

# BENCHMARKING SELF-SUPERVISED VISION TRANSFORMERS IN ASTRONOMY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This work does not describe a novel method. Instead, it aims to extend the success of self-supervised pre-training on natural images to astronomical data. To address the lack of comprehensive benchmarks in astronomy, we first curate an unlabeled pre-training dataset and multiple datasets for typical astronomical tasks. Through extensive experiments, we demonstrate that our pre-training scheme has the following advantages. Representation transferability: pre-training followed by fine-tuning can not only significantly boost performance but also reduce training epochs compared to from-scratch training on downstream tasks (*e.g.*, improve 12% accuracy and reduce 83% epochs in galaxy classification), mirroring trends in natural image domains. Cross-instrument generalization: our pre-trained model generalizes across telescope instruments and outperforms domain models. Domain-specific pre-training data value: In-domain pre-training data further improves model performance and surpasses the results trained on general datasets such as ImageNet and other domain datasets. Furthermore, we explore Vision Transformers (ViTs) scaling in astronomy via parameter and data variation to offer insights and experiences for vision foundation model development in astronomy.

## 1 INTRODUCTION

The pre-training + transfer learning paradigm has revolutionized computer vision by significantly improving downstream tasks performance (Li et al., 2022; Fang et al., 2023) while substantially reducing computational and annotation costs (He et al., 2019; Dosovitskiy et al., 2020; He et al., 2022). Self-supervised learning methods, *e.g.*, contrastive learning (He et al., 2020; Chen et al., 2021; Chen & He, 2021; Grill et al., 2020) and masked image modeling (He et al., 2022; Chen et al., 2020; Bao et al., 2021; Xie et al., 2022; El-Nouby et al., 2024), eliminate manual annotation requirements, allowing pre-training on unprecedented scales (millions (Deng et al., 2009) to billions (Singh et al., 2023; Schuhmann et al., 2022) of images). However, despite its transformative impact on natural images, this paradigm remains largely unexplored in astronomy.

Astronomical data present two fundamental tensions that motivate our research. First, modern instruments generate petabyte-scale observational data (Dey et al., 2019; Mellier et al., 2024; Zhan, 2021), while curated annotations remain scarce. Second, intrinsic domain gaps – such as low signal-to-noise ratios (Sharma et al., 2020) and instrument-specific responses (Dey et al., 2019) – challenge direct transfer of natural image pre-training strategies. These disparities compel astronomers to spend years developing customized processing pipelines for individual instruments (Walmsley et al., 2020), underscoring the urgent need for foundation models with cross-instrument generalization capabilities.

As a first step, we examine three prerequisites for adapting self-supervised learning to astronomy:

(i) Pre-training method compatibility. Astronomical images share structural similarities with natural images, featuring three-band ( $g, r, z$ ) analogous to ( $r, g, b$ ) channels. This permits application of established pre-training methods that have been proven effective on natural images without obstacles. After pre-training, we fine-tune the models to adapt to task with a specific format, *e.g.*, five-band ( $u, g, r, i, z$ ) inputs for photometric redshift estimation task.

(ii) Pre-training data curation. Current astronomy lacks systematically curated pre-training datasets. To our knowledge, AstroSSL (Stein et al., 2022), updated version of Hayat et al. (2021);

054 Stein et al. (2021), represents the *only one*<sup>1</sup> that provides 76 million unlabeled images (named  
 055 `Astro-76M`) for self-supervised pre-training. However, the visual content in `Astro-76M` suffers  
 056 from morphological homogeneity (Figure 1), which is not conducive to learning powerful representa-  
 057 tion (Fan et al., 2021; Fang et al., 2022; Huang et al., 2022). We address this issue by supplementing  
 058 a subset with diversified image content (detailed description in Sec. 3).

059 (iii) Standard evaluation. Unlike natural images with standard benchmarks (*e.g.*, ImageNet-  
 060 1K (Deng et al., 2009), COCO (Lin et al., 2014), and ADE20K (Zhou et al., 2017)), vision models  
 061 for astronomy lack standard evaluation. Existing work employs ad hoc metrics (Hayat et al., 2021;  
 062 Hausen & Robertson, 2020; Lin et al., 2021; Zhang et al., 2022; Lanusse et al., 2023), hindering  
 063 comparative analysis. We build the first comprehensive benchmarks for typical astronomical tasks  
 064 that can serve as standard pre-training testbed.

065 Having completed the aforementioned preparations, we focus on investigating three key questions:  
 066

- 067 • Transferability: Can the pre-training + transfer learning paradigm generalize to astronomy?
- 068 • Cross-instrument generalization: Do pre-trained models achieve cross-telescope adaptability?
- 069 • Domain-specific data value: Does natural image pre-training (*e.g.*, ImageNet (Deng et al., 2009))  
 070 benefit astronomical tasks? Does in-domain pre-training data further improve performance?

071 Our empirical results demonstrate the advantages of pre-training in astronomical contexts:  
 072

- 073 • Our pre-trained models achieve substantial performance enhancements over from scratch, *e.g.*,  
 074 increasing morphology classification accuracy by 12% (Tables 2-4).
- 075 • Our pre-trained models demonstrate robust cross-instrument generalization capabilities, even  
 076 outperforming strong domain models (Table 5).
- 077 • Domain-specific pre-training substantially exceeds natural image pre-training in efficacy, particu-  
 078 larly for redshift estimation task: pre-training on astronomical data reduces prediction error by  
 079 64.6% than ImageNet pre-training (Table 6).
- 080 • Our pre-trained models show substantial improvements over existing ones, due to our enhanced  
 081 pre-training dataset and efforts in astronomy-specific optimizations.

082 Our findings reveal that the pre-training + transfer learning paradigm’s success in natural images  
 083 is not merely replicable but amplifiable in scientific domains through curating pre-training data,  
 084 designing comprehensive evaluation benchmarks and configuring model and training procedures.

## 085 2 RELATED WORKS

086 **Pre-training and Transfer Learning.** Modern vision tasks adopt a pre-training + transfer learning  
 087 paradigm: a general-purpose, task-agnostic backbone is pre-trained through supervised (He et al.,  
 088 2016; Dosovitskiy et al., 2020) or self-supervised learning (He et al., 2020; Chen et al., 2021; He  
 089 et al., 2022), whose structure is later modified and adapted to the downstream tasks. Self-supervised  
 090 pre-training, eliminates manual annotation requirements while enabling unprecedented scalability –  
 091 leverage millions (Deng et al., 2009) even billions (Singh et al., 2023; Fan et al., 2025) of training  
 092 samples, making it a powerful alternative to supervised training. The key advantage of pre-training is  
 093 that it can not only significantly boost performance but also reduce training costs on downstream tasks.  
 094 Our work investigates whether these established benefits can be effectively extended to astronomical  
 095 data, where unique observational characteristics present both opportunities and challenges.

096 **Pre-training Data.** The efficacy of deep learning models fundamentally depends on training data  
 097 quality and scale (Goyal et al., 2021; Chen et al., 2023). Within the pre-training paradigm, massive  
 098 unlabeled data like ImageNet (Deng et al., 2009) and LAION-5B (Schuhmann et al., 2022) serve  
 099 as foundational resources, enabling the development of vision foundation models through scalable  
 100 self-supervised learning. To date, the astronomy community possesses only one self-supervised  
 101 pre-training benchmark – the `Astro-76M` dataset (Stein et al., 2022), where focus on galaxies  
 102 with accurate morphological classification by selecting targets with sufficient pixel numbers and  
 103 low redshift values (Figure 1). Models pre-trained on such low-diversity dataset exhibit bias toward  
 104 “clean” morphological features (Stein et al., 2022; Lanusse et al., 2023), hindering their ability to  
 105

106 <sup>1</sup> Angeloudi et al. (2024) provides a larger dataset, but neither releases pre-trained models nor conducts  
 107 benchmark studies.

Table 1: **Summary of model architectures, pre-train data, and train recipes for C-MAE.** We employ models of varying scales and pre-training data sizes to investigate their scaling properties. We calculate the channel mean and standard deviation of the DESI-2M dataset for normalization, and search for the optimal hyperparameters, like base learning rate (blr) and weight decay (wd).

encoder	decoder embed, depth, heads	parameters	pre-train data	data augmentation					
ViT-B/L/H	512, 1, 16	100/300/600M	DESI-2M	rand size+crop+flip					
input size	mask ratio	optimizer	$(\beta_1, \beta_2)$	batch size	wd	blr	epochs	precision	GPUs
224×224	75%	AdamW	(0.9, 0.95)	4096	0.05	2e-4	800	fp16	64A100

process real-world observational data containing low signal-to-noise ratios objects. In this work, we build a new pre-training dataset by supplementing a subset with diversified image content. Previous finding (Fan et al., 2025) demonstrates that pre-training data distribution significantly impacts learned representations, underscoring the necessity of our augmented pre-training data.

**Downstream Task Benchmarking.** The computer vision community has established high-quality benchmark datasets, such as ImageNet-1K (Deng et al., 2009), COCO (Lin et al., 2014), and ADE20K (Zhou et al., 2017), to accurately evaluate model transfer capabilities. Astronomy, as a discovery-driven science, prioritizes scientific objectives over standardized methodological evaluations. For example, existing work (Stein et al., 2022; Hausen & Robertson, 2020; Lin et al., 2021; Zhang et al., 2022; Lanusse et al., 2023) predominantly employ ad hoc metrics tailored to individual instrument or object, precluding meaningful cross-study comparisons. To bridge this gap, our work establishes the first standard pre-training testbed and performs comprehensive evaluation for astronomical foundation models. We hope that these fundamental efforts on pre-training data and task-specific benchmarks can facilitate the development of artificial intelligence in astronomy.

### 3 PRE-TRAINING

Table 1 summarizes the pre-training settings, introduced next.

**Pre-training method.** In this work, we aim to extend the success of self-supervised pre-training on natural images to astronomical data, rather than designing a new method. Contrastive learning (Chen et al., 2021; Chen & He, 2021) will consume nearly  $2\times$  training resources, we consider using the more simple and efficient masked image modeling methods (He et al., 2022; Xie et al., 2022; Bao et al., 2021). Our pre-training follows the standard MAE (He et al., 2022) which takes the original image pixels as reconstruction target, as there was no image tokenizer available, prior to our work. During pre-training, we randomly mask a large subset (*e.g.*, 75%) of image patch and input the remaining visible patches to encoder. The encoded patches and mask tokens are processed by a lightweight decoder that reconstructs the original image in pixels. Building on the standard MAE framework and its application to data processing for the Chinese Space Station Telescope (CSST) (Zhan, 2021), we name our pre-training **C-MAE**.

**Architecture.** We use an asymmetric encoder-decoder architecture, similar to MAE. The encoder is a vanilla ViT (Dosovitskiy et al., 2020) with different parameters, to explore scaling behavior. The decoder adopts a lightweight design, which can further accelerate training.

**Pre-training data.** Contemporary astronomy is characterized by an abundance of observational data from large-scale sky surveys, *e.g.*, the Dark Energy Spectroscopic Instrument (DESI (Dey et al., 2019)) projects to accumulate 10 PB of data by the conclusion of 2025 and the CSST is expected to produce annual datasets exceeding 30 PB. However, current astronomy lacks systematically curated pre-training datasets. To the best of our knowledge, AstroSSL (Stein et al., 2022) currently represent the *sole one* in the literature that systematically curated 76 million unlabeled astronomical images (sampled from DESI, designated as Astro-76M) for self-supervised backbone pre-training via contrastive learning (He et al., 2020). This foundational effort has subsequently informed multiple research trajectories: LVM (Fu et al., 2024) adopted the Astro-76M corpus for autoencoder-based representation learning, while AstroCLIP (Lanusse et al., 2023) extended the paradigm through image vs. spectrum contrastive objectives with additional 0.2M pair data.

162 However, the `Astro-76M` dataset exhibits limited  
 163 semantic complexity in its image content  
 164 (Figure 1, row-1), where celestial objects pri-  
 165 marily manifest low-variability morphological  
 166 profiles. We posit that such structural homogene-  
 167 ity creates detrimental conditions for learning  
 168 strong representations (Fan et al., 2021; Fang  
 169 et al., 2022), as models tend to develop trivial  
 170 solutions through low-level interpolation, even  
 171 with a high mask ratio (Huang et al., 2022).

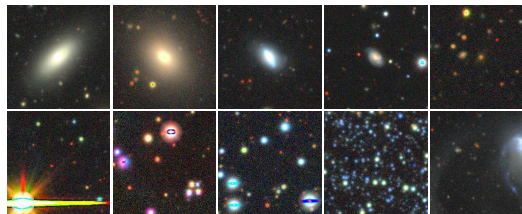


Figure 1: **Example training images.** The first row of images are from `Astro-76M`. The second row is from the 0.5M subset we added.

172 In this work, we conduct a new pre-train dataset  
 173 that sample a subset (1.5M) from `Astro-76M`, along with supplementing a subset (0.5M) with more  
 174 complex image content (row-2 in Figure 1). We refer to the 2M unlabeled data as `DESI-2M`. Table 3  
 175 quantifies the contribution of the 0.5M diverse samples. Further details on dataset construction are  
 176 provided in the Appendix A. Moreover, we uniformly sample 1M and 0.5M data from `DESI-2M`,  
 177 creating `DESI-1M` and `DESI-0.5M`, respectively. The consistent data distribution ensures fair  
 178 comparisons of data scaling experiments. We employ `DESI-0.5M` to search for pre-training hyper-  
 179 parameters, while using `DESI-1M` for default model pre-training under the optimized configuration.

## 181 4 EXPERIMENTS

182 We develop C-MAE variants with parameter counts of 100M, 300M, and 600M, to test the scaling  
 183 capability. These variants follow standard ViTs architectures (Dosovitskiy et al., 2020) and scaling  
 184 rules (El-Nouby et al., 2024). The detailed model configurations are shown in Appendix B. We  
 185 conduct extensive experiments on typical astronomical tasks, and also evaluate models on more  
 186 realistic and valuable cross-telescope settings. We will open source pre-train data, task-specific  
 187 datasets, and our code to foster future research.

### 190 4.1 CAN THE PRE-TRAINING + TRANSFER LEARNING PARADIGM GENERALIZE TO 191 ASTRONOMY?

192 Our insight is that for this paradigm to be valid, pre-training must satisfy two criteria: (1) achieving  
 193 superior results with less training cost (*e.g.*, fewer training epochs or less annotated data) compared to  
 194 trained from scratch; (2) pre-trained representations can be transferred to various downstream tasks.  
 195 Next, we demonstrate that C-MAE has these capabilities.

#### 197 4.1.1 GALAXY MORPHOLOGY CLASSIFICATION

199 **Dataset.** We collect 20K galaxy images from DESI. These images share the same source telescope as  
 200 our pre-training data, and we named `galaxy-desi`. Its eight categories are derived from a series of  
 201 questions and answers organized by the famous Galaxy Zoo 2 project (Willett et al., 2013), *e.g.*,  
 202 round elliptical, cigar-shaped elliptical, barred spiral. We split the training and test set into a 0.8/0.2  
 203 ratio. Moreover, we randomly sample subsets from the training set for few-shot experiments.

204 **Baseline.** We train classic ResNet (He et al., 2016) and ViT from scratch as baselines. Additionally,  
 205 the baseline also includes three publicly available pre-trained models: LVM (Fu et al., 2024),  
 206 AstroSSL (Stein et al., 2022), and AstroCLIP (Lanusse et al., 2023). All the three models were  
 207 pre-trained on `Astro-76M` dataset. Differently, LVM adopts a symmetrical autoencoder framework,  
 208 while AstroSSL and AstroCLIP use contrastive learning method. We search for learning rate, weight  
 209 decay, drop path rate, and epochs, for each model size (18, 50, 101, B, L, H) and for each model type  
 210 (ResNet, ViT, Swin). These hyper-parameters are included in the Appendix B.

211 **Results.** As detailed in Table 2, our C-MAE achieves significant improvement by 12% points  
 212 compared to scratch models (*e.g.*, 87.23% vs. 74.74% from ResNet-101). This performance gain is  
 213 attributed to pre-training, as only 50 epochs of fine-tuning were required (*vs.* 300 epochs from scratch).  
 214 Other pre-trained models (LVM, AstroSSL, and AstroCLIP) also demonstrate higher accuracy and  
 215 lower training costs than trained from scratch. This confirms the effectiveness of pre-training +  
 transfer learning paradigm in astronomical data.

