



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

Neurips 2023
Paper Reading by Yixuan Zhang



- 作者介绍
- 研究背景
- 本文方法
- 实验效果
- 总结反思



作者介绍

3

M. Jehanzeb Mirza^{†1}

Leonid Karlinsky²

Wei Lin¹

Mateusz Kozinski¹

Horst Possegger¹

Rogerio Feris²

Horst Bischof¹

¹Institute for Computer Graphics and Vision, TU Graz, Austria.

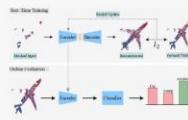
²MIT-IBM Watson AI Lab, USA.



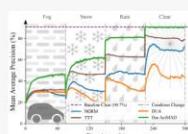
TAP: Targeted Prompting for Task Adaptive Generation of Textual Training Instances for Visual Classification
M. Jehanzeb Mirza, Leonid Karlinsky, Wei Lin, Horst Possegger, Rogerio Feris, Horst Bischof
ArXiv pre-print (Under Review)
[Paper]



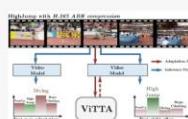
LaFTer: Label-Free Tuning of Zero-shot Classifier using Language and Unlabeled Image Collections
M. Jehanzeb Mirza, Leonid Karlinsky, Wei Lin, Mateusz Kozinski, Horst Possegger, Rogerio Feris, Horst Bischof
NeurIPS 2023
[Paper]



MATE: Masked Autoencoders are Online 3D Test-Time Learners
*M. Jehanzeb Mirza, *Inkyu Shin, *Wei Lin, Andreas Schriebl, Kunyang Sun, Jaesung Choe,
Mateusz Kozinski, Horst Possegger, In So Kweon, Kun-Jin Yoon, Horst Bischof (*Equal Contribution)
ICCV 2023
[Paper | Code]



ActMAD: Activation Matching to Align Distributions for Test-Time-Training
M. Jehanzeb Mirza, Pol Jane Soneira, Wei Lin, Mateusz Kozinski, Horst Possegger, Horst Bischof
CVPR 2023
[Paper | Project Page | Code]



Video Test-Time Adaptation for Action Recognition
*Wei Lin, *M. Jehanzeb Mirza, Mateusz Kozinski, Horst Possegger, Hilde Kuehne, Horst Bischof (*Equal Contribution)
CVPR 2023
[Paper | Code]

智能多媒体内容计算实验室

Intelligent Multimedia Content Computing Lab



- 作者介绍
- 研究背景
- 本文方法
- 实验效果
- 总结反思



研究背景

5

- Large-scale Vision&Language (VL) Models
 - Zero-shot downstream tasks
 - CLIP, ALIGN,...
 - Improve image-text representation alignment
 - off-the-shelf object detectors
 - cross-attention and additional objective functions ITM / MIM
 - finer-grained interactions, modern Hopfield networks, optimal transport distillation, cycle consistency , hierarchical feature alignment

研究背景

6

□ 无监督微调VLM

■ Unsupervised prompt learning (UPL)

- 构造伪标签，learnable prompt
- <https://arxiv.org/abs/2204.03649>

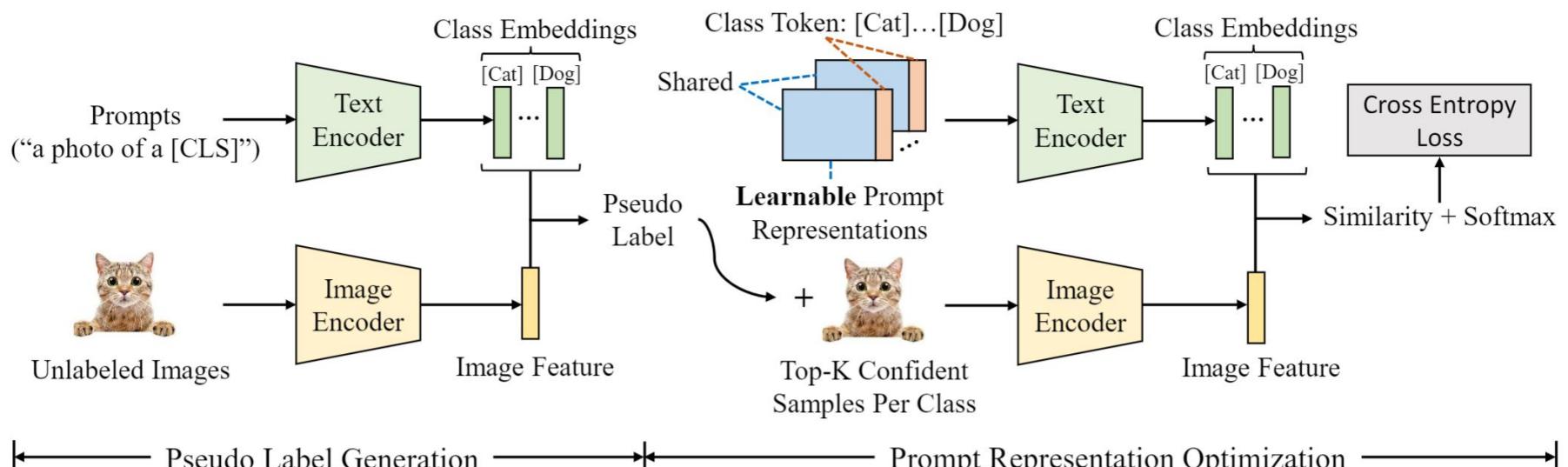


Figure 2: Overview of the proposed unsupervised prompt learning (UPL) framework. Our UPL mainly contains two parts, namely pseudo label generation and prompt representation optimization. We first use CLIP with a simple prompt (e.g., “a photo of a [CLS]”) to generate pseudo labels for target datasets and select top- K confident samples per class for subsequent training. Then we define a learnable prompt representation which is optimized on selected pseudo-labeled samples. For inference, we simply swap out the hand-crafted prompts with the well-optimized prompt representations.



研究背景

7

□ 无监督微调VLM

■ Improving Zero-Shot Models with Label Distribution Priors (CLIP-PR)

- 伪标签 + 分布约束

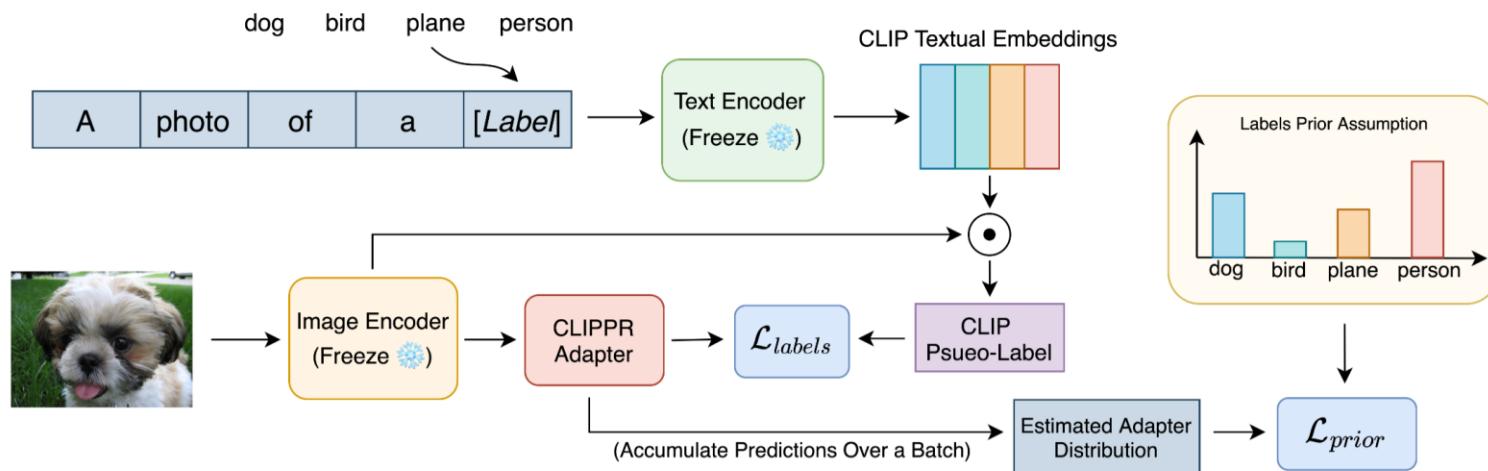


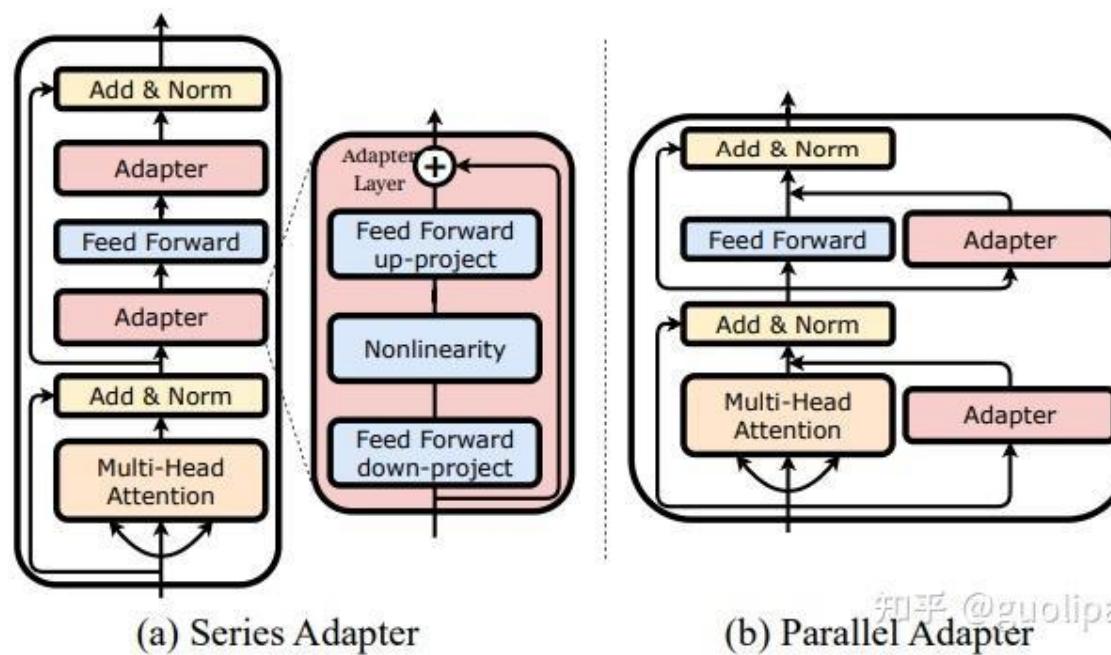
Figure 1. An illustration of our method. Our method trains an **adapter module** on top of a frozen V&L model image encoder. The adapter is trained with two competing objectives: (i) Predicting labels close to the original V&L model zero-shot predictions. (ii) Predicting a labels distribution similar to the given prior distribution. Together, these two objectives adapt the original zero-shot predictions to the distributional prior, resulting in better performance.

研究背景

8

□ 预训练大模型的微调方法

- Fine-tuning
- Parameter efficient fine-tuning (PEFT)
 - Adapter



研究背景

9

□ 预训练大模型的微调方法

- Fine-tuning
- Parameter efficient fine-tuning (PEFT)
 - Adapter
 - Low-Rank Adaption(LoRA)

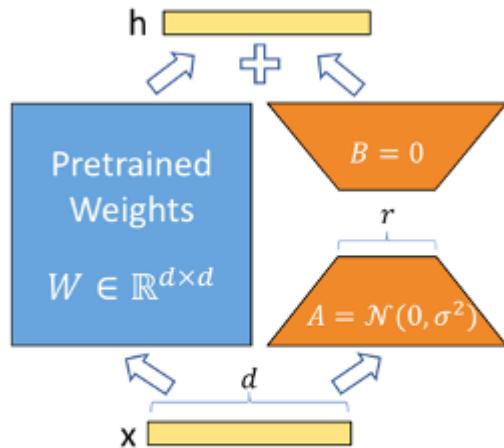


Figure 1: Our reparametrization. We only train A and B .

$$h = W_0x + \Delta Wx = W_0x + BAx$$

■ 模型是过参数化的，依赖于这个低的内在维度 (low intrinsic dimension) 去做任务适配。低秩自适应 (LoRA) 方法对大模型进行微调。

- Adapter由于增加了模型的深度从而额外增加了模型推理的延时
- Prompt较难训练，同时减少了模型的可用序列长度

研究背景

10

□ 预训练大模型的微调方法

- Fine-tuning
- Parameter efficient fine-tuning (PEFT)
 - Adapter
 - LoRA
 - Prompt tuning : text prompting

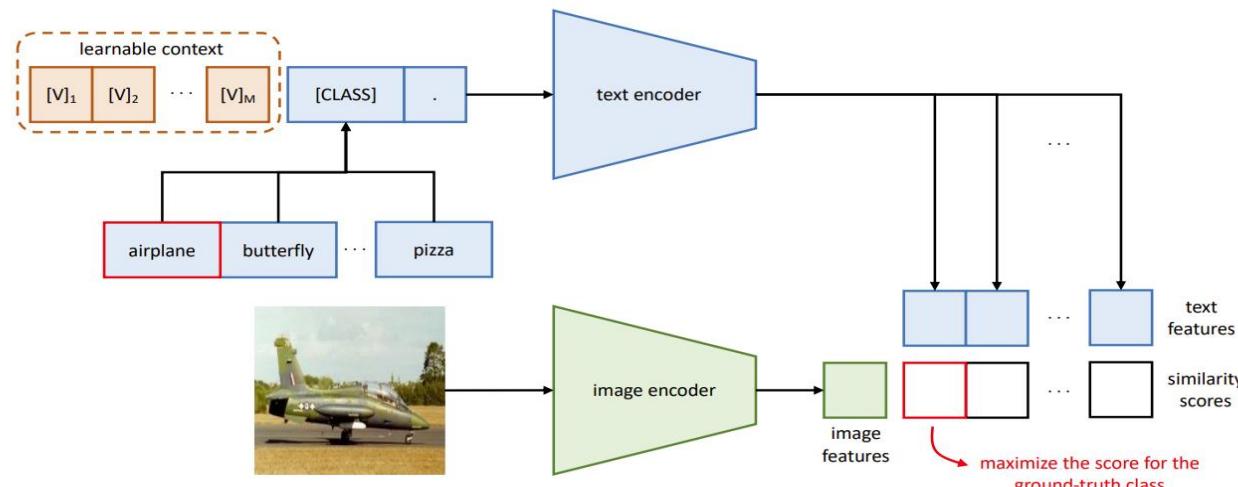


Fig. 2 Overview of Context Optimization (CoOp). The main idea is to model a prompt's context using a set of learnable vectors, which can be optimized through minimizing the classification loss. Two designs are proposed: one is unified context, which shares the same context vectors with all classes; and the other is class-specific context, which learns for each class a specific set of context vectors.

研究背景

11

□ 预训练大模型的微调方法

- Fine-tuning
- Parameter efficient fine-tuning (PEFT)
 - Adapter
 - LoRA
 - Prompt tuning : text prompting

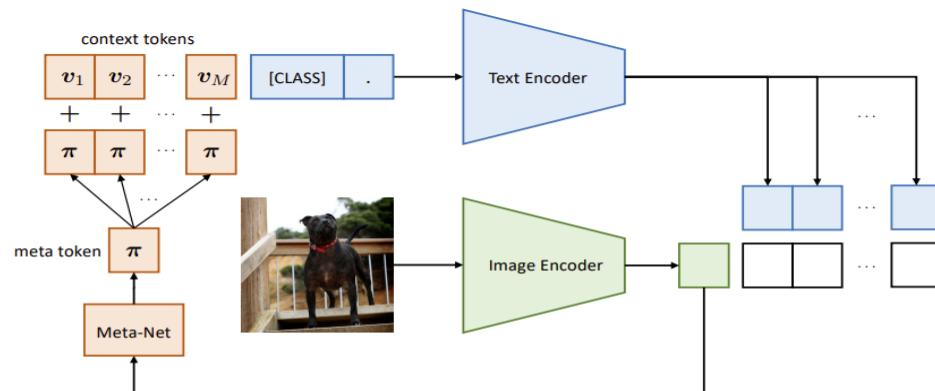


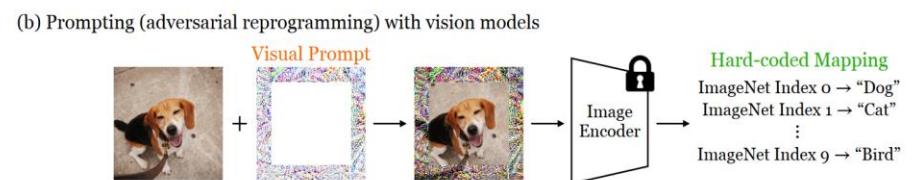
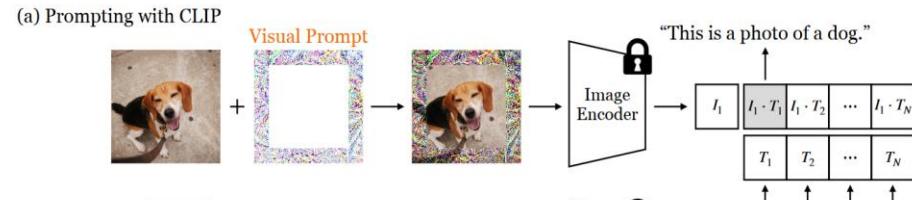
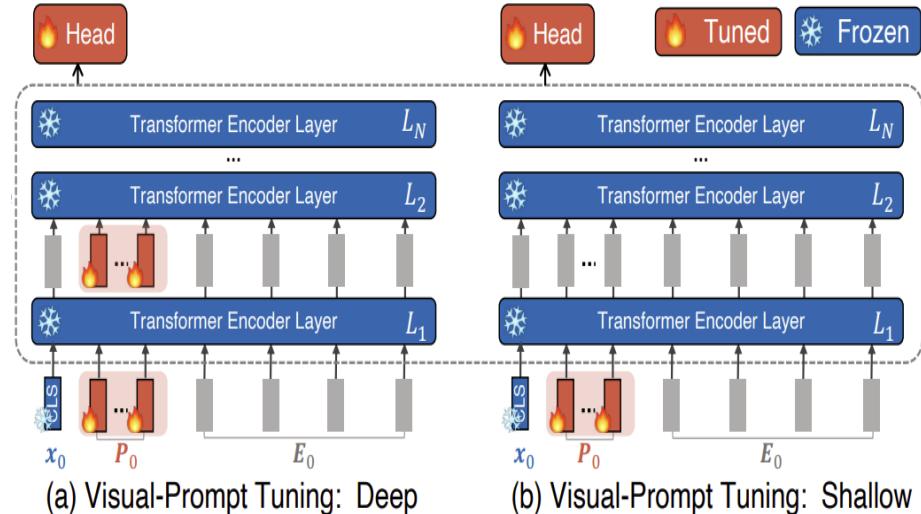
Figure 2. Our approach, Conditional Context Optimization (Co-CoOp), consists of two learnable components: a set of context vectors and a lightweight neural network (Meta-Net) that generates for each image an input-conditional token.

研究背景

12

□ 预训练大模型的微调方法

- Fine-tuning
- Parameter efficient fine-tuning (PEFT)
 - Adapter
 - LoRA
 - Prompt tuning : visual prompting



研究背景

13

- 半监督学习
- Fixmatch

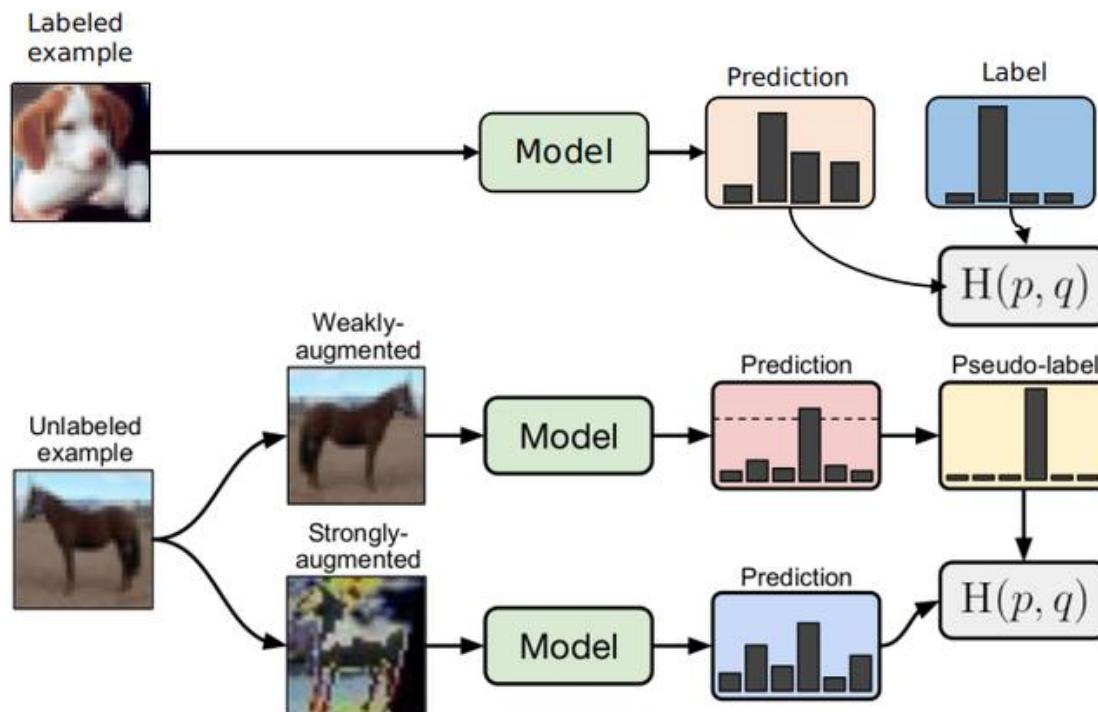


Image
Embedding



Text
Embedding

Western gull
in lake.

Cardinal in
ocean.

Fish crow
in forest.

Gadwall crow
in bamboo.

Training

Waterbird or
Landbird?



Classifier

Diagnosing &
Rectifying

可以通过文本数据上微调零样本图像分类器

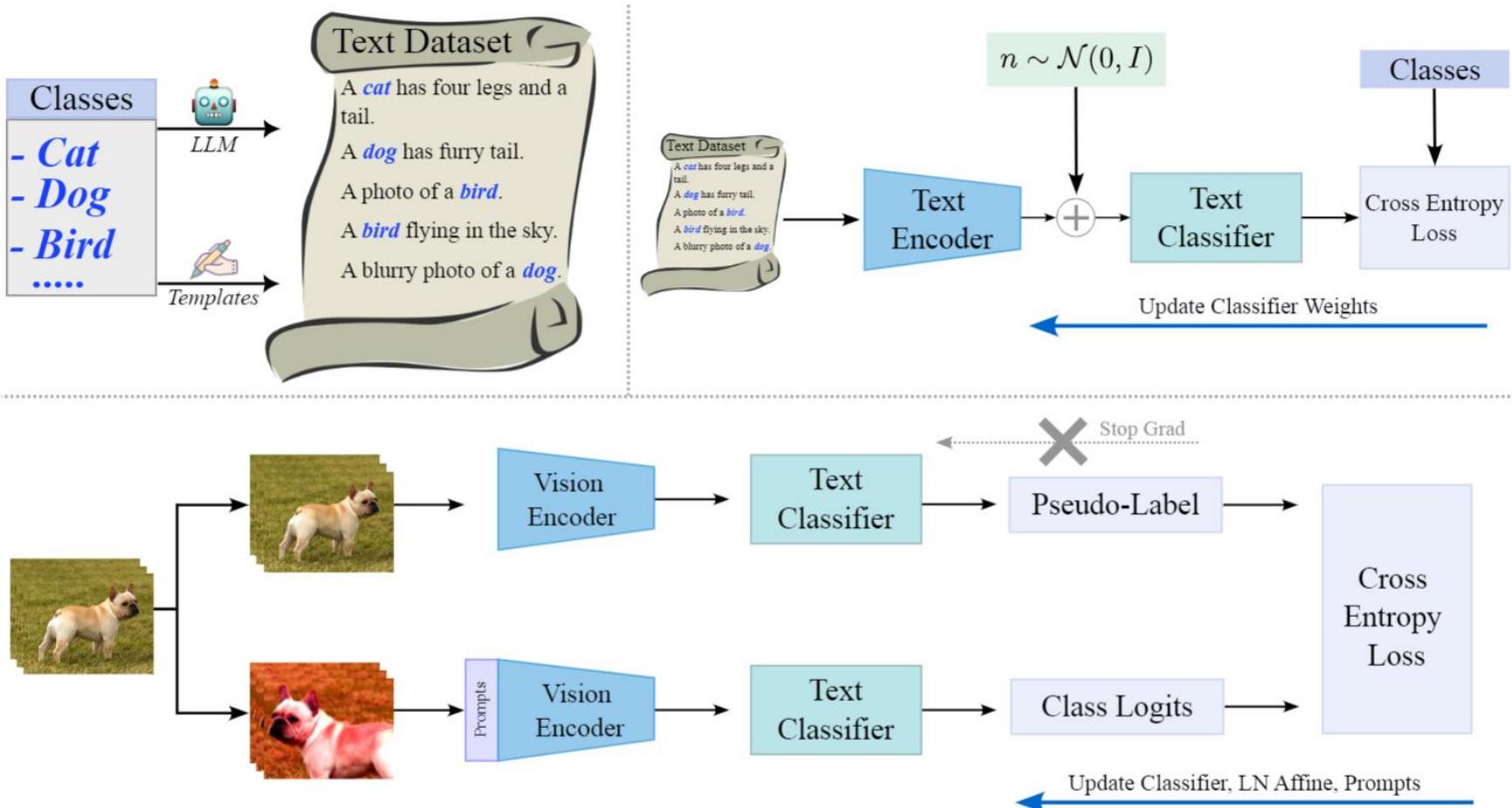


- 作者介绍
- 研究背景
- 本文方法
- 实验效果
- 总结反思



本文方法

16





本文方法

17

□ step1: LLM构造数据集

◎ 采用的LLM模型: GPT3

◎ 50 descriptions for each category: GPT3生成+clip模板

- Describe what a **category** looks like.
- How can you identify a **category**?
- What does a **category** look like?
- Describe an image from the internet of a **category**.
- A caption of an image of a **category**?
- A **quail** is a small game bird with a rounded body and a small head.
- A **quail** is a small, plump bird with a round body and a short tail.
- A **quail** can be identified by its plump body, short legs, and small head with a pointed beak.
- A **quail** can be identified by its small, rounded body and short tail.
- A **quail** looks like a small chicken.
- A **quail** is a small, crested game bird.
- This image is of a **quail** in a natural setting.
- In the image, there is a brown and white **quail** perched on a branch.
- A **quail** hiding in some foliage.
- A young **quail** pecks at the ground in search of food.
- a bad photo of the **category**.
- a **category** in a video game.
- a origami **category**.
- a photo of the small **category**.
- art of the **category**.
- a photo of the large **category**.
- itap of a **category**.

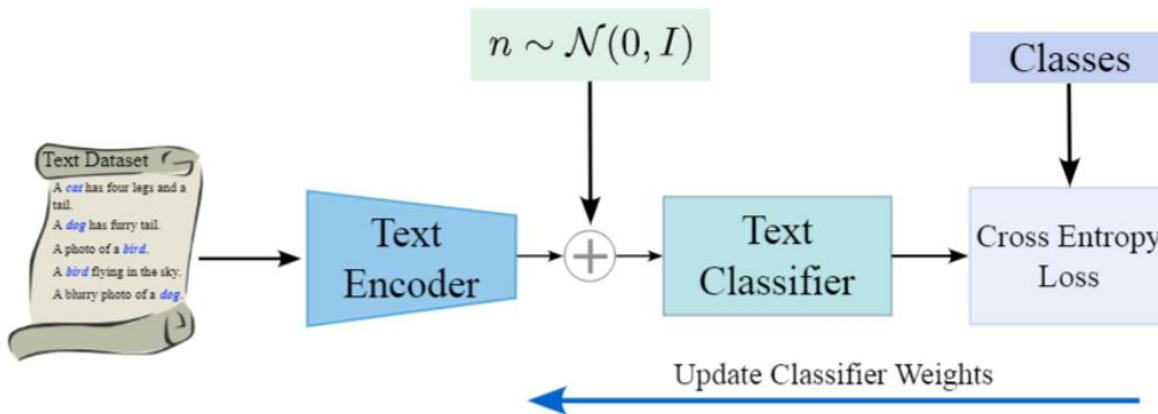
本文方法

18

□ step2: 根据文本数据训练分类器

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

where θ is the parameter of the classifier. Training of the text-only classifier is very efficient. For example, 3000 epochs of training the classifier on the data set of 130000 text sentences, representing the 1000 classes of the ImageNet [53] dataset is completed in ~ 120 seconds on an NVIDIA 3090 graphics card.



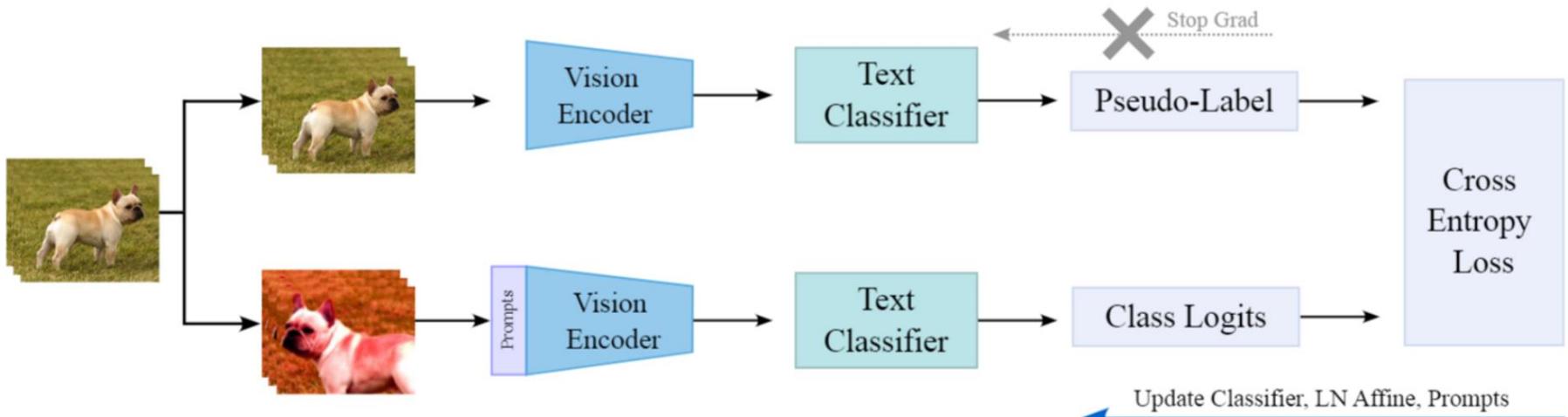
本文方法

19

□ Step3: Unsupervised Finetuning on Target Domain Images

◎ 伪标签生成

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



The network is trained with the smoothed cross entropy loss \mathcal{L}_{SCE} [52]. We denote the set of unlabelled target domain images by D , and formalize the training objective as

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

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



- 作者介绍
- 研究背景
- 本文方法
- 实验效果
- 总结反思



消融实验

25

	Clip	w/o Aug	w/o Prompts	w/o Affine	w/o Cls	w/o Stop Grad	LaFTer
CIFAR-10	88.8	93.5	92.5	94.6	94.1	94.7	95.8
CIFAR-100	64.2	72.6	71.7	72.9	67.4	73.2	74.6
UCF-101	61.0	65.3	66.1	63.6	63.4	67.5	68.2
EuroSat	45.1	63.2	61.2	64.2	60.9	69.2	73.9
ImageNet-A	29.6	29.9	29.8	30.8	30.1	31.1	31.5
ImageNet-S	40.6	40.3	40.7	41.1	40.9	42.1	42.7
ImageNet-R	65.8	68.0	67.7	66.8	66.1	71.8	72.6
Average	56.4	61.8	61.4	62.0	60.4	64.2	65.6

Table 4: Top-1 Accuracy (%) for our LaFTer while ablating the various critical design choices in our methodology. For each of these experiments, we disable one component from our framework and test the resulting method. Aug: Augmentations, Cls: Classifier.



- 作者介绍
- 研究背景
- 本文方法
- 实验效果
- 总结反思



总结反思

27

- 总结
 - ◎ 更好利用图文匹配的方法 --- 文本端优化、视觉端受益
- 思考
 - ◎ 任务设定

总结反思

28

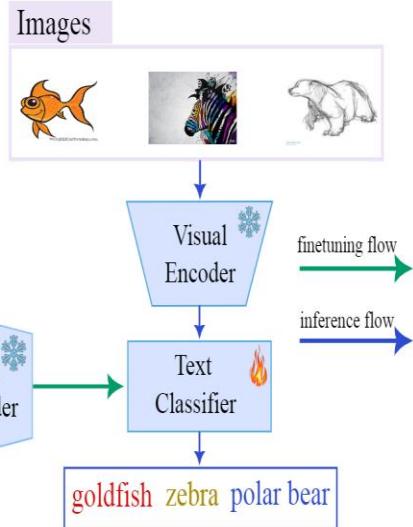
TAP: Targeted Prompts

- Describe what a *cartoon* rendition of a **goldfish** looks like.
- Describe what a *deviantart* rendition of a **zebra** looks like
- Describe what a *sketch* rendition of a **polar bear** looks like

LLM

Targeted Descriptions

- A *cartoon* rendition of a **goldfish** is typically illustrated as a bright yellow, orange, or white fish with large eyes and a wide smile.
- A *deviantart* rendition of a **zebra** usually takes the form of an exaggerated, cartoon-like version of the animal.
- A *sketch* rendition of a **polar bear** might be a simple drawing of a white bear with small black eyes, small ears, a large snout, and a long, shaggy coat.

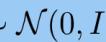


Text-Only Pre-Training

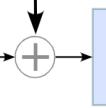
Targeted Class Descriptions

$$n \sim \mathcal{N}(0, I)$$

Text Encoder



Text Classifier



+

Cross Entropy Loss

Update Classifier Weights

Zero-Shot Classification



Image Encoder

Text Classifier

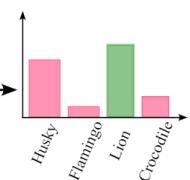


Fig. 2: The targeted descriptions generated by the LLM are automatically matched with the class names to generate the text dataset. This dataset is used to train a text classifier in a *supervised* setting (top). The trained text classifier is used to classify the visual data (bottom).



Thanks for Attention!