

---

# LaFTer: Label-Free Tuning of Zero-shot Classifier using Language and Unlabeled Image Collections

---

## Abstract

1 Recently, large-scale pre-trained Vision and Language (VL) models have set a new  
2 state-of-the-art (SOTA) in zero-shot visual classification enabling open-vocabulary  
3 recognition of potentially unlimited set of categories defined as simple language  
4 prompts. However, despite these great advances, the performance of these zero-  
5 shot classifiers still falls short of the results of dedicated (closed category set)  
6 classifiers trained with supervised fine-tuning. In this paper we show, for the  
7 first time, how to reduce this gap without any labels and without any paired VL  
8 data, using an unlabeled image collection and a set of texts auto-generated using a  
9 Large Language Model (LLM) describing the categories of interest and effectively  
10 substituting labeled visual instances of those categories.

## 11 1 Introduction

12 Vision and Language (VL) models [1–3] recently became the de-facto standard for generalized zero-  
13 shot learning enabling recognition of arbitrary (open) set of categories provided with just their text  
14 descriptions and without requiring any additional data or training. However, this incredible flexibility  
15 comes at a performance cost and all the VL models, including the most widely used CLIP [1], still  
16 require additional supervised training (e.g. tuning its vision encoder) to compete with the closed  
17 set (of target categories) supervised training methods. Naturally, this incurs undesirable adaptation  
18 costs for those, otherwise very versatile and flexible, foundation models. Image annotations are often  
19 expensive to collect, especially for legacy vision systems, such as traffic surveillance, quality control,  
20 or security.

21 In this paper we propose LaFTer, an approach for completely label-free parameter efficient finetuning  
22 of VL models to a set of target classes. Our goal is to substitute the need for expensive-to-obtain  
23 image annotations with unsupervised finetuning of VL models. We show that since the VL models  
24 share a common text-image embedding space (due to their contrastive learning objective), there is  
25 a possibility to train using embeddings of samples from one modality (e.g., auto-labeled text) and  
26 then successfully apply what we trained to classify embeddings of samples from the other modality  
27 (e.g., unlabeled images). More specifically, we show that it is possible to train a neural network  
28 (e.g., a classifier) to classify text instances that can be successfully used to classify images, showing  
29 successful cross-modal transfer. Instead of collecting labeled visual instances, in our label-free  
30 LaFTer approach, we are mining text descriptions of the target categories by prompting an LLM (e.g.,  
31 GPT-3 [4]) and combining them with handcrafted prompts, as shown in Figure 1. After creating such  
32 a text dataset, we train a neural network to classify each text instance (sentence) in it to the source  
33 class label that was used to produce the instance. This text classifier can be readily used to classify  
34 images when used on top of a CLIP visual encoder. Furthermore, we take advantage of this text-only  
35 pre-trained classifier by employing it in a pseudo-labeling pipeline (inspired by FixMatch [5]), to  
36 further finetune the CLIP vision encoder on an unlabeled image collection. To reduce overfitting  
37 and keep the finetuning parameter efficient, we make use of Visual Prompt Tuning [6] combined  
38 with adapting the affine transformations (scale and shift) of the normalization layers in the otherwise  
39 frozen network.

## 40 2 LaFTer

41 CLIP [1] consists of a vision encoder and a text encoder, which project images and texts to a common  
42 embedding space. It has been trained on a very large volume (400M) of image-text pairs to align the  
43 embedding of each image to the embedding of its corresponding text, at the same time pushing it away

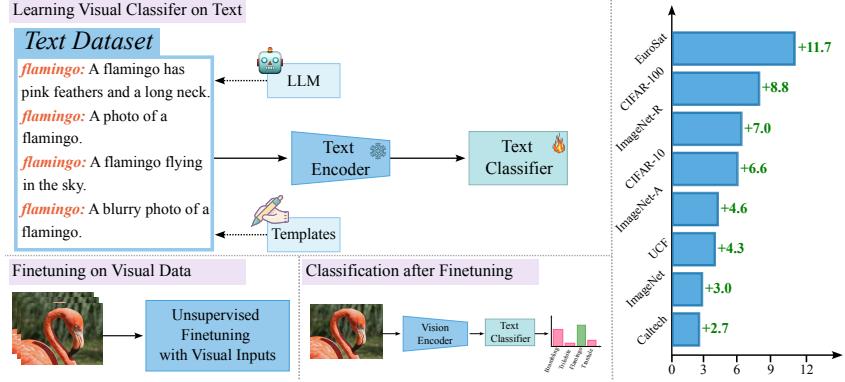


Figure 1: LaFTer proposes to first train a classifier on a natural language text dataset mined in a controlled manner from a set of target classes by generating descriptions for each class label using an LLM and mixing them with handcrafted templates. The training objective is to classify each description to the correct (source) class name (top-left). In the second stage, LaFTer employs the text-only classifier to generate pseudo-labels on the unlabeled data to further finetune the vision encoder in a parameter-efficient manner (bottom-left). Finally, the finetuned visual encoder and text classifier is used for eventual classification (bottom-middle).

44 from the embeddings of unrelated texts corresponding to other images. In [1], it was demonstrated  
 45 that the text and vision encoders enable effective zero-shot image classification. Given the set of class  
 46 names  $C$ , the text encoder  $u$  is used to generate the embedding  $u_c$  of a prompt for each class  $c \in C$ ,  
 47 typically of the form ‘A photo of ...’, completed with the class name. Each test image  $x$  is encoded by  
 48 the vision encoder  $v$  and classified using its cosine similarity  $\cos$  to the text embedding of each class.  
 49 The likelihood of a predicted class  $\hat{c}$  is computed with a softmax,

$$l_{\hat{c}}(x) = \frac{e^{\cos(u_{\hat{c}}, v(x))/\tau}}{\sum_{c \in C} e^{\cos(u_c, v(x))/\tau}}, \quad (1)$$

50 where  $\tau$  denotes the temperature constant.

51 The zero-shot classifier (1) does not require any training data, but is typically outperformed by  
 52 networks trained on data from the target domain. In the next section 2.1, we describe a technique to  
 53 train a visual classifier on purely textual data, conveniently generated by Large Language Models. In  
 54 Section 2.2, we propose an unsupervised training setup that lets us benefit from unlabelled data from  
 55 the target domain, when such data is available. Combining these techniques significantly reduces  
 56 the performance gap to the supervised classifier and is even competitive with few-shot approaches  
 57 despite not using any supervision. These results are provided in the Supplementary.

## 58 2.1 Learning an Image Classifier using Text

59 Aligning the text and image embedding spaces is the key idea underlying Vision Language models. It  
 60 gives CLIP and similar methods, their main advantages: the capacity to be trained on an extremely  
 61 large volume of data available on the Internet and the resulting effectiveness as zero-shot classifiers.  
 62 Recently, Zhang et al. [7] explore yet another advantage of the alignment between text and image  
 63 embedding spaces - it allows diagnosing and rectifying vision models spawn from the VL model  
 64 vision encoder by using the other modality. In particular, it hints that it might be possible to finetune  
 65 a zero-shot image classifier on textual data.

66 While acquiring and annotating images requires manual labor, a training set of annotated texts can  
 67 be constructed automatically, using a Large Language Model, like the GPT-3 [4]. Generating text  
 68 using an LLM constitutes a convenient alternative to text mining, as LLMs represent extremely  
 69 large text corpora on which they were trained. Moreover, prompting LLMs is much more efficient  
 70 than searching a large body of text for the words of interest. To construct our training set, we take  
 71 inspiration from [8] and for each class  $c$ , we prompted GPT-3 with queries of the following kind:  
 72 ‘Describe what a  $[c]$  looks like.’ We repeated the prompting for each class  $c$ , in the dataset and  
 73 complemented the resulting set of synthetic texts with text generated procedurally, using the same  
 74 hand-crafted templates as for constructing prompts for the zero-shot classifier, for example, ‘A photo  
 75 of a  $[c]$ .’. We defer the full list of prompts (to LLMs) and templates to the supplementary material.

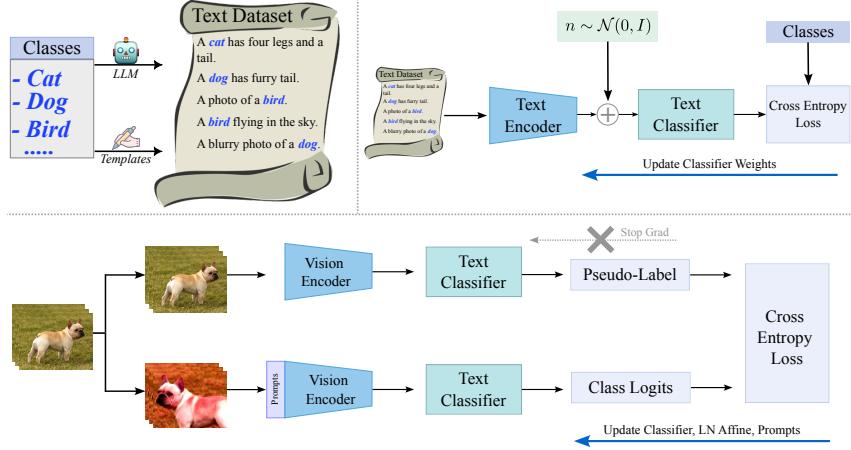


Figure 2: Overview of our LaFTer. (top) Given a set of class labels, we generate a data set of short texts by prompting a Large Language Model (LLM) multiple times with each class name. We compute embeddings of these texts using CLIP text encoder. This lets us train a neural network, the *Text Classifier*, to infer the class used to prompt the LLM from the embedding of the text it generated. Even though the *Text Classifier* has been trained exclusively on text, it performs well in classifying image embeddings generated by CLIP vision encoder. (bottom) We further take advantage of the *Text Classifier* by leveraging it in a pseudo-labeling setup to finetune the VL model.

76 The capacity to finetune the classifier in a supervised text classification setting lifts the architectural  
 77 constraints imposed by the zero-shot setup. More precisely, the class likelihoods no longer need to  
 78 be produced according to (1), but can instead be computed by a trainable classifier  $f$ , of arbitrary  
 79 architecture. We train  $f$  with the Smoothed Cross Entropy loss. To regularize training, we add  
 80 Gaussian Noise  $n \sim \mathcal{N}(0, I)$  to the  $l_2$ -normalized feature vector from the clip text encoder. Formally,  
 81 given a text training set  $T$  consisted of pairs of text fragments  $t$  and class labels  $c$ , the training  
 82 objective is:

$$\min_{\theta} \sum_{\substack{(t,c) \in T \\ n \sim \mathcal{N}(0, I)}} \mathcal{L}_{\text{SCE}}\left(f_{\theta}\left(\frac{u(t)}{\|u(t)\|} + n\right), c\right), \quad (2)$$

83 where  $\theta$  is the parameter of the classifier. Training of the text-only classifier is very efficient. For  
 84 example, 3000 epochs of training the classifier on the data set of 130000 text sentences, representing  
 85 the 1000 classes of the ImageNet [9] dataset is completed in  $\sim 120$  seconds on an NVIDIA 3090  
 86 graphics card. As demonstrated in Supplementary, training the classifier on our dataset of synthetic  
 87 texts yields a network that matches or outperforms the zero-shot classifier, even though the classifier  
 88 is initialized randomly.

## 89 2.2 Unsupervised Finetuning on Target Domain Images

90 Text-only training does not require any image data. However, in many applications, unlabeled  
 91 images from the target domain are readily available, or can be acquired at a low cost. Given a set of  
 92 unlabeled images we propose to take advantage of the text-only pre-trained classifier and use it in a  
 93 pseudo-labeling pipeline on top of the vision encoder as demonstrated at the bottom part of Figure 2.

94 Inspired by Fixmatch [5], for each training image  $x$ , we generate two views: the weakly-augmented  
 95 view  $\alpha_w(x)$  and the strongly-augmented view  $\alpha_h(x)$ , where  $\alpha$  denotes a stochastic augmentation  
 96 function. Contrary to Fixmatch, in our unsupervised finetuning we set  $\alpha_w$  as an identity transfor-  
 97 mation. For  $\alpha_h$  we use the augmentations proposed in [10]. The weakly-augmented view serves to  
 98 generate a pseudo-label. To that end, it is passed through the vision encoder  $v$  and the text classifier  
 99  $f$ , yielding a vector of class probabilities  $p$ . Class probabilities give rise to pseudo labels, that we  
 100 denote by  $\hat{c}(x)$ . Formally,

$$\hat{c}(x) = \arg \max_{c \in C} p_c, \quad \text{where} \quad p = f(v(\alpha_w(x))). \quad (3)$$

101  $\hat{c}(x)$  is not differentiable, and we do not backpropagate the error signal through it during training. The  
 102 pseudo labels  $p_{\hat{c}}$ , are used as ground truth for training the network on the heavily-augmented views.  
 103 Since the vision encoders in the two branches share the weights, the pseudo-labels are generated in  
 104 an *online* manner and are constantly refined as adaptation is progressing.

105 When processing the heavily augmented views  $\alpha_h(x)$ , we augment the vision encoder by visual  
 106 prompt tuning [6]. That is, we append randomly initialized, learnable parameters (a matrix of size  
 107  $N_P \times d_v$  where  $N_P$  is the number of prompts and  $d_v$  is the channel dimension of the vision encoder)  
 108 to the input of the vision transformer after the initial embedding layer. This helps the network account  
 109 for the heavy augmentation of the input and, as we show in Supplementary, boosts the performance  
 110 of the finetuned classifier. We denote the vision encoder with prompting by  $v^p$ . The prediction for  
 111 the heavily-augmented image is obtained as  $f(v^p(\alpha_h(x)))$ .

112 The network is trained with the smoothed cross entropy loss  $\mathcal{L}_{SCE}$ . We denote the set of unlabelled  
 113 target domain images by  $D$ , and formalize the training objective as

$$\min_{\theta, \eta} \sum_{x \in D} \mathcal{L}_{SCE}\left(f_{\theta}\left(v_{\eta}^p(\alpha_h(x))\right), \hat{c}(x)\right), \quad (4)$$

114 where  $\theta$  and  $\eta$  denote the finetuned parameters of the classifier and of the vision encoder, respectively.  
 115 The parameters  $\eta$  are the visual prompts and the scale and shift (affine) parameters of the normalization  
 116 layers of the vision encoder. The selection of  $\eta$  is motivated by keeping the adaptation *parameter-*  
 117 *efficient*. For LaFTer, the number of trainable parameters are less than 0.4% of the entire model  
 118 parameters, making adaptation extremely lightweight.

### 119 3 Conclusion

120 We propose a completely label-free finetuning method for Vision-Language models by first showing  
 121 cross-modality transfer and learning a classifier on natural language inputs which can successfully  
 122 classify visual data. Later, we leverage this text-only pre-trained classifier in our pseudo-labeling  
 123 pipeline to further finetune the VL models in a parameter efficient manner. We extensively evaluate  
 124 our LaFTer and achieve state-of-the-art results when comparing to other methods in an unsupervised  
 125 finetuning paradigm, while also performing favorably in comparison to methods relying on few-shot  
 126 supervised learning routines. We refer to the supplementary material for related work, detailed  
 127 experimental results and discussion.

### 128 References

- 129 [1] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,  
 130 P. Mishkin, J. Clark, *et al.*, “Learning Transferable Visual Models from Natural Language  
 131 Supervision,” in *Proc. ICML*, 2021.
- 132 [2] Y. Li, F. Liang, L. Zhao, Y. Cui, W. Ouyang, J. Shao, F. Yu, and J. Yan, “Supervision  
 133 Exists Everywhere: A Data Efficient Contrastive Language-Image Pre-training Paradigm,”  
 134 *arXiv:2110.05208*, 2021.
- 135 [3] N. Mu, A. Kirillov, D. Wagner, and S. Xie, “Slip: Self-supervision Meets Language-image  
 136 Pre-training,” *arXiv:2112.12750*, 2021.
- 137 [4] T. B. Brown *et al.*, “Language Models are Few-Shot Learners,” *arXiv:2005.14165*, 2020.
- 138 [5] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Ku-  
 139 rakin, and C.-L. Li, “Fixmatch: Simplifying Semi-supervised Learning with Consistency and  
 140 Confidence,” in *NeurIPS*, 2020.
- 141 [6] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim, “Visual  
 142 Prompt Tuning,” in *Proc. ECCV*, 2022.
- 143 [7] Y. Zhang, J. Z. HaoChen, S.-C. Huang, K.-C. Wang, J. Zou, and S. Yeung, “DrML: Diagnosing  
 144 and Rectifying Vision Models using Language,” in *Proc. ICLR*, 2023. [Online]. Available:  
 145 <https://openreview.net/forum?id=losu6IAaPeB>.
- 146 [8] S. Pratt, R. Liu, and A. Farhadi, “What does a Platypus Look Like? Generating Customized  
 147 Prompts for Zero-shot Image Classification,” *arXiv:2209.03320*, 2022.
- 148 [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A Large-scale  
 149 Hierarchical Image Database,” in *Proc. CVPR*, 2009.
- 150 [10] X. Chen and K. He, “Exploring Simple Siamese Representation Learning,” in *Proc. CVPR*,  
 151 2021.