Table 2: **Galaxy morphology classification on galaxy-desi**. Our C-MAE achieves significant improvements in comparison to both scratch models and existing pre-trained models, achieving state-of-the-art results. AE: autoencoder, CL: contrastive learning, MIM: masked image modeling.

method	backbone	pre-training	parameters	epochs	accuracy (%)
scratch	Res-18	-	11.18M	300	71.64
	Res-50	-	23.52M	300	73.89
	Res-101	-	42.52M	300	74.74
	ViT-Ti	-	5.53M	300	72.25
	ViT-S	-	21.67M	300	73.94
	ViT-B	-	85.80M	300	67.27
	ViT-L	-	303.31M	300	73.02
LVM	Swin-T	Astro-76M, AE	56.53M	50	84.63
AstroSSL	Res-50	Astro-76M, CL	23.52M	50	84.18
AstroCLIP	ViT-L	Astro-76M, CL	303.31M	50	82.57
<b>C-MAE (ours)</b>	ViT-B	DESI-1M, MIM	85.80M	50	<u>87.23</u>
	ViT-L	DESI-1M, MIM	303.31M	50	88.38
	ViT-H	DESI-1M, MIM	681.27M	50	<b>89.10</b>

Surprisingly, C-MAE outperforms models pre-trained on Astro-76M by 2.6%–4.5% accuracy using significantly less data. This large performance improvements can be explained by two factors.

(1) Our augmented pre-training data. Recalling our analysis in Sec. 3, the images in Astro-76M exhibit low-variation morphological attributes. Such simplistic image content is not conducive to learning powerful visual representations. In contrast, our DESI-1M dataset contains richer and more diverse image data. The 3.0% accuracy gains in C (with diverse samples) versus B – with identical data volume and the base Astro-0.5M samples in Table 3 – directly quantifies the contribution of the 0.5M diverse samples, using ViT-B as encoder and C-MAE for pre-training in both cases.

(2) Our effort on astronomy-specific optimizations. We ablate key designs, details in Appendix C, and observe differences from natural images pre-training. We argue that these gaps are mainly due to the low signal-to-noise ratio of astronomical images. These redesigns significantly improve the learned representation, leading to 1.8% accuracy gains. In contrast, existing models (LVM, AstroSSL, and AstroCLIP) directly adopt the pre-training approaches developed for natural images.

Figure 2 shows few-shot results. Our C-MAE pre-trained ViT-B achieves 69.55% accuracy with only 800 labeled images, surpassing the best scratch-trained model (67.27% accuracy with 16,000 labeled images). This demonstrates C-MAE’s capability to reduce the need for labeled data in downstream tasks. Furthermore, C-MAE consistently outperforms LVM across different training data sizes, and scaling the backbone leads to higher accuracy.

Table 3: **Quantitative contribution** of the 0.5M diverse samples using three different pre-training data. **A**: 0.5M images randomly sampled from Astro-76M; **B**: A + another 0.5M images from Astro-76M, for a total of 1M samples; **C**: A + the 0.5M diverse samples. C differs from the DESI-1M dataset, which can be roughly regarded as “Astro-0.75M + 0.25M diverse samples”.

pre-training data	A	B	C
accuracy (%)	83.17	84.39	87.46

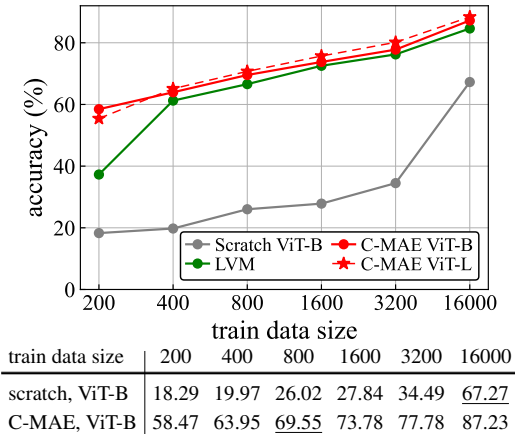


Figure 2: **Few-shot morphology classification** on galaxy-desi test set. Our C-MAE pre-training achieves higher accuracy than trained from scratch while using less labeled data, e.g., 800 vs. 16000.

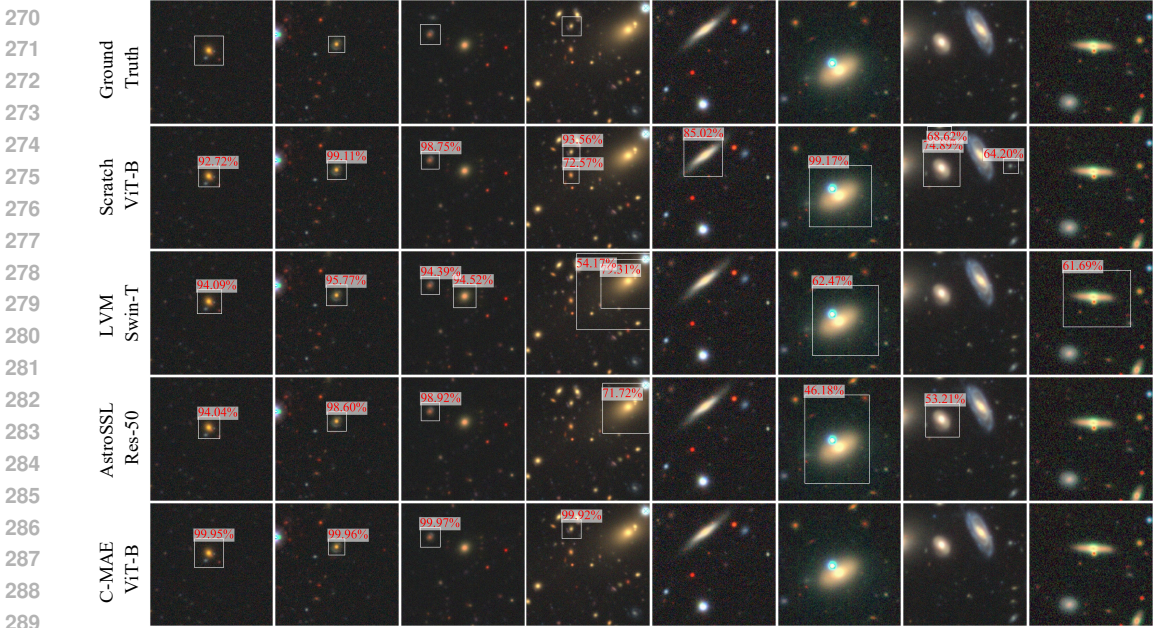


Figure 3: **Visualization results** on `neuralens-desi` test set. All baseline methods (2-4 rows) exhibit systematic output false positives results (3-8 columns), our C-MAE shows no such issue.

Table 4: **Lensing detection** on `neuralens-desi`. Pre-trained models outperform from-scratch using fewer training iterations, confirming the transferability to object detection. The visualized outputs are shown in Figure 3.

method	backbone	iters.	AP	AP <sup>50</sup>	AP <sup>75</sup>	AP <sup>S</sup>	AP <sup>M</sup>	AR@1	AR@10
scratch	Res-50	64K	21.29	45.21	16.87	24.01	25.06	37.34	57.65
	ViT-B	64K	28.95	44.36	32.37	32.87	35.76	29.26	55.42
LVM	Swin-T	32K	31.38	59.65	27.90	36.99	37.30	38.64	61.16
AstroSSL	Res-50	32K	31.76	58.37	29.95	31.69	37.56	39.58	60.44
AstroCLIP	ViT-L	32K	27.53	54.84	23.49	31.14	33.25	34.94	56.38
<b>C-MAE (ours)</b>	ViT-B	32K	39.88	49.22	45.70	38.84	81.74	41.83	42.54
	ViT-L	32K	41.31	49.37	47.04	39.72	85.19	43.39	43.74
	ViT-H	32K	<b>42.62</b>	51.46	49.43	41.97	87.64	45.10	45.68

#### 4.1.2 STRONG GRAVITATIONAL LENSING DETECTION

Strong gravitational lenses, caused by foreground massive galaxies or clusters, distort the light from background sources. Their detection (Jia et al., 2022b; Li et al., 2024) is challenging due to their extremely low spatial density compared to ordinary galaxies. Consequently, detecting them requires an extremely low false-positive rate to be scientifically useful.

