



# Visual Tuning

曹耘宁  
2023/6/8



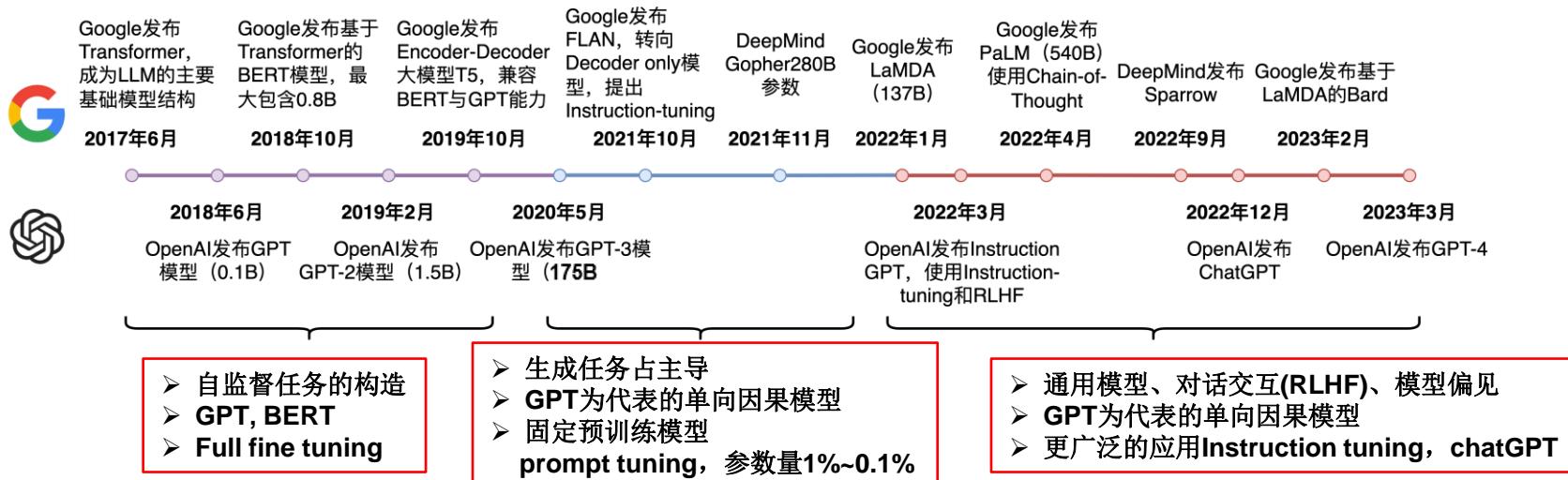
- 研究背景
- Prompt tuning
- Adapter tuning
- Parameter tuning
- 总结



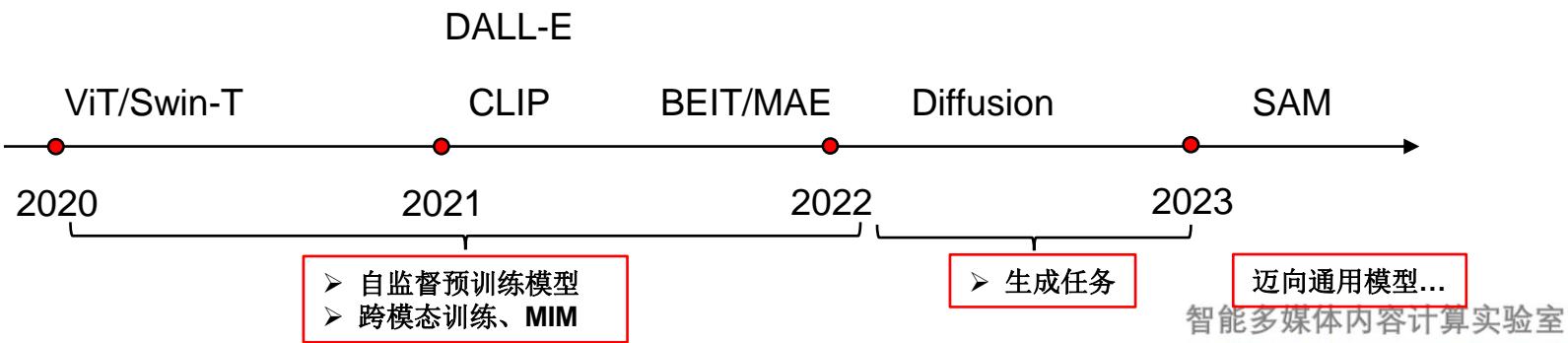
# 大模型的发展

3

## □ 大语言模型LLM



## □ 视觉大模型/视觉语言模型VLM



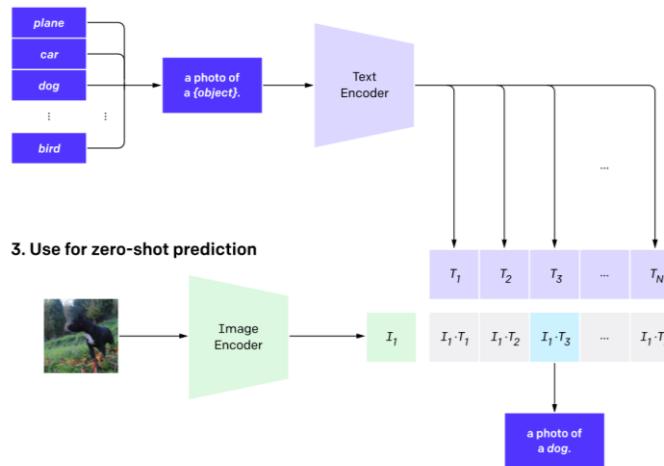


# 视觉语言模型

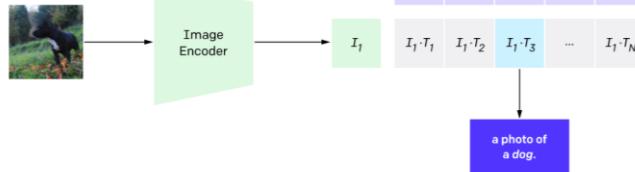
4

- CLIP, 视觉语言模型
- 后续工作致力于将CLIP微调到下游任务

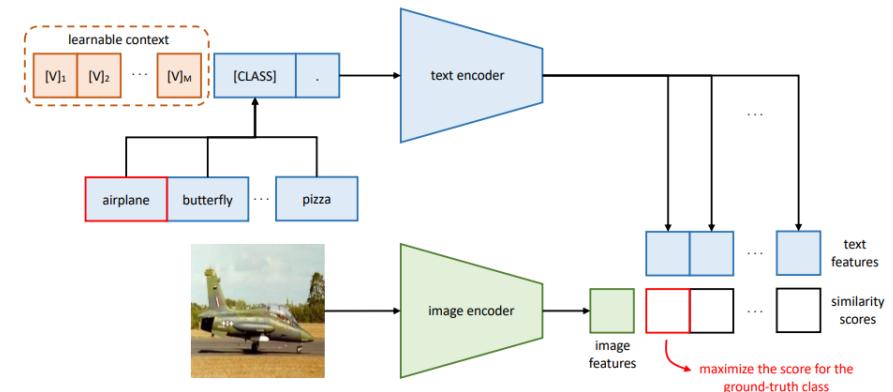
## 2. Create dataset classifier from label text



## 3. Use for zero-shot prediction



CLIP



**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.

CoOp



# 视觉大模型

5

- SAM设计了适用于分割任务的多种prompt形式，实现交互式分割、多模态分割
- 近期工作使用SAM生成伪标签、对SAM做adapter tuning

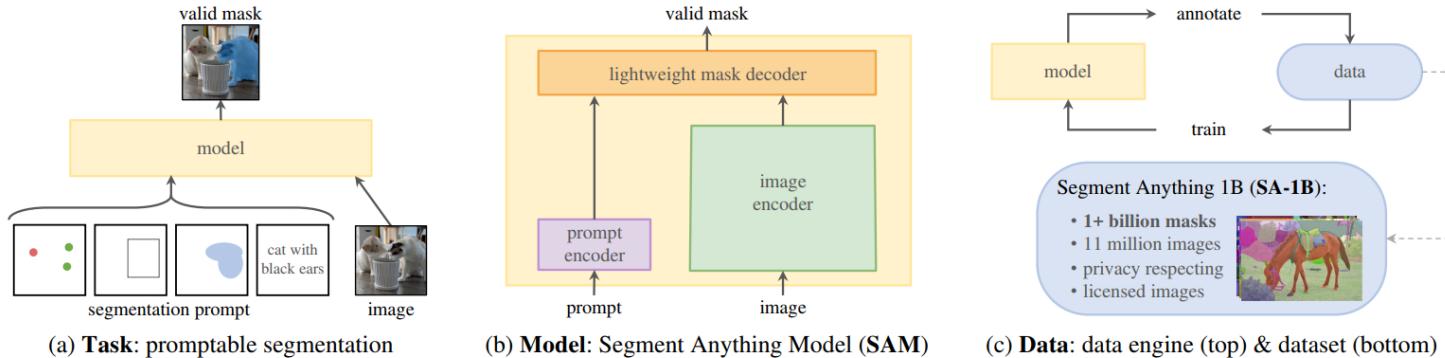


Fig. 1: Overview of the SA project, including task, model, and data. The figure is borrowed from the original paper [20].

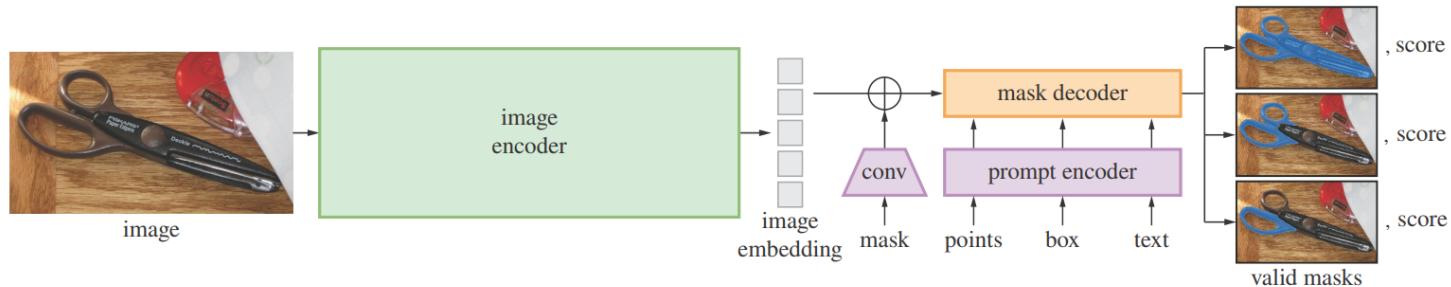


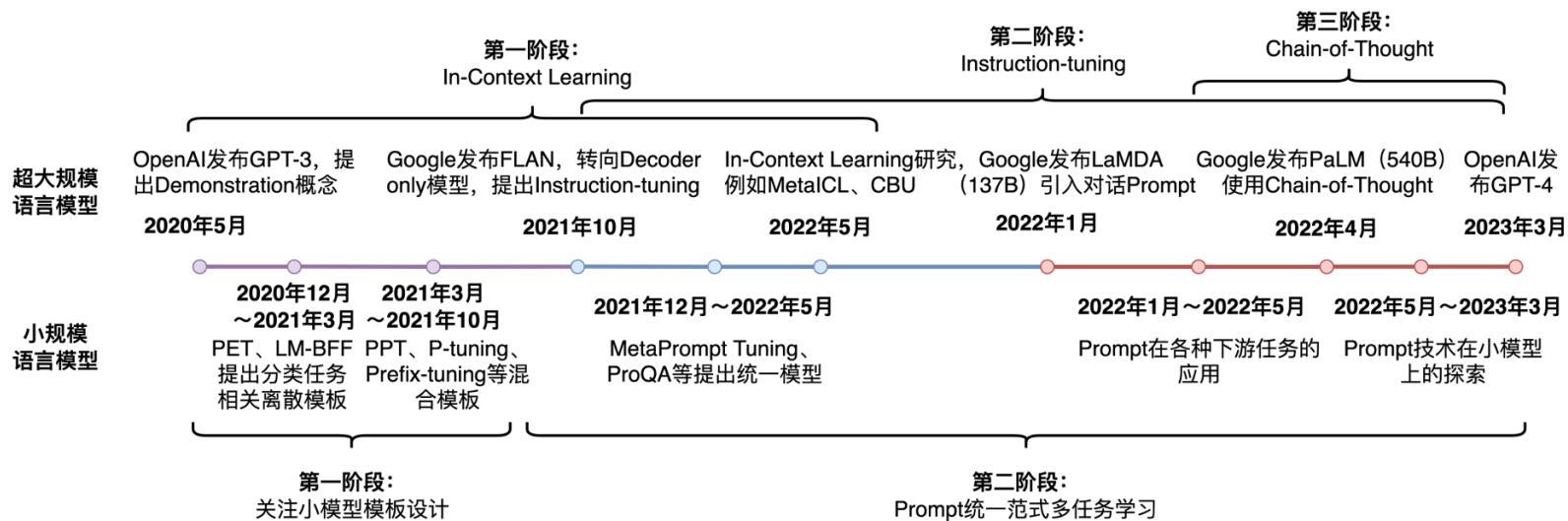
Fig. 2: Overall structure of SAM from the original paper [20].



# 语言模型微调技术

6

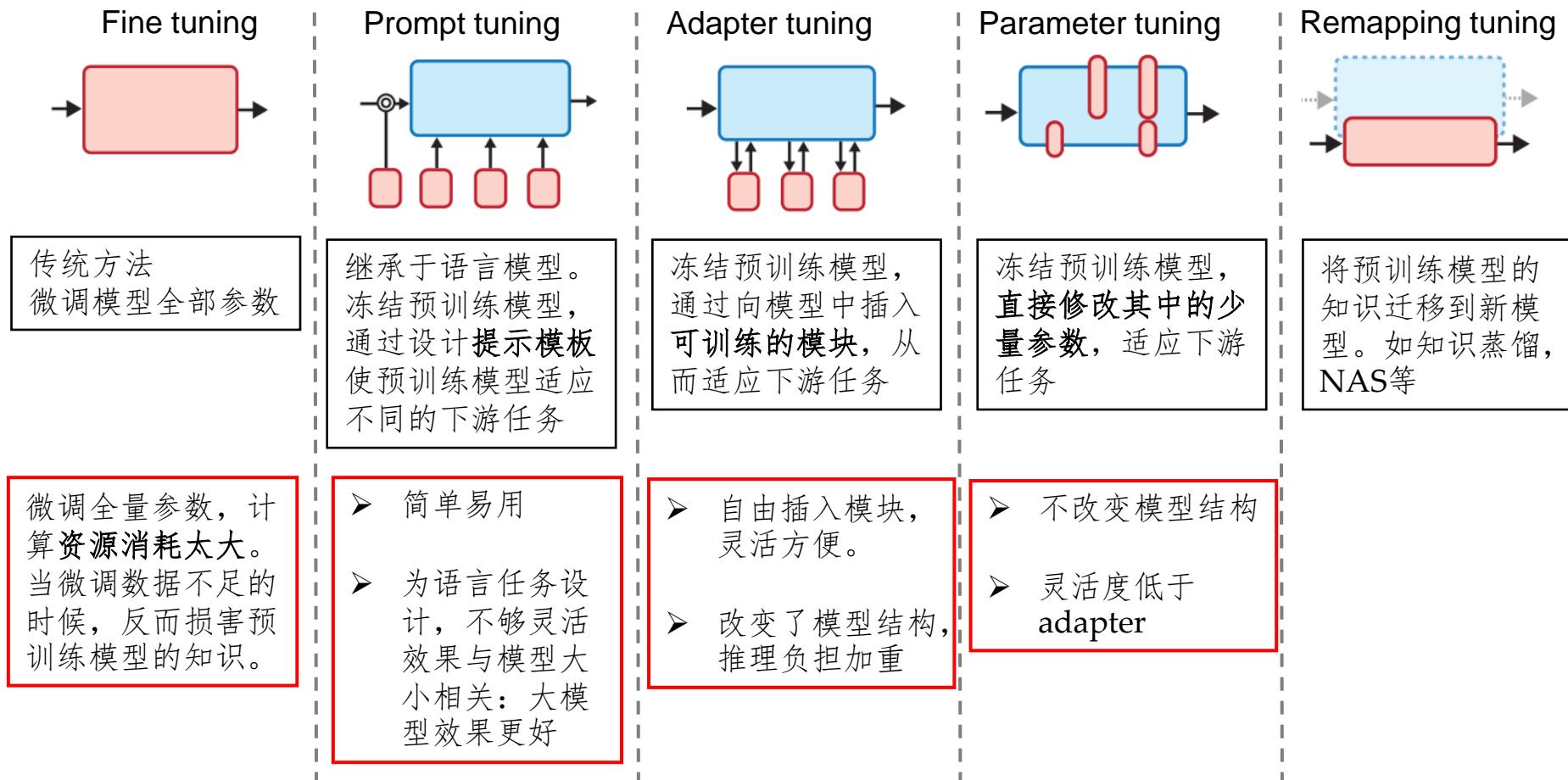
- 手工设计prompt 模板
- Prompt tuning





# 视觉微调技术分类

7





- 研究背景
- Prompt tuning
  - Visual
  - Language
  - Multi-model
- Adapter tuning
- Parameter tuning
- 总结

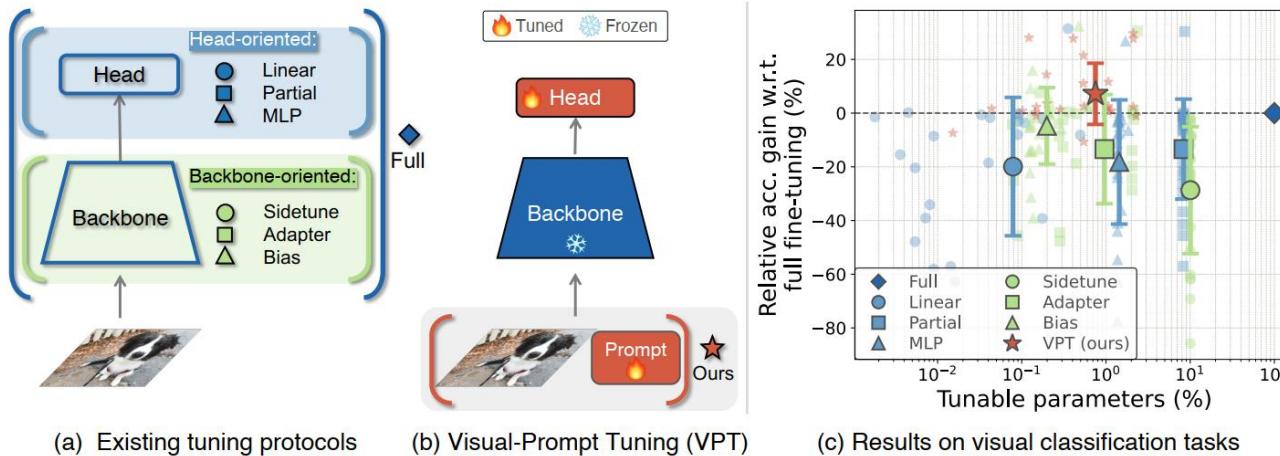
# Visual-driven

9

Menglin Jia<sup>\*1,2</sup>, Luming Tang<sup>\*1</sup>  
 Bor-Chun Chen<sup>2</sup>, Claire Cardie<sup>1</sup>, Serge Belongie<sup>3</sup>  
 Bharath Hariharan<sup>1</sup>, and Ser-Nam Lim<sup>2</sup>

<sup>1</sup>Cornell University<sup>2</sup>Meta AI<sup>3</sup>University of Copenhagen

- 早期工作，添加prompt把预训练ViT模型适应到不同数据集

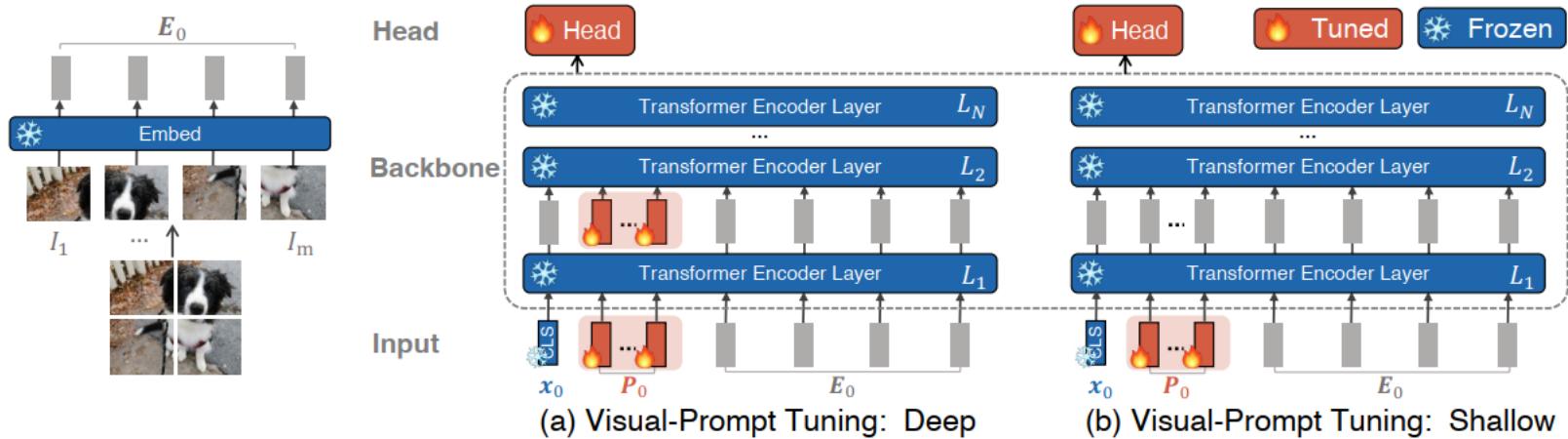


**Fig. 1.** Visual-Prompt Tuning (VPT) *vs.* other transfer learning methods. (a) Current transfer learning protocols are grouped based on the tuning scope: Full fine-tuning, Head-oriented, and Backbone-oriented approaches. (b) VPT instead adds extra parameters in the input space. (c) Performance of different methods on a wide range of downstream classification tasks adapting a pre-trained ViT-B backbone, with mean and standard deviation annotated. VPT outperforms Full fine-tuning 20 out of 24 cases while using less than 1% of all model parameters

# Visual-driven

10

## □ 具体网络结构



**Fig. 2.** Overview of our proposed Visual-Prompt Tuning. We explore two variants: (a) prepend a set of learnable parameters to each Transformer encoder layer's input (VPT-DEEP); (b) only insert the prompt parameters to the first layer's input (VPT-SHALLOW). During training on downstream tasks, only the parameters of prompts and linear head are updated while the whole Transformer encoder is frozen.



# Visual-driven

11

## □ 实验效果

ViT-B/16 (85.8M)		Total params	Scope		Extra params	FGVC		VTAB-1k		
			Input	Backbone		5	7	Natural	Specialized	Structured
	Total # of tasks					5	7	4	4	8
(a)	FULL	24.02×		✓		88.54	75.88	83.36	47.64	
	LINEAR	1.02×				79.32 (0)	68.93 (1)	77.16 (1)	26.84 (0)	
(b)	PARTIAL-1	3.00×				82.63 (0)	69.44 (2)	78.53 (0)	34.17 (0)	
	MLP-3	1.35×			✓	79.80 (0)	67.80 (2)	72.83 (0)	30.62 (0)	
	SIDETUNE	3.69×		✓	✓	78.35 (0)	58.21 (0)	68.12 (0)	23.41 (0)	
(c)	BIAS	1.05×		✓		88.41 (3)	73.30 (3)	78.25 (0)	44.09 (2)	
	ADAPTER	1.23×		✓	✓	85.66 (2)	70.39 (4)	77.11 (0)	33.43 (0)	
(ours)	VPT-SHALLOW	1.04×		✓	✓	84.62 (1)	76.81 (4)	79.66 (0)	46.98 (4)	
	VPT-DEEP	1.18×			✓	<b>89.11 (4)</b>	<b>78.48 (6)</b>	<b>82.43 (2)</b>	<b>54.98 (8)</b>	

# Language-driven

Zixian Guo<sup>1,2\*</sup> Bowen Dong<sup>1</sup> Zhilong Ji<sup>2</sup> Jinfeng Bai<sup>2</sup> Yiwen Guo<sup>4</sup> Wangmeng Zuo<sup>1,2(✉)</sup>

<sup>1</sup>Harbin Institute of Technology <sup>2</sup>Tomorrow Advancing Life <sup>3</sup>Pazhou Lab, Guangzhou <sup>4</sup>Independent Researcher

zixian\_guo@foxmail.com cndongsky@gmail.com zhilongji@hotmail.com

jfbai.bit@gmail.com guoyiwen89@gmail.com wmuo@hit.edu.cn

12

- 早期经典方法都是language-driven, 如CoOp, CoCoOp等
- CLIP Prompt tuning **without image**?
- 合理性在于, CLIP已经实现了视觉-语言特征对齐,(text + text encoder)与(image + image encoder)可以相互替换。当目标数据集改变, 只要把语言特征调整到目标域, 再替换image encoder即可提升性能。

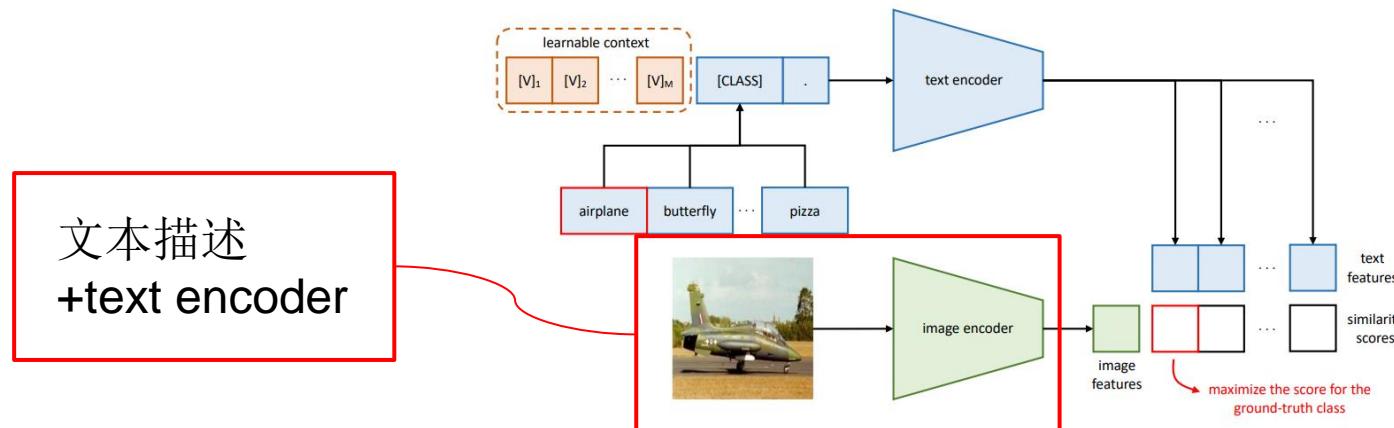


Figure 2: Overview of context optimization (CoOp).

# Language-driven

- 本文针对多标签分类（更关注标签间的结构关系），提出在训练过程不使用图像，用text description+text encoder替代image+image encoder
- 把text encoder适应到数据集标签所在的域

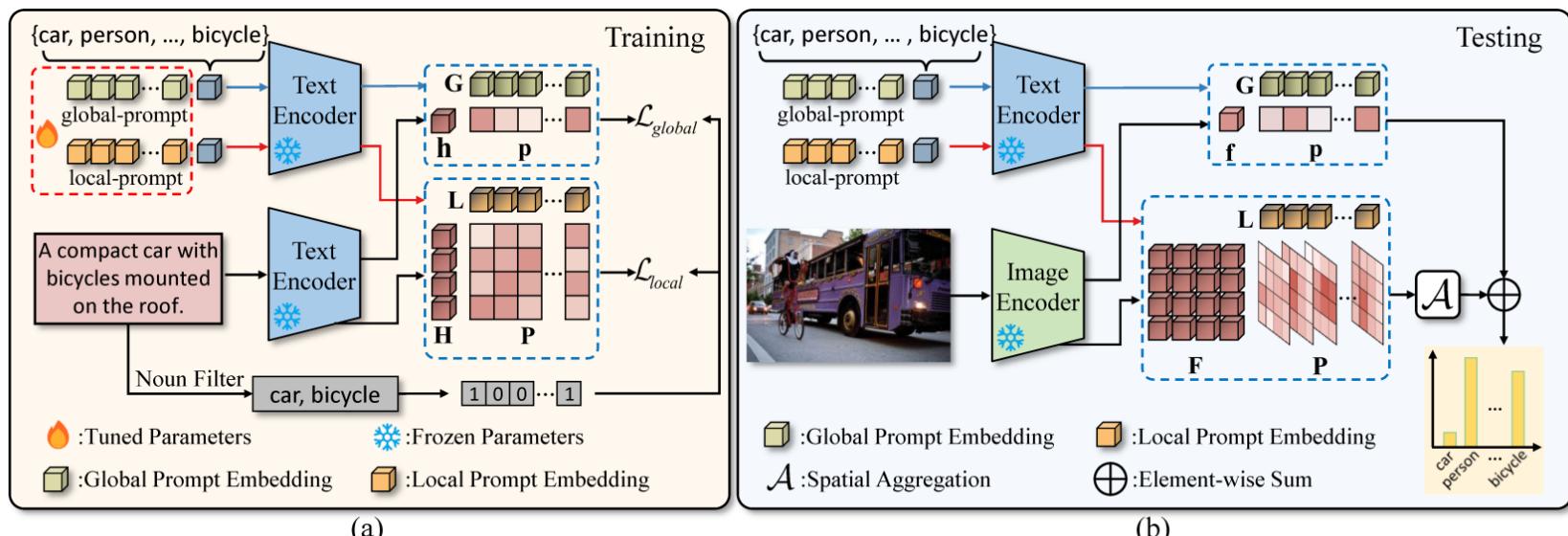


Figure 2. Training and testing pipeline of our proposed Text-as-Image (TaI) prompting, where we use text descriptions instead of labeled images to train the prompts. (a) During training, we use two identical text encoders from pre-trained CLIP to extract the global & local class embeddings (**G&L**) and overall & sequential text embeddings (**h&H**) respectively from the prompts and text description. The corresponding cosine similarity (**p&P**) between the embeddings are guided by the derived pseudo labels with ranking loss. (b) During testing, we replace the input from text descriptions to images. The global and local class embeddings can discriminate target classes from global & local image features (**f&F**). The final classification results are obtained by merging the scores of the two branches.



# Text Description

14

- 要求
  - 完整描述一张图的内容
  - 文本描述覆盖目标数据集的所有类别
- 直接使用MS-COCO、OpenImage等数据集的caption标注作为语料
- 类别名称的近义词词典

```
{'dog', 'pup', 'puppy', 'doggy'}  
{'person', 'people', 'man', 'woman', 'human'}  
{'bicycle', 'bike', 'cycle'}  
{'car', 'taxi', 'automobile'}  
{'boat', 'raft', 'dinghy'}  
...
```

- 从语料库中筛选出包含至少一个类别名称的句子，用来训练

# Visualization

15

## □ Text prompt tuning (local)

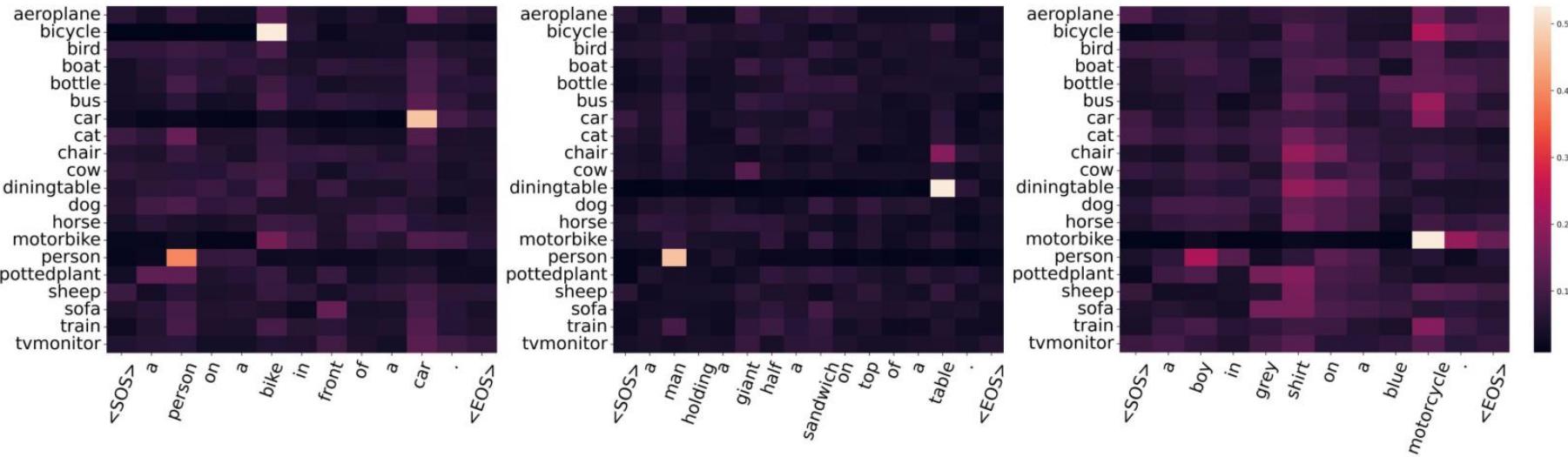
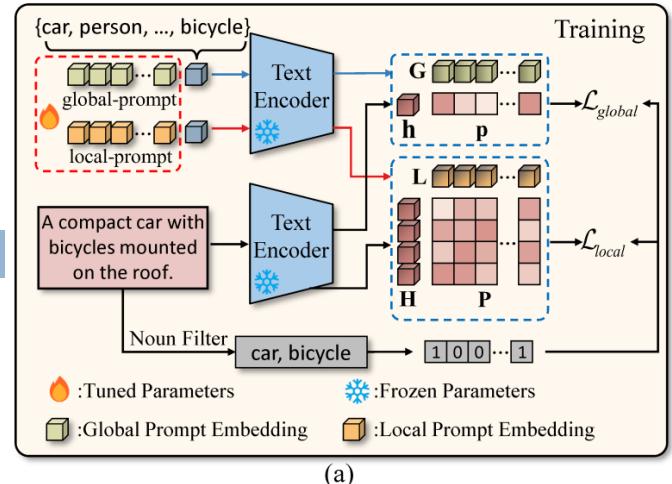


Figure 3. Visualization of correlations  $\mathbf{P}$  between the local class embedding  $\mathbf{L}$  and sequential token feature from texts. Each class embedding clearly correlates to words that describe the corresponding class (shown in highlight regions) rather than the global  $\langle \text{EOS} \rangle$  token.

## □ Few-shot multi-label learning

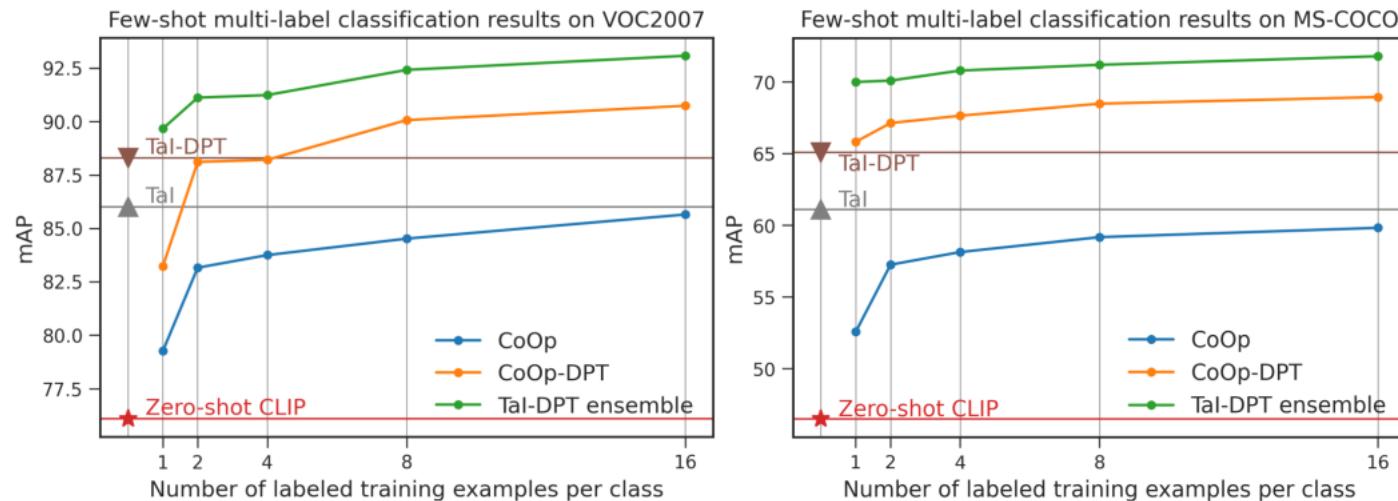


Figure 6. Comparison of different methods in few-shot multi-label recognition on VOC2007 and MS-COCO. Our zero-shot Tai-DPT can achieve comparable results with methods trained by 16-shot labeled image samples. And learned prompt ensemble proofs the complementarity between images and texts.



# Language-driven

17

## □ 性能对比

Table 1. Comparison with zero-shot methods on VOC2007, MS-COCO, and NUS-WIDE. Our proposed TaI-DPT outperforms CLIP [24] by a large margin on all datasets.

Method	DPT	VOC2007	MS-COCO	NUSWIDE
ZSCLIP	✗	76.2	47.3	36.4
	✓	77.3	49.7	37.4
TaI	✗	86.0	61.1	44.9
	✓	<b>88.3</b>	<b>65.1</b>	<b>46.5</b>

Table 3. Results of integrating our TaI-DPT with partial-label multi-label recognition method based on pre-trained CLIP. Our approach further improves the frontier performance of DualCoOp [28]. \* indicates the results based on our own reproduction.

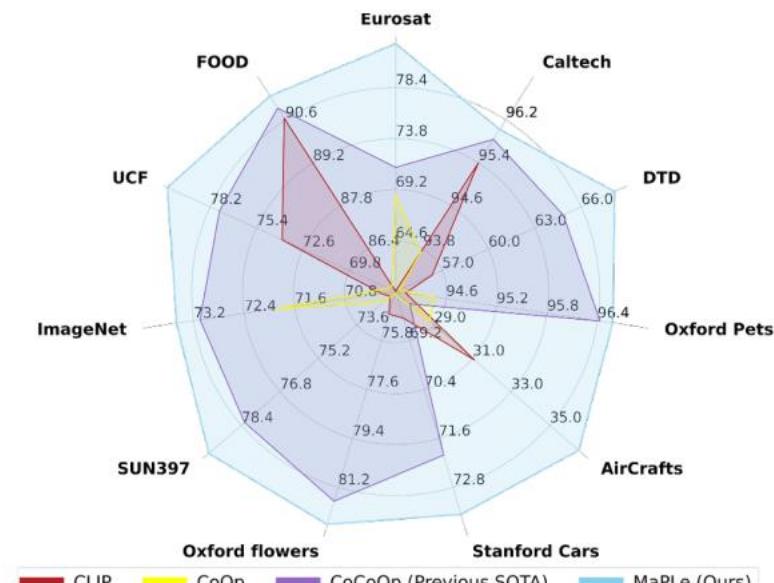
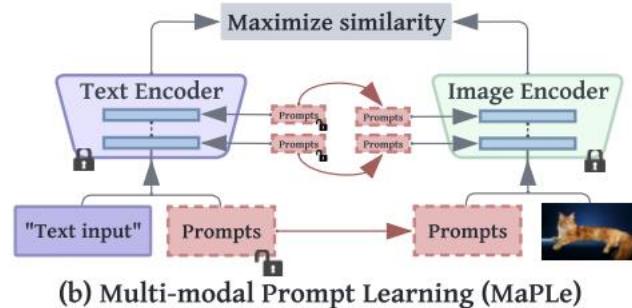
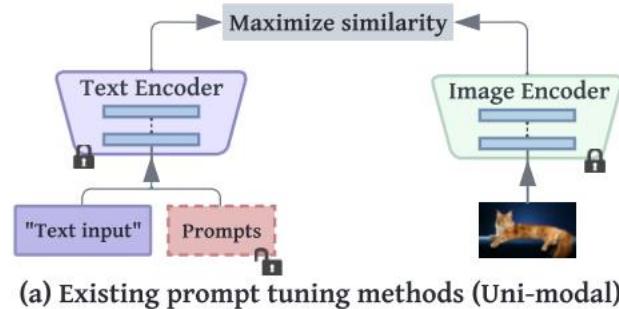
Datasets	Method	10%	20%	30%	40%	50%	60%	70%	80%	90%	Avg.
MS-COCO	SARB [23]	71.2	75.0	77.1	78.3	78.9	79.6	79.8	80.5	80.5	77.9
	DualCoOp [28]	78.7	80.9	81.7	82.0	82.5	82.7	82.8	83.0	83.1	81.9
	DualCoOp*	81.0	82.3	82.9	83.4	83.5	83.9	84.0	84.1	84.3	83.3
	+TaI-DPT	<b>81.5</b>	<b>82.6</b>	<b>83.3</b>	<b>83.7</b>	<b>83.9</b>	<b>84.0</b>	<b>84.2</b>	<b>84.4</b>	<b>84.5</b>	<b>83.6</b>
PascalVOC 2007	SARB [23]	83.5	88.6	90.7	91.4	91.9	92.2	92.6	92.8	92.9	90.7
	DualCoOp [28]	90.3	92.2	92.8	93.3	93.6	93.9	94.0	94.1	94.2	93.2
	DualCoOp*	91.4	93.8	93.8	94.3	94.6	94.7	94.8	94.9	94.9	94.1
	+TaI-DPT	<b>93.3</b>	<b>94.6</b>	<b>94.8</b>	<b>94.9</b>	<b>95.1</b>	<b>95.0</b>	<b>95.1</b>	<b>95.3</b>	<b>95.5</b>	<b>94.8</b>
NUS-WIDE	DualCoOp*	54.0	56.2	56.9	57.4	57.9	57.9	57.6	58.2	58.8	57.2
	+TaI-DPT	<b>56.4</b>	<b>57.9</b>	<b>57.8</b>	<b>58.1</b>	<b>58.5</b>	<b>58.8</b>	<b>58.6</b>	<b>59.1</b>	<b>59.4</b>	<b>58.3</b>

# Multi-modal

Muhammad Uzair Khattak<sup>1</sup> Hanoona Rasheed<sