2024). This task is evaluated using the mean squared error (MSE).
 - (d) Redshift Estimation (Sims): Simulated time series from ZTF (Muthukrishna et al., 2019; PLAsTiCC Modelers, 2019) and LSST (Narayan & ELAsTiCC Team, 2023) with corresponding spectroscopic redshifts. This task is evaluated similarly to the last one using the MSE. We note that this task does not share any data with the training task that also uses simulated data.
 - (e) Zero-shot Redshift Estimation (Sims): Special testing scenario of the above task. The training data includes ZTF and evaluation is performed **zero-shot on LSST data**. This task is evaluated using the MSE.
3. Periodic Data Downstream Tasks (out of domain)
 - (a) Periodic Classification (e.g. Audenaert et al., 2021): **Real time series from Kepler** (Audenaert et al., 2021) with corresponding class labels. This task is evaluated using the macro-averaged AUROC metric.

For most benchmarks on real data, we propose two scenarios. The `Limited` scenario limits the amount of data available for fine-tuning to simulate the development process for new telescopes where labeled data will not be readily available.

In this scenario, the rest of the data is provided without labels. We choose to limit this scenario to 512 labeled objects for each task which is the amount of data that can be labeled in a few months on new telescopes. Due to the high variability in selecting such a small number of samples, we evaluate over 5 randomly selected `Limited` data samples. Incorporating unlabeled data into the model fine-tuning process is beyond the scope of this work, but we hope that future researchers explore ways to accomplish this. The `Full` mode includes all data. The zero-shot redshift and Kepler benchmarks do not include the `Limited` testing scenario, the latter is an out-of-domain benchmark that is not a direct target for our FMs.

A.1. Simulated Tasks

Table 3 shows the number of objects per supernova class in our datasets. The supernova classes are described further in previous work (Gupta et al., 2024; Muthukrishna et al., 2022). The classification tasks are used for pretraining and redshift task is used for downstream zero-shot evaluation. Because classification is used for pretraining, we ensure that there is no overlap between the data for both tasks. The LSST zero-shot redshift estimation task restricts evaluation to light curves that are in the ZTF redshift range. We do this because, as seen in Fig. 2, LSST will observe a far greater range of redshifts. The goal of this work is to leverage existing data for future models, and the current prior only exists for data from ZTF. Thus, we find it reasonable to perform zero-shot evaluation solely on LSST. Evaluating in this manner also seems to limit the applicability of these zero-shot models, however astronomers may know the rough redshift range of an object and thus can choose to use these models when they see fit.

A.2. ZTF Real Data

Table 4 shows the number of light curves from each class in our dataset of real ZTF light curves. We define anomalies to be transients from any of the following classes: TDE, Ca-rich, ILRT, LBV, LRN, SLSN-I, SLSN-II, SN Ia-91T, SN Ia-91bg, SN Ibn, SN Ic-BL, SN Icn, and TDE. These objects are specifically chosen because of their low observation rates and limited human understanding.

Fig. 3 shows the redshift distribution in our dataset of real ZTF objects. This redshift range is similar to that of the simulated ZTF data as seen in 2 [top]. Unlike with the simulated redshift estimation task, we do not use any real light curves for pretraining and reserve them solely for evaluation. Thus, all tasks on ZTF real data share a roughly identical pool of observed objects. Further information about the real ZTF data used in this work can be found in (Rehemtulla et al., 2024).

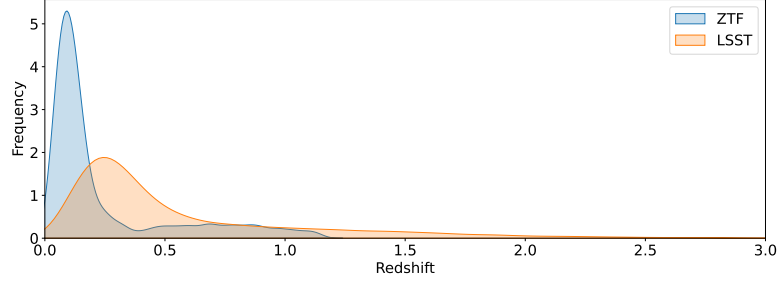


Figure 2. The range of redshifts for objects in our datasets. After fine-tuning a domain-agnostic FM to estimate redshifts for ZTF objects, we evaluate its zero-shot LSST performance only using transients in the ZTF range.

Task	SNIa	SNIa-91bg	SNIax	SNIb/c	SNI	TDE	SLSN-I	AGN	Total
ZTF Classification	9436	10663	10681	6769	31193	9260	10451	8627	87080
LSST Classification	8427	6079	8298	7664	9465	9686	6947	7822	64388
ZTF Redshifting	967	1140	1124	702	3213	932	1053	869	10000
LSST Redshifting	1236	1238	1226	1286	1248	1279	1253	1234	9990
Zero-shot LSST	265	273	535	258	333	215	446	271	2596

Table 3. The amount of simulated labeled data pairs available for the different tasks on simulated data. There is no overlap between the data in the classification tasks and the redshifting tasks. The zero-shot scenario for LSST only includes transients from LSST which are in the redshift range of ZTF.

Task	SNIa	SNIb/c	SNI	Anomaly	Total
Classification	771 (107)	2828 (350)	148 (12)	0	3747
Redshift Estimation	771 (107)	2828 (350)	148 (12)	0	3747
Anomaly Detection	771 (107)	2828 (350)	148 (12)	0 (38)	3785

Table 4. The amount of real labeled data pairs available for the different LSST and ZTF tasks used in this work. The Limited versions of these tasks use a random sample of 512 objects. The number in parenthesis represents the amount of data provided in the evaluation set.

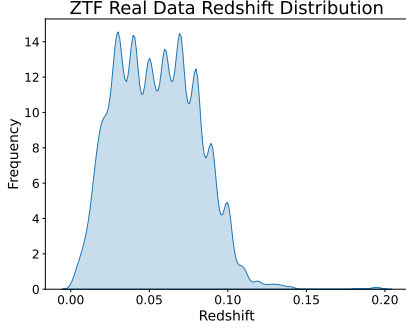


Figure 3. Redshift distribution for our dataset of real ZTF objects

A.3. Zero-Shot LSST Redshifting

The motivation behind zero-shot prediction is to directly reuse our understanding of ZTF to train models for LSST. Thus, when we evaluate zero-shot FM performance, we restrict the evaluation to LSST light-curves in the same redshift range as ZTF. LSST will observe a much broader redshift range than ZTF (as seen in Fig. 2) and its imperative to maximize discovery from all new transients. However, it is unreasonable to expect that zero-shot models will be able to perform well on data that is outside the range of prior telescopes; to build more general models, fine-tuning a pretrained FM (as done in Table 1) or using transfer learning (Gupta & Muthukrishna, 2025) is a better option.

A.4. Kepler Tasks

We take the labeled Kepler (Koch et al., 2010; Borucki et al., 2010) light curve training set of variable stars from (Audenaert et al., 2021). There are eight different classes, ranging from stochastic to periodic variability. The number of light curves per class is shown in Table 5. Detailed class descriptions of each class can be found in Audenaert et al. (2021).

Each of the light curves consists 1024 roughly evenly spaced brightness measurements at a 30 minute sampling rate. Because ZTF and LSST have a much lower sampling rate and fewer measurements per time series, we bin the *Kepler* light curves by taking mean of every five consecutive observations, resulting in 205 measurements per light curve. This more closely matches the number of measurements in our pretraining sample.

B. Model and Training Details

B.1. Classifier Training

Our neural network classifier is trained to take a vector of telescope observations as input, with each row being a distinct observation in the form $[\lambda_p, t_i, f_i]$, where λ_p represents

the median passband wavelength of the observation, t_i represents the time since the first observation in days, and f_i represents the flux (measured brightness) of the event. We scale the flux by dividing the measured fluxes by 500, a value chosen close to the mean of our various datasets. This limited preprocessing allows for the usage of our models in real time, however real-time evaluation is beyond the scope of this work. This input method (Huang et al., 2023; Gupta et al., 2024; Gupta & Muthukrishna, 2025) and specifically allows for the usage of the same model across surveys, something not facilitated by many previous input methods.

B.2. Adversarial Training Algorithm

Our adversarial pretraining is summarized in Algorithm 1. Here, X_i denotes the input light curve, c_i is its class label, and $O_i \in \{\text{ZTF}, \text{LSST}\}$ indicates the observatory. The latent representation is extracted as $L_i = C_L(X_i)$, where C_L is the penultimate layer of the classifier. The categorical cross-entropy loss is denoted by $H(p, q)$, where p is a predicted distribution and q is a target one-hot vector.

Algorithm 1 Adversarial Training

Require: Dataset $\{(X_i, c_i, O_i)\}_{i=1}^N$: light curves, class labels, and observatory labels

Require: Classifier C , Discriminator D

- 1: Initialize C and D with random weights
- 2: **repeat**
- 3: **// Step 1: Train the discriminator**
- 4: Freeze the classifier C
- 5: For each sample, compute latent representation $L_i = C_L(X_i)$
- 6: Compute discriminator loss:
- 7: Update D to minimize \mathcal{L}_D
- 8: **// Step 2: Train the classifier**
- 9: Freeze the discriminator D
- 10: Compute classifier loss with adversarial objective:

$$\mathcal{L}_D = H(D(L_i), O_i)$$

$$\mathcal{L}_C = H(C(X_i), c_i) - H(D(C_L(X_i)), O_i)$$

- 11: Update C to minimize \mathcal{L}_C
 - 12: **until** convergence
-

B.3. Contrastive Training Algorithm

For our supervised contrastive loss, we use the contrastive objective proposed in Chen et al. (2020) for model pretraining. It is formally defined as follows:

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{L}_i, \mathbf{L}_j) / \tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{L}_i, \mathbf{L}_k) / \tau)} \quad (1)$$

Task	Aperiodic	Constant	Contact	DSCT_BCEP	Eclipse	GDOR_SPB	Instr	RR_CEP	Solar
Kepler Classification	831	1000	2260	772	974	630	1171	63	1800

Table 5. The number of labeled data pairs available for each Kepler task used in this work.

where: $\mathbf{L}_i, \mathbf{L}_j$ are the latent representations from the classifier, $\text{sim}(\mathbf{L}_i, \mathbf{L}_j)$ denotes cosine similarity: $\text{sim}(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|}$, $\tau > 0$ is a temperature parameter that scales the similarity scores, and the loss is computed for all pairs (i, j) where X_i and X_j share the same class label.

The total supervised contrastive loss is computed by summing $\ell_{i,j}$ over all valid positive pairs in a batch. This encourages latent vectors from the same class to be close together, while implicitly pushing apart representations from other classes.

Our contrastive pretraining is summarized in Algorithm 2. We set $\tau = 0.5$ similar to the default set in previous work (Chen et al., 2020).

Algorithm 2 Supervised Contrastive Training

Require: Training set $\{(X_i, c_i)\}_{i=1}^N$; light curves, class labels

Require: Classifier C , temperature parameter τ

- 1: Initialize C with random weights
- 2: **repeat**
- 3: Compute classification loss

$$\mathcal{L}_C = H(C(X_i), c_i)$$

- 4: Compute latent representations $\mathbf{L}_i = C_L(X_i)$
- 5: Initialize total contrastive loss $\mathcal{L}_{\text{SCL}} \leftarrow 0$, counter $M \leftarrow 0$
- 6: **for** each anchor sample $i \in \{1, \dots, N\}$ **do**
- 7: Let $P(i) = \{j \neq i : c_j = c_i\}$
- 8: **for** each $j \in P(i)$ **do**
- 9: Compute pairwise contrastive loss:

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{L}_i, \mathbf{L}_j) / \tau)}{\sum_{k \neq i} \exp(\text{sim}(\mathbf{L}_i, \mathbf{L}_k) / \tau)}$$

- 10: Accumulate loss: $\mathcal{L}_{\text{SCL}} \leftarrow \mathcal{L}_{\text{SCL}} + \ell_{i,j}$
 - 11: Increment counter: $M \leftarrow M + 1$
 - 12: **end for**
 - 13: **end for**
 - 14: Compute mean contrastive loss: $\mathcal{L}_{\text{SCL}} \leftarrow \mathcal{L}_{\text{SCL}} / M$
 - 15: Compute total loss: $\mathcal{L} \leftarrow \mathcal{L}_{\text{SCL}} + \mathcal{L}_C$
 - 16: Update C to minimize \mathcal{L}
 - 17: **until** convergence
-

B.4. Fine-Tuning

B.4.1. GENERAL TASKS

We usually freeze the entire model when fine-tuning for downstream tasks. However, freezing the initial model does not work effectively for novel domains (most notably Kepler). Thus, when applying our foundation model to novel data, we unfreeze the initial neural network layers to assist in domain generalization. These techniques are motivated by past research done in transfer learning and they have been shown to improve performance (Gupta & Muthukrishna, 2025).

B.4.2. ANOMALY DETECTION

For anomaly detection, which is neither a regression or classification task, we use a classifier-based approach (Gupta et al., 2024), where a classifier is trained on a set of *normal* data. The penultimate layer of this classifier is then used as a latent space for anomaly detection and an isolation forest (Liu et al., 2008) is trained using this latent space to detect anomalies. This method has state-of-the-art (SoTA) performance for anomaly detection on real data² (Perez-Carrasco et al., 2023; Gupta et al., 2024).

B.5. Architecture Details

Our FMs are built with a recurrent neural network architecture containing Gated Recurrent Units (GRU; Cho et al., 2014). We chose to use GRUs because they are shown to be more effective than RNNs and have quicker training times than LSTMs (Chung et al., 2014). Further, neural network and GRU-based models have worked effectively in training past models for time-domain astronomy (e.g. Boone, 2019; Gupta et al., 2024; Muthukrishna et al., 2019). The provided code has more details describing the exact architecture.

When fine-tuning our FMs on downstream tasks for real data, we freeze the foundation model and leave the MLP unfrozen. When training on tasks from new observatories, we further unfreeze the initial layers. These decisions were based on rough hyperparameter searches. More precise tuning is beyond the scope of this work and the main contributions stand as long we remain consistent across different model types. The provided code also has more details on which layers are frozen.

²This dataset uses human-defined features extracted from light curves. However, they are not standardized across telescopes and thus we opt to use raw time series for our models.

We train and fine-tune our models using the Adam optimizer (Kingma & Ba, 2017) and stop training when the validation loss has not decreased for 5 epochs. Our classifier-based, adversarial, and contrastive models take roughly 10, 20, and 45 minutes to converge on a standard V100 GPU respectively. The final experiments required an estimated 15-25 GPU hours, however ideation and experimentation required considerably more.

Fig. 4 shows the training loss as a function of epoch for our models. Contrastive and adversarial models end at a worse cross-entropy loss but are able to simultaneously optimize for more complex loss functions. All models perform well on downstream tasks, as seen in Table 1, which shows that there is no clear metric to predict how well an FM will perform on downstream tasks. This opens up future research to use different techniques to achieve better performance.

C. Qualitative Latent Space Analysis

Fig. 5 and 6 show a UMAP (McInnes et al., 2020) visualization of the penultimate layer of our neural network classifiers. As seen, contrastive and adversarial models help unify the distributions of ZTF and LSST in the latent space down to the class level.

D. Further Analysis

D.1. Adversarial vs. Contrastive Loss

Incorporating the specialized training techniques proposed in this work does not improve model performance on downstream tasks (Table 1), which makes sense because the specialized loss functions are not designed for direct downstream tasks. In the case of the adversarial loss, however, we see a significant decrease in performance. We believe this is because the adversarial model lacks direct supervision and is forced to learn an implicit relationship between the two telescopes while training, whereas the contrastive model is given an explicit relationship through the supervised contrastive loss.

On zero-shot tasks, the adversarial and contrastive models outperform a classifier, showing that the unification in the latent space is indeed meaningful. The contrastive model slightly outperforms the adversarial model, similar to Table 1, and we think this is again due to the increased supervision provided during model training. Training with both losses simultaneously also does not improve performance in comparison to a purely contrastive model.

D.2. Performance on Simulations

We observe that fine-tuned performance on real-data lags the performance on the same tasks for simulations. In the redshift task for real data, our best model achieves an R^2 score

of 0.431 ± 0.053 , while the same metric for simulations is 0.580 ± 0.011 . Models can leverage these physics-based simulations as effective starting points but still require labeled real data to perform well. This gap between real data and simulations is why there is a significant performance gap between the Limited and Full evaluation scenarios. In other words, models trained on simulations need to be fine-tuned on real data to work well.

D.3. Anomaly Detection

Anomaly detection is the only task in which we do not see a performance improvement from the Limited to Full settings (Table 1). By definition, anomalies are objects that astronomers find interesting. By using human-defined simulations to pretrain FMs, they are naturally equipped to detect specifically what humans find interesting. The gap between simulations and real data is what anomaly detection pipelines are trying to fill. Further analysis of anomaly detection specifically is out of the scope of this work, however we hope that future researchers analyze the nature of this task and how domain expertise can be incorporated into it.

D.4. Are Classifiers the Best Foundation Models?

For our Classifier foundation models, we find that classification performance (as AUROC or crossentropy) is an effective metric for model selection, i.e. better performing classifiers perform better on downstream tasks. However, this method of model selection does not generalize to our domain-agnostic FMs (most notably the contrastive model). These models perform worse on the classification objective because they optimize for a more complex loss function (as seen in Fig. 4). However, they still perform well on downstream tasks, as noted in Table 1. This result shows that there is significant room for growth and that novel methods could incorporate information not captured by the models described in this work. Ultimately, this line of research does not end with classifiers.

For example, while we believe that supervised training is a promising direction for foundation models, at some point the complexity of physics-informed simulations may render it difficult to directly incorporate this information into deep learning models. Thus, using unsupervised methods alongside class-based supervision in model pretraining is also an important research direction. This is one of the many ways we see scientific foundation models expanding beyond classifiers with more complex pretraining methods.

E. kNN Zero-Shot Estimation

Aside from fine-tuning an MLP, we also evaluate using a k Nearest Neighbors approach for zero-shot estimation (Zhang et al., 2024; Parker et al., 2024). To perform zero-

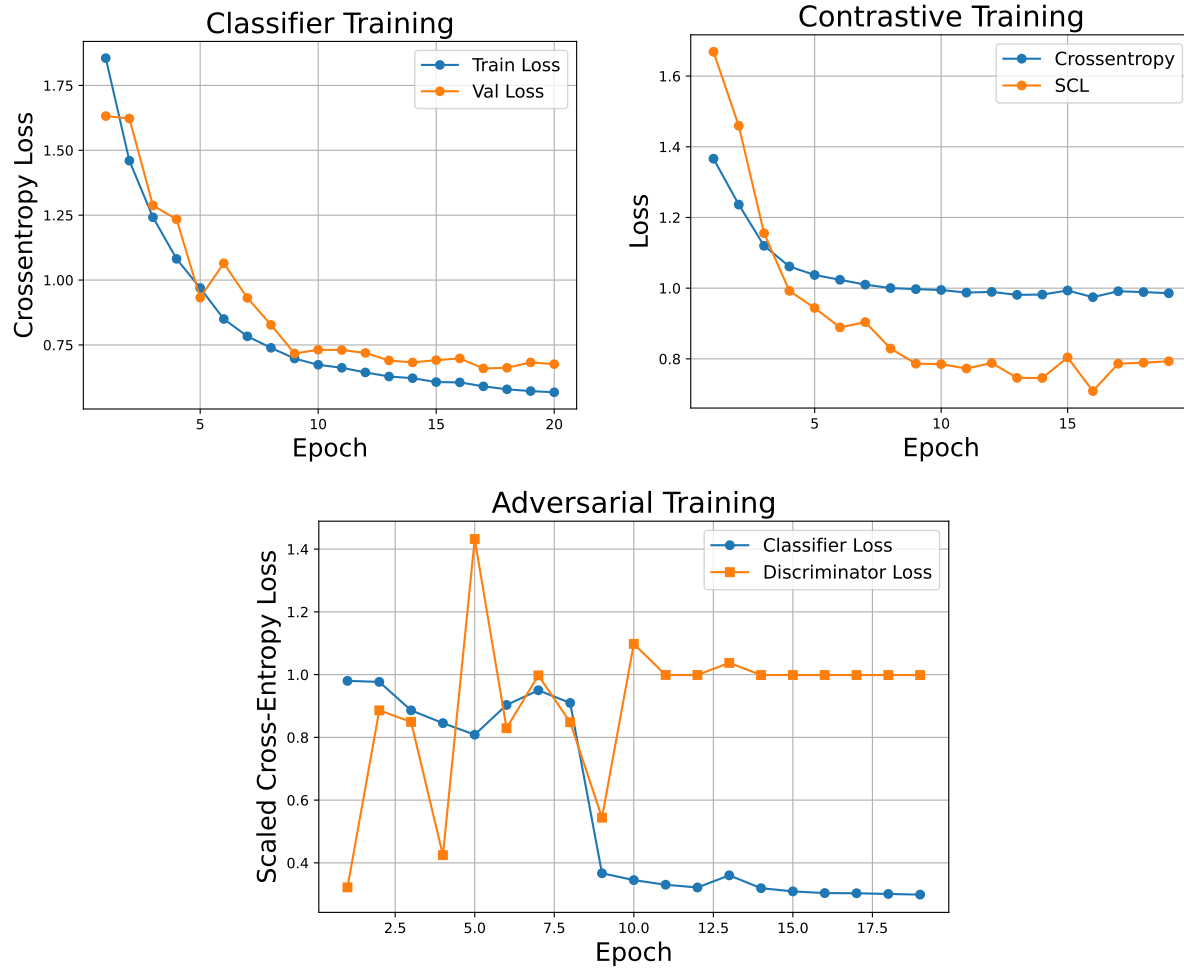


Figure 4. Loss as a function of training epoch for each FM in this work.

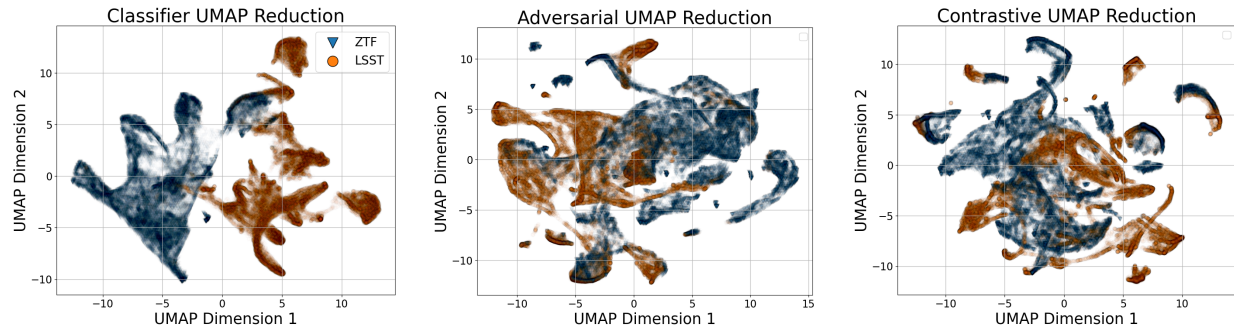


Figure 5. UMAP representations of the final layer of each foundation model. As seen, both the contrastive and adversarially trained models are able to unify LSST and ZTF transients into the same latent space, unlike a classifier which has distinct clusters.

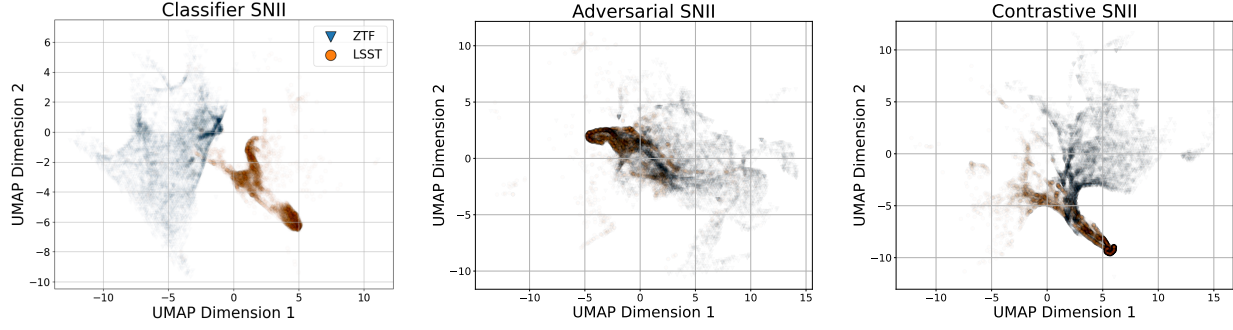


Figure 6. UMAP representations of the final layer of each foundation model only for a single class of supernova (Type II). As seen, both the contrastive and adversarially trained models are able to unify LSST and ZTF transients even at a class level.

shot redshift estimation for an LSST object, we first find the $k = 100$ closest ZTF embeddings to the LSST light curve embedding in the latent space. Then, we use the distance-weighted average of the corresponding redshifts to estimate the final redshift of the LSST object. As seen in Table 2, this zero-shot estimation method performs worse than using a directly trained MLP.