**Dataset.** Following LVM (Fu et al., 2024), we create a dataset from the `NeuraLens` database (Huang et al., 2021), providing center coordinates of candidate lensing systems. We generate  $256 \times 256$  pixel cutouts (3-channel  $g, r, z$  composite at 0.262 arcsec/pixel) centered on each candidate with a  $[-256, 256]$  pixel random offset for augmentation. Bounding-box annotations were created using OpenCV and Labelme<sup>2</sup> tools. The resulting dataset, named `neuralens-desi`, contains 11.6K images (6K that show strong lensing) and is divided into training and test sets with 88%/12% ratio.

**Implementation.** We use Detectron2 (Wu et al., 2019) with the influential Mask R-CNN framework (He et al., 2017). For ViT and Swin backbones, we follow ViTDet (Li et al., 2022), which is specially designed for vision transformer. We train ResNet and ViT/Swin backbones with SGD and AdamW (Loshchilov & Hutter, 2019), respectively. The input size is  $512 \times 512$ , augmented during training by large-scale jitter (Ghiasi et al., 2021) with a scale range of  $[0.1, 2.0]$ . All methods differ only in their backbone (type and pre-training), which ensures alignment of low-level details.

<sup>2</sup><https://opencv.org/>, <https://labelme.io/>

Table 5: **Transferring pre-trained models to SDSS telescope.** Photo-Net is a strong baseline and significantly outperforms other pre-trained models. Our C-MAE outperforms the ensemble results. †: pre-trains backbone by contrastive learning; ‡: combines 6 models.

method	backbone	cross	galaxy-sdss		redshift-sdss			
		telescope	epochs	accuracy (%)	epochs	$\Delta z \downarrow$	$\sigma_{\text{MAD}}$	$\eta$ (%)
scratch	Res-50	✗	100	85.46	50	16.31e-4	2.57e-2	0.92
	ViT-B	✗	100	87.38	50	4.32e-4	2.43e-2	0.83
SC-Net <sup>†</sup>	6conv+2fc	✗	50	89.13	-	-	-	-
Photo-Net <sup>‡</sup>	Inception	✗	-	-	-	1.70e-4	1.43e-2	1.26
LVM	Swin-T	✓	50	93.48	50	13.82e-4	40.81e-2	0.48
AstroSSL	Res-50	✓	50	95.26	50	7.37e-4	1.98e-2	1.11
AstroCLIP	ViT-L	✓	50	93.03	50	3.47e-4	2.38e-2	0.86
C-MAE (ours)	ViT-B	✓	50	<u>96.00</u>	50	<u>0.97e-4</u>	2.31e-2	1.09
	ViT-L	✓	50	95.76	50	0.77e-4	1.42e-2	1.57
	ViT-H	✓	90	<b>96.02</b>	85	<b>0.51e-4</b>	0.56e-2	2.43

**Results.** Table 4 reports results in COCO format (Lin et al., 2014), including AP (averaged over IoU thresholds), AP<sup>50</sup>, AP<sup>75</sup>, AP<sup>S</sup>, AP<sup>M</sup> (AP at different scales), and AR@1, AR@10. Pre-trained models (except AstroCLIP) outperformed models trained from scratch, requiring fewer iterations (32K vs. 64K). Our C-MAE achieved the best performance, with a +10.9 point improvement in AP and +12.6 in AR@1 compared to scratch-trained ViT-B (37% and 43% relative improvement). Example outputs are visualized in Figure 3. While other baselines exhibit systematic false positives (3-8 cols in Figure 3), C-MAE did not. This is of great practical significance for the subsequent astronomical data analysis due to the extremely unbalance distribution of strong gravitational lensing.

#### 4.2 DO PRE-TRAINED MODELS ACHIEVE CROSS-TELESCOPE ADAPTABILITY?

Cross-telescope data processing is vital for modern astronomical research, which involves integrating massive heterogeneous datasets from various observatories to support multi-wavelength and multi-messenger analysis. This study explores whether a pre-trained model on one telescope’s data (data from DESI (Dey et al., 2019)) can perform well on another’s (data from SDSS (York et al., 2000)), addressing a key issue in cross-device data processing.

**Dataset.** We evaluate models on two different tasks.

- The `galaxy-sdss` is a classification dataset (Zhang et al., 2022) derived from SDSS data release (Alam et al., 2015), comprising 5 classes, 23,037 training images, and 5,754 validation images. Validation accuracy is reported.
- The `redshift-sdss` was constructed using a sample of 12 major galaxies from the SDSS data release (Alam et al., 2015), accessible via the CasJobs interface<sup>3</sup>. After applying the SQL query from (Pasquet et al., 2019) and excluding objects with poor spectroscopic measurements ( $Z_{\text{WARNING}} \neq 0$ ), 41,650 galaxies were selected with  $r$ -band-dereddened Petrosian Magnitude  $r < 17.77$  (survey completeness limit) and reliable spectroscopic redshifts ( $0.01 < z < 0.3$ ). The input images are  $224 \times 224$  pixels (0.396 arcsec/pixel) with 5-band ( $u, g, r, i, z$ ), retrieved from the SDSS Science Archive server using the Astroquery package (Ginsburg et al., 2019). The dataset is divided into a training set (10,100 images) and a test set (31,550 images). Results are evaluated using common statistical metrics (Hogg et al., 2000): residuals  $\Delta z$ , median absolute deviation  $\sigma_{\text{MAD}}$ , and outliers fraction  $\eta$ .

**Baseline.** We build two additional baselines that were trained on SDSS (without cross-telescope). Specifically, SC-Net (Zhang et al., 2022) pre-trains backbone on the `galaxy-sdss` training set by contrastive learning (Chen & He, 2021) and then conducts supervised classification; Photo-Net (Pasquet et al., 2019) is an ensemble model integrating six sub-models for redshift estimation. We replicate their reported results. For other baselines and our C-MAE, we attach a linear regressor to the final features of the backbone network for adaptation to the redshift estimation task. The loss function computes the mean squared error between the prediction and truth redshift.

<sup>3</sup><https://skyserver.sdss.org/CasJobs/>

Table 6: **Compare pre-training strategies** with ViT-B backbone. Although ImageNet pre-training can help learn astronomical tasks, astronomical data pre-training can further improve performance. IN-1K and DESI-1M have similar amount of training data.

pre-train	galaxy-desi	galaxy-sdss	neuralens-desi		redshift-sdss		
	accuracy	accuracy	AP	AR@1	$\Delta z \downarrow$	$\sigma_{MAD}$	$\eta$ (%)
none (random init.)	67.27	87.38	28.95	29.26	4.32e-4	2.43e-2	0.83
IN-1K, supervised	86.44	94.96	33.69	39.30	2.74e-4	0.75e-2	1.75
IN-1K, MAE	82.82	94.35	35.10	40.85	13.64e-4	1.93e-2	1.07
DESI-1M, C-MAE	87.23	96.00	39.88	41.83	0.97e-4	2.31e-2	1.09

**Results.** Table 5 presents the comparison results. For the `galaxy-sdss` classification (5-class,  $\sim 23K$  images), all pre-trained models demonstrated significant performance improvements, with the weakest AstroCLIP outperforming SC-Net by 4% accuracy (93.03% vs. 89.13%).

However, for the more challenging redshift estimation task (limited training data), pre-trained models (LVM, AstroSSL, and AstroCLIP) performed markedly worse than Photo-Net and even underperformed trained from scratch (*e.g.*,  $7.37e-4$  vs.  $4.32e-4$ ). This highlights the difficulty in transferring across telescopes. Notably, our C-MAE (single-model result) achieved a residual of  $0.97e-4$ , outperforming Photo-Net by 43% relative improvement, demonstrating that a well-designed pre-trained model can effectively generalize across different telescopes. Figure 4 illustrates the visualization of redshift predictions.

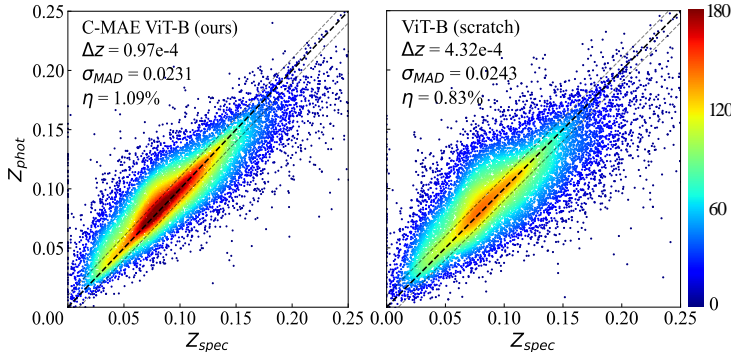


Figure 4: **Photo-z estimates on redshift-sdss test set.** Compare the fine-tuned C-MAE (left) to the equivalent scratch network (right). Figure 4 illustrates the visualization of redshift predictions.

#### 4.3 DOES NATURAL IMAGE PRE-TRAINING BENEFIT ASTRONOMICAL TASKS?

The ImageNet-1K (IN-1K) dataset (Deng et al., 2009) is widely adopted in computer vision, owing to its potential to improve downstream task performance (Dosovitskiy et al., 2020; He et al., 2022; 2020). We transfer representations learned on IN-1K to astronomy tasks, as shown in Table 6, where evaluates both supervised and self-supervised models, the latter using MAE for pre-training.

Results demonstrate that supervised pre-training on IN-1K yields significant improvements over no-pretraining baseline, achieving gains of 19% on `galaxy-desi`, 4.7% AP on `neuralens-desi`, and reducing  $\Delta z$  by 37% on `redshift-sdss`. In contrast, MAE pretraining exhibits mixed efficacy: while it surpasses the no-pretraining in lensing detection, it produces bad result on redshift estimation, far worse than the non-pretrained model.

We hypothesize that supervised pre-training fosters discriminative feature representations critical for cross-domain adaptation, whereas MAE’s focus on reconstructing local image patches prioritizes spatially oriented features (Li et al., 2022; Park et al., 2023), which can benefit spatial tasks such as detection but are less transferable to regression-based tasks such as redshift estimation.

Notably, despite the DESI-1M dataset’s similar scale to IN-1K (1M vs. 1.28M images), our C-MAE achieves superior performance in all tasks. This highlights that while natural image pre-training provides a foundational advantage, domain-specific pre-training on astronomical data offers substantial additional benefits. These findings underscore the interplay between pre-training objectives and downstream task requirements in domain shift scenarios.

We have included more results from natural image pre-training (DINO-v2 (Oquab et al., 2023) and CLIP (Radford et al., 2021)) in the Appendix C. Although these models utilized larger pre-training datasets or more sophisticated training procedures, they still underperform compared to our C-MAE, which further strengthens our conclusion.

Table 7: **Scaling pre-training data.** Increasing pre-training data from 0.5M to 2M, all models show better results, but tend to saturate. Both large scale image samples and sufficient training iterations are important for pre-training. The lower  $\Delta z$  is better.

pre-train data	epochs	accuracy on galaxy-desi			$\Delta z$ on redshift-sdss		
		ViT-B	ViT-L	ViT-H	ViT-B	ViT-L	ViT-H
same pre-training epochs							
DESI-0.5M	800	85.95	87.05	87.43	1.72e-4	1.14e-4	1.29e-4
DESI-1M	800	87.23	88.38	89.10	0.97e-4	0.77e-4	0.51e-4
DESI-2M	800	87.65	88.41	89.24	1.13e-4	0.63e-4	0.33e-4
DESI-1M	1600	87.38	88.32	89.17	1.04e-4	0.76e-4	0.60e-4
DESI-2M	400	86.38	87.94	88.52	1.38e-4	0.89e-4	0.92e-4

#### 4.4 SCALING PRE-TRAINING DATA.

We conduct ablation studies on key components of our pre-training framework (*e.g.*, decoder design, mask ratios, and training epochs), observing behaviors differ from MAE’s image pre-training (He et al., 2022) (see Appendix C for details).

This section focuses on evaluating how scaling pre-training data impacts downstream performance. Our prior experiments primarily used 1 million images for pre-training. Table 7 now examines the effects of data scale on galaxy morphology classification and redshift estimation. While performance improves when increasing pre-training data from 0.5M to 2M, we observe diminishing returns (*cf.* DESI-1M *vs.* DESI-2M with 800 epochs). This aligns with prior findings (Singh et al., 2023; Xie et al., 2023) showing limited benefits of data scaling in vision models – for instance, MAE pre-trained ViT-L achieves only 1.0% improvement despite  $14\times$  more pre-training data and 512px resolution (Singh et al., 2023), unlike the strong scaling law observed in language models (Brown et al., 2020; Kaplan et al., 2020). Developing vision learning methods that effectively utilize large-scale data remains an open challenge.

To further determine whether performance gains are due to more samples or more iterations, we keep the total number of seen samples constant across different data scales by adjusting training epochs. Extending DESI-1M training from 800 to 1,600 epochs shows modest gains but still underperformed DESI-2M trained for 800 epochs (same seen samples). This suggests that both data scale and sufficient training iterations are crucial for representation quality. The performance gap between 400-epoch DESI-2M and 800-epoch DESI-1M further confirms this.

## 5 DISCUSSION AND CONCLUSION

**Societal impacts.** Pre-trained models may inherit biases from the pre-training data and generate inexistent content. Our work does not sufficiently evaluate C-MAE on imbalanced data distributions, nor investigate whether and how they are biased as well.

**Limitations and future work.** First, we did not design a new self-supervised learning method but utilized existing ones. Future work could consider incorporating astronomical data characteristics (*e.g.*, equipment noise, target uncertainty) into the pre-training scheme; Second, during the curating of the pre-training data, we can further consider its characteristics to provide a more systematic data diversity control scheme; Finally, in addition to fine-tuning, we also explored linear probing and visual prompt tuning (Jia et al., 2022a; Wang et al., 2025), but these approaches did not achieve satisfactory results. Future work should focus on designing more efficient algorithms specifically tailored for astronomy to reduce the adaptation cost of foundation models.

**Conclusion.** This work extends self-supervised pre-training to astronomical data. We create benchmarks by curating an unlabeled pre-training dataset and task-specific datasets. Key findings include: Pre-training improves downstream task performance and reduces training costs; Pre-trained models generalize well across telescopes; In-domain pre-training data enhance model adaptation. We further explore the parameters and data scaling capacity. We hope that this work can provide useful insights and experiences for developing vision foundation models in astronomy.

## REFERENCES

- 486  
487  
488 Shadab Alam, Franco D Albareti, Carlos Allende Prieto, Friedrich Anders, Scott F Anderson, Timothy  
489 Anderton, Brett H Andrews, Eric Armengaud, Éric Aubourg, Stephen Bailey, et al. The eleventh  
490 and twelfth data releases of the sloan digital sky survey: final data from sdss-iii. *The Astrophysical  
491 Journal Supplement Series*, 219(1):12, 2015.
- 492  
493 Eirini Angeloudi, Jeroen Audenaert, Micah Bowles, Benjamin M Boyd, David Chemaly, Brian  
494 Cherinka, Ioana Ciucă, Miles Cranmer, Aaron Do, Matthew Grayling, et al. The multimodal uni-  
495 verse: enabling large-scale machine learning with 100 tb of astronomical scientific data. *Advances  
496 in Neural Information Processing Systems*, 37:57841–57913, 2024.
- 497  
498 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint  
499 arXiv:1607.06450*, 2016.
- 500  
501 Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers.  
502 *arXiv preprint arXiv:2106.08254*, 2021.
- 503  
504 Daniel Bolya, Po-Yao Huang, Peize Sun, Jang Hyun Cho, Andrea Madotto, Chen Wei, Tengyu Ma,  
505 Jiale Zhi, Jathushan Rajasegaran, Hanoona Rasheed, et al. Perception encoder: The best visual  
506 embeddings are not at the output of the network. *arXiv preprint arXiv:2504.13181*, 2025.
- 507  
508 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
509 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
510 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 511  
512 Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever.  
513 Generative pretraining from pixels. In *International conference on machine learning*, pp. 1691–  
514 1703. PMLR, 2020.
- 515  
516 Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian  
517 Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver,  
518 Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish Thapliyal, James  
519 Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme,  
520 Andreas Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. Pali: A jointly-  
521 scaled multilingual language-image model, 2023. URL <https://arxiv.org/abs/2209.06794>.
- 522  
523 Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of  
524 the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15750–15758, 2021.
- 525  
526 Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision  
527 transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp.  
528 9640–9649, 2021.
- 529  
530 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
531 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
532 pp. 248–255. Ieee, 2009.
- 533  
534 Arjun Dey, David J Schlegel, Dustin Lang, Robert Blum, Kaylan Burleigh, Xiaohui Fan, Joseph R  
535 Findlay, Doug Finkbeiner, David Herrera, Stéphanie Juneau, et al. Overview of the desi legacy  
536 imaging surveys. *The Astronomical Journal*, 157(5):168, 2019.
- 537  
538 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
539 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An  
540 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint  
541 arXiv:2010.11929*, 2020.
- 542  
543 Alaaeldin El-Nouby, Michal Klein, Shuangfei Zhai, Miguel Angel Bautista, Alexander Toshev,  
544 Vaishaal Shankar, Joshua M Susskind, and Armand Joulin. Scalable pre-training of large autore-  
545 gressive image models. *arXiv preprint arXiv:2401.08541*, 2024.

- 540 David Fan, Shengbang Tong, Jiachen Zhu, Koustuv Sinha, Zhuang Liu, Xinlei Chen, Michael  
541 Rabbat, Nicolas Ballas, Yann LeCun, Amir Bar, and Saining Xie. Scaling language-free visual  
542 representation learning. *arXiv preprint arXiv:2504.01017*, 2025.
- 543
- 544 Lijie Fan, Sijia Liu, Pin-Yu Chen, Gaoyuan Zhang, and Chuang Gan. When does contrastive learning  
545 preserve adversarial robustness from pretraining to finetuning? *Advances in neural information  
546 processing systems*, 34:21480–21492, 2021.
- 547
- 548 Alex Fang, Gabriel Ilharco, Mitchell Wortsman, Yuhao Wan, Vaishaal Shankar, Achal Dave, and  
549 Ludwig Schmidt. Data determines distributional robustness in contrastive language image pre-  
550 training (clip). In *International Conference on Machine Learning*, pp. 6216–6234. PMLR, 2022.
- 551
- 552 Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong  
553 Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale.  
554 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.  
19358–19369, 2023.
- 555
- 556 Mingxiang Fu, Yu Song, Jiameng Lv, Liang Cao, Peng Jia, Nan Li, Xiangru Li, Jifeng Liu, A-  
557 Li Luo, Bo Qiu, et al. A versatile framework for analyzing galaxy image data by implanting  
558 human-in-the-loop on a large vision model. *Chinese Physics C*, 2024.
- 559
- 560 Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and  
561 Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation.  
562 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.  
2918–2928, 2021.
- 563
- 564 Adam Ginsburg, Brigitta M Sipőcz, CE Brasseur, Philip S Cowperthwaite, Matthew W Craig,  
565 Christoph Deil, Austen M Groener, James Guillochon, Giannina Guzman, Simon Liedtke, et al.  
566 Astroquery: an astronomical web-querying package in python. *The Astronomical Journal*, 157(3):  
98, 2019.
- 567
- 568 Priya Goyal, Mathilde Caron, Benjamin Lefaudeux, Min Xu, Pengchao Wang, Vivek Pai, Mannat  
569 Singh, Vitaliy Liptchinsky, Ishan Misra, Armand Joulin, and Piotr Bojanowski. Self-supervised  
570 pretraining of visual features in the wild, 2021. URL [https://arxiv.org/abs/2103.  
571 01988](https://arxiv.org/abs/2103.01988).
- 572
- 573 Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena  
574 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,  
575 et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural  
information processing systems*, 33:21271–21284, 2020.
- 576
- 577 Ryan Hausen and Brant E Robertson. Morpheus: A deep learning framework for the pixel-level  
578 analysis of astronomical image data. *The Astrophysical Journal Supplement Series*, 248(1):20,  
2020.
- 579
- 580 Md Abul Hayat, George Stein, Peter Harrington, Zarija Lukić, and Mustafa Mustafa. Self-supervised  
581 representation learning for astronomical images. *The Astrophysical Journal Letters*, 911(2):L33,  
582 2021.
- 583
- 584 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
585 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
586 pp. 770–778, 2016.
- 587
- 588 Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *Proceedings of the  
IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- 589
- 590 Kaiming He, Ross Girshick, and Piotr Dollar. Rethinking imagenet pre-training. In *Proceedings of  
the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- 591
- 592 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
593 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on  
computer vision and pattern recognition*, pp. 9729–9738, 2020.

- 594 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked  
595 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer  
596 vision and pattern recognition*, pp. 16000–16009, 2022.  
597
- 598 David W Hogg, Judith G Cohen, and Roger Blandford. The caltech faint galaxy redshift survey. xii.  
599 clustering of galaxies. *The Astrophysical Journal*, 545(1):32, 2000.  
600
- 601 Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze,  
602 and Christoph Feichtenhofer. Masked autoencoders that listen. *Advances in Neural Information  
603 Processing Systems*, 35:28708–28720, 2022.
- 604 Xiaosheng Huang, Christopher Storfer, A Gu, V Ravi, A Pilon, W Sheu, R Venguswamy, S Banka,  
605 A Dey, M Landriau, et al. Discovering new strong gravitational lenses in the desi legacy imaging  
606 surveys. *The Astrophysical Journal*, 909(1):27, 2021.  
607
- 608 Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and  
609 Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision*, pp. 709–727.  
610 Springer, 2022a.
- 611 Peng Jia, Ruiqi Sun, Nan Li, Yu Song, Runyu Ning, Hongyan Wei, and Rui Luo. Detection of  
612 strongly lensed arcs in galaxy clusters with transformers. *The Astronomical Journal*, 165(1):26,  
613 2022b.  
614
- 615 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott  
616 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.  
617 *arXiv preprint arXiv:2001.08361*, 2020.  
618
- 619 Francois Lanusse, Liam Holden Parker, Siavash Golkar, Alberto Bietti, Miles Cranmer, Michael  
620 Eickenberg, Geraud Krawezik, Michael McCabe, Ruben Ohana, Mariel Pettee, et al. Astroclip:  
621 Cross-modal pre-training for astronomical foundation models. In *NeurIPS 2023 AI for Science  
622 Workshop*, 2023.
- 623 Xu Li, Ruiqi Sun, Jiameng Lv, Peng Jia, Nan Li, Chengliang Wei, Hu Zou, Xinzhong Er, Yun Chen,  
624 Zhang Ban, et al. Csst strong-lensing preparation: A framework for detecting strong lenses in the  
625 multicolor imaging survey by the china survey space telescope (csst). *The Astronomical Journal*,  
626 167(6):264, 2024.  
627
- 628 Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer  
629 backbones for object detection. In *European conference on computer vision*, pp. 280–296. Springer,  
630 2022.
- 631 Joshua Yao-Yu Lin, Song-Mao Liao, Hung-Jin Huang, Wei-Ting Kuo, and Olivia Hsuan-Min  
632 Ou. Galaxy morphological classification with efficient vision transformer. *arXiv preprint  
633 arXiv:2110.01024*, 2021.  
634
- 635 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
636 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–  
637 ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings,  
638 Part V 13*, pp. 740–755. Springer, 2014.  
639
- 640 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-  
641 ence on Learning Representations*, 2019.
- 642 Y Mellier, Abdurrouf Abdurrouf, JA Acevedo Barroso, A Achúcarro, J Adamek, R Adam, GE Addi-  
643 son, N Aghanim, M Aguena, V Ajani, et al. Euclid. i. overview of the euclid mission. *Astronomy  
644 & Astrophysics*, 2024.  
645
- 646 Igor Molybog, Peter Albert, Moya Chen, Zachary DeVito, David Esiobu, Naman Goyal, Punit Singh  
647 Koura, Sharan Narang, Andrew Poulton, Ruan Silva, et al. A theory on adam instability in  
large-scale machine learning. *arXiv preprint arXiv:2304.09871*, 2023.

- 648 Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov,  
649 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao  
650 Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran,  
651 Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut,  
652 Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision,  
653 2023.
- 654 Namuk Park, Wonjae Kim, Byeongho Heo, Taekyung Kim, and Sangdoon Yun. What do self-  
655 supervised vision transformers learn? In *The Eleventh International Conference on Learning*  
656 *Representations*, 2023.
- 657  
658 Johanna Pasquet, Emmanuel Bertin, Marie Treyer, Stéphane Arnouts, and Dominique Fouchez.  
659 Photometric redshifts from sdss images using a convolutional neural network. *Astronomy &*  
660 *Astrophysics*, 621:A26, 2019.
- 661 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
662 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
663 models from natural language supervision. In *International conference on machine learning*, pp.  
664 8748–8763. PmLR, 2021.
- 665 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi  
666 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An  
667 open large-scale dataset for training next generation image-text models. *Advances in neural*  
668 *information processing systems*, 35:25278–25294, 2022.
- 669  
670 Kaushal Sharma, Ajit Kembhavi, Aniruddha Kembhavi, T Sivarani, Sheelu Abraham, and Kaustubh  
671 Vaghmare. Application of convolutional neural networks for stellar spectral classification. *Monthly*  
672 *Notices of the Royal Astronomical Society*, 491(2):2280–2300, 2020.
- 673 Mannat Singh, Quentin Duval, Kalyan Vasudev Alwala, Haoqi Fan, Vaibhav Aggarwal, Aaron Ad-  
674 cock, Armand Joulin, Piotr Dollár, Christoph Feichtenhofer, Ross Girshick, et al. The effectiveness  
675 of mae pre-pretraining for billion-scale pretraining. In *Proceedings of the IEEE/CVF international*  
676 *conference on computer vision*, pp. 5484–5494, 2023.
- 677 George Stein, Peter Harrington, Jacqueline Blaum, Tomislav Medan, and Zarija Lukic. Self-  
678 supervised similarity search for large scientific datasets. *arXiv preprint arXiv:2110.13151*, 2021.
- 679  
680 George Stein, Jacqueline Blaum, Peter Harrington, Tomislav Medan, and Zarija Lukić. Mining for  
681 strong gravitational lenses with self-supervised learning. *The Astrophysical Journal*, 932(2):107,  
682 2022.
- 683 Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going  
684 deeper with image transformers. In *Proceedings of the IEEE/CVF international conference on*  
685 *computer vision*, pp. 32–42, 2021.
- 686  
687 Hugo Touvron, Matthieu Cord, Alaaeldin El-Nouby, Jakob Verbeek, and Hervé Jégou. Three things  
688 everyone should know about vision transformers. In *European Conference on Computer Vision*,  
689 pp. 497–515. Springer, 2022.
- 690 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
691 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*  
692 *systems*, 30, 2017.
- 693 Mike Walmsley, Lewis Smith, Chris Lintott, Yarin Gal, Steven Bamford, Hugh Dickinson, Lucy  
694 Fortson, Sandor Kruk, Karen Masters, Claudia Scarlata, et al. Galaxy zoo: probabilistic morphology  
695 through bayesian cnns and active learning. *Monthly Notices of the Royal Astronomical Society*,  
696 491(2):1554–1574, 2020.
- 697  
698 Yuzhu Wang, Manni Duan, and Shu Kong. Attention to the burstiness in visual prompt tuning! In  
699 *International Conference on Computer Vision (ICCV)*, 2025.
- 700  
701 Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer.  
Masked feature prediction for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF*  
*conference on computer vision and pattern recognition*, pp. 14668–14678, 2022.

702 Kyle W Willett, Chris J Lintott, Steven P Bamford, Karen L Masters, Brooke D Simmons, Kevin RV  
703 Casteels, Edward M Edmondson, Lucy F Fortson, Sugata Kaviraj, William C Keel, et al. Galaxy  
704 zoo 2: detailed morphological classifications for 304 122 galaxies from the sloan digital sky survey.  
705 *Monthly Notices of the Royal Astronomical Society*, 435(4):2835–2860, 2013.

706 Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2.  
707 <https://github.com/facebookresearch/detectron2>, 2019.

708  
709 Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu.  
710 Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF*  
711 *conference on computer vision and pattern recognition*, pp. 9653–9663, 2022.

712 Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Yixuan Wei, Qi Dai, and Han Hu. On data scaling  
713 in masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
714 *Pattern Recognition*, pp. 10365–10374, 2023.

715  
716 Donald G York, Jennifer Adelman, John E Anderson Jr, Scott F Anderson, James Annis, Neta A  
717 Bahcall, JA Bakken, Robert Barkhouser, Steven Bastian, Eileen Berman, et al. The sloan digital  
718 sky survey: Technical summary. *The Astronomical Journal*, 120(3):1579, 2000.

719 Hu Zhan. The wide-field multiband imaging and slitless spectroscopy survey to be carried out by  
720 the survey space telescope of china manned space program. *Chinese Science Bulletin (Chinese*  
721 *Version)*, 66(11):1290–1298, 2021.

722 Zhirui Zhang, Zhiqiang Zou, Nan Li, and Yanli Chen. Classifying galaxy morphologies with few-shot  
723 learning. *Research in Astronomy and Astrophysics*, 22(5):055002, 2022.

724  
725 Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene  
726 parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and*  
727 *pattern recognition*, pp. 633–641, 2017.

728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A PRE-TRAIN DATASET

The dataset is derived from the January 2021 public release of Dark Energy Spectroscopic Instrument Legacy Imaging Surveys Data Release 9 (DESI DR9<sup>4</sup>). Following the data processing pipeline described in `Astro-76M` (Stein et al., 2022), we first filter stellar objects from the Tractor Catalog<sup>5</sup> by selecting only sources whose best-fit morphological models deviate from the stellar point spread function.

From this preprocessed catalog, we initially extract 1.5 million galaxies, generating  $3 \times 256 \times 256$  pixel cutouts (0.262 arcsec/pixel resolution) centered on each galaxy’s equatorial coordinates using DESI DR9 JPEG imaging data. To augment the dataset, we supplement these with 0.5 million additional cutouts from background regions. Furthermore, we improve dataset quality by removing samples containing edge artifacts where blackened regions exceed 10% of the image area – a process demonstrated to reduce loss spikes (Molybog et al., 2023) and enhance training stability during model pre-training.

Pre-train data critically impact representation quality. The field of natural image is undergoing a shift in pre-train data: transitioning from the classic ImageNet (Deng et al., 2009) to the emerging MetaCLIP datasets (Fan et al., 2025; Bolya et al., 2025). While this work has primarily enhanced pre-train data through a straightforward implementation, future research could explore more sophisticated data curation strategies for higher-quality datasets, potentially even integrating multimodal inputs from multiple observational sources.

## B IMPLEMENTATION DETAILS

**Architectures.** Our C-MAE adopts an asymmetric encoder-decoder architecture. The encoder uses the standard ViT (Dosovitskiy et al., 2020) architectures that have a stack of Transformer blocks (Vaswani et al., 2017). Each block consists of a multi-head self-attention layer and an MLP layer with LayerNorm (Ba et al., 2016). Refer to Table 8 for details about the encoder. The decoder adopts a lightweight design (ablation in Table 9) and is discarded after pre-training. Moreover, we do not use [CLS] token (He et al., 2022) during pre-training, and treat the global average pooling on the image tokens sequence as input for the task head.

**Hyper-parameters.** We search for the learning rate, weight decay, drop path rate, and epochs, for each model size (18, 50, 101, B, L, H), each model type (ResNet, ViT, Swin), and each downstream task. The hyper-parameters are included in Tables 10, 11 and 12.

Table 8: **Encoder architectures.** These variants adopt scale rules in El-Nouby et al. (2024).

architecture	Layers	Patch size	Embedding dim	MLP size	Heads	Parameters
ViT-Base	12	16	768	3,072	12	86M
ViT-Large	24	16	1,024	4,096	16	303M
ViT-Huge	24	14	1,536	6,144	16	680M

## C ADDITIONAL ABLATION AND RESULTS

We ablate some basic components related to C-MAE pre-training in Table 9 and Figure 5, and observe some phenomena different from those in natural images, introduced next. The default setting is 800-epoch pre-training C-MAE on `DESI-1M` with ViT-B as encoder, and reports fine-tuning and linear probing accuracy (%) on `galaxy-desi` test set.

**Decoder design.** The decoder consists of a set of Transformer blocks (Vaswani et al., 2017) and can be flexibly designed in a manner that independent of the encoder. Table 9 studies its depth (number of Transformer blocks) and width (embedding dim). C-MAE needs a shallow decoder (1~4 blocks), both for linear probing and fine-tuning. Increasing the depth will produce very bad results; for example, when the depth is 12, the accuracy is reduced by half in linear probing. This

<sup>4</sup><https://www.legacysurvey.org/dr9>

<sup>5</sup><https://www.legacysurvey.org/dr9/description/#the-tractor-catalogs>

Table 9: **Ablation of decoder design.** We report top-1 accuracy (%) on `galaxy-desi` under the protocol of linear probing and fine-tuning. The default settings are marked by grey.

(a) decoder depth.			(b) decoder width.		
depth	linear probing	fine tuning	width	linear probing	fine tuning
1	<b>53.21</b>	<b>87.23</b>	256	45.19	87.52
2	44.82	87.21	384	45.45	87.65
4	53.10	87.12	<b>512</b>	<b>53.21</b>	87.23
8	32.76	86.05	640	48.61	<b>87.72</b>
12	27.27	85.42	768	47.23	87.69

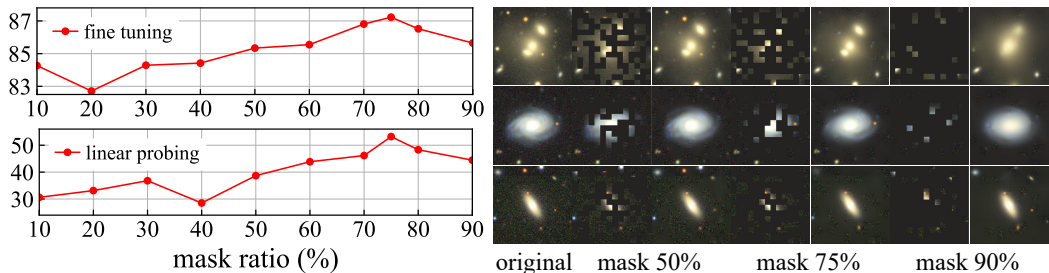


Figure 5: **Mask ratio.** Left: 75% mask ratio works significantly better than others for both fine-tuning (top) and linear probing (bottom). The y-axes are accuracy on `galaxy-desi` test set. Right: example reconstruction results on `galaxy-desi` test images. The C-MAE is pre-trained with a mask ratio of 75% but applied on inputs with other mask ratios. Although different from the ground truth, it is semantically reasonable.

phenomenon is different from observations in natural images, which MAE (He et al., 2022) works with a sufficiently deep decoder. We infer that this difference is caused by the different content complexity of astronomical images and natural images: compared to natural images, astronomical images have simpler contents, where celestial objects manifest low-variability morphological profiles, so it does not require a deep decoder to reconstruct the missing contents. Table 9b varies decoder width. We use 512-d by default, which performs well under linear probing and fine-tuning.

**Mask ratio.** Figure 5 left shows the influence of the mask ratio. Surprisingly, linear probing and fine-tuning exhibit similar trends with varying mask ratios. The optimal mask ratio, 75%, works significantly better than other values. This is different from the observation of MAE in natural images (ImageNet-1K), where a wide range of mask ratios (40~80%) work well for fine-tuning. Figure 5 right shows some reconstruction results using 75% mask ratio during pre-training.

**Training epochs.** Figure 6 varies the training epochs. The fine-tuning accuracy improves steadily with more training epochs, which aligns with MAE in natural image. From another perspective, the fact that more training epochs can give further improvement means that the model converges more slowly. We assume that this is due to the encoder only sees 25% of patches per image, which speeds up training but also loses some meaningful image content.

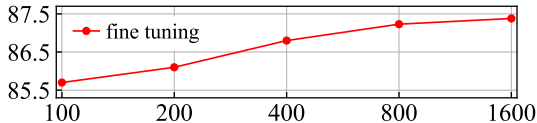


Figure 6: **Training epochs.** Each point is a full training schedule. The y-axis is accuracy on `galaxy-desi` test set.

**Some ineffective tricks.** Beyond standard MAE (He et al., 2022), we also tried some other tricks during pre-training, including using a fixed random patch projection (Chen et al., 2021) layer to embed the image patches, learnable position embedding (Dosovitskiy et al., 2020; Chen et al., 2021), LayerScale (Touvron et al., 2021; 2022), and converting MaskFeat (Wei et al., 2022) into a regularization loss, but did not observe any benefit on performance or training stability.

**Astronomy-specific optimizations** involve two aspects: (1) We enhanced the pre-training data by supplementing it with a richer set of images, and during pre-training, we removed training samples containing artifacts. This process has been demonstrated to reduce loss spikes (Molybog et al., 2023)

and enhance training stability. (2) We redesigned the decoder in an astronomical context. Compared to the optimal configuration for natural images, our redesign achieved a 1.8% accuracy gain.

**More results of natural image pre-training.** We evaluate DINO-v2 (Oquab et al., 2023) and CLIP (Radford et al., 2021) on morphological classification and redshift estimation tasks and show results to the table below, using ViT-B as encoder in all cases. Note that the publicly available DINO-v2 model was pre-trained on LVD-142M, and CLIP on LAION-2B.

pre-train	IN-1K, sup	IN-1K, MAE	LVD-142M, DINO-v2	LAION-2B, CLIP	DESI-1M, C-MAE
accuracy (%)	86.44	82.82	87.15	86.78	87.23
$\Delta z \downarrow$	2.74e-4	13.64e-4	3.59e-4	8.25e-4	0.97e-4

These larger scale natural image models ultimately underperform our DESI-1M. This performance gap demonstrates the necessity of pre-training on astronomical data for astronomical tasks.

## D CODE AND DEMO

**Code.** We include our self-contained codebase (refer to the zip file `Code-CMAE`) as a part of the supplementary material. Please refer to `README.md` for instructions how to use the code. We do not include model weights in the supplementary material as they are too large ( $>100\text{MB}$ ) that exceed the space limit. We will open source pre-train dataset, task datasets, and our code to foster research.

**License.** We release open-source code under the MIT license to foster future research in this field.

**Requirement.** Running our Python code requires some common packages, such as PyTorch, TorchVision, and timm. Please refer to `Code-CMAE/README.md` for more details.

**Demo.** We use Jupyter Notebook to create three demos, including evaluating image classification and detection results and displaying the redshift prediction. See `demo-accuracy-eval.ipynb`, `demo-detection-eval.ipynb`, and `demo-redshift-eval.ipynb` for more details.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

Table 10: Fine tuning C-MAE on different tasks. Multiple values in a cell are for different models.

config	galaxy-desi	neuralens-desi	galaxy-sdss	redshift-sdss
optimizer	AdamW	AdamW	AdamW	AdamW
momentum	$\beta_1, \beta_2 = 0.9, 0.999$	$\beta_1, \beta_2 = 0.9, 0.999$	$\beta_1, \beta_2 = 0.9, 0.999$	$\beta_1, \beta_2 = 0.9, 0.999$
weight decay	0.05, 0.5, 0.5	0.1	0.3, 0.1, 0.05	0.5, 0.1, 0.5
batch size	64, 64, 32	64	64, 64, 32	64, 64, 32
learning rate	1.5e-3, 2e-3, 1e-3	5e-4	1.5e-3, 2e-3, 1e-3	1.5e-3, 1e-3, 1.5e-3
lr schedule	cosine decay	multi step	cosine decay	cosine decay
layer-wise lr decay	0.65, 0.75, 0.75	0.7, 0.9, 0.9	0.65, 0.75, 0.65	0.65, 0.65, 0.75
training epochs	50	32K iters	50, 50, 90	50, 50, 85
warmup epochs	5	500 iters	5	5, 5, 10
drop path	0.1	0.1	0.1	0.1
EMA	0.9999	-	0.9999	0.9999

Table 11: Training ResNet and ViT from scratch on galaxy-desi.

config	ResNet-18/50/101	ViT-Ti/S/B/L
optimizer	SGD	AdamW
optimizer momentum	0.9, w/ Nesterov	$\beta_1, \beta_2 = 0.9, 0.999$
weight decay	5e-3	5e-2 (Ti/S), 1e-2 (B/L)
batch size	128	128
learning rate	1e-3	1e-3 (Ti/S), 1e-4 (B/L)
learning rate schedule	multi step, [160, 230, 290]	cosine decay
training epochs	300	300
warmup epochs	0	10

Table 12: Fine tuning LVM, AstroSSL, and AstroCLIP on galaxy-desi.

config	LVM	AstroSSL	AstroCLIP
	Swin-T	ResNet-50	ViT-L
optimizer	AdamW	SGD	AdamW
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.999$	0.9, w/ Nesterov	$\beta_1, \beta_2 = 0.9, 0.999$
weight decay	1e-3	1e-3	1e-3
batch size	64	128	64
learning rate	1e-4	1e-2	1e-4
learning rate schedule	cosine decay	multi step, [30, 45]	cosine decay
layer-wise lr decay	0.65	-	0.75
training epochs	50	50	50
warmup epochs	5	5	5
drop path	0.1	-	0.1
EMA	0.9999	0.9999	0.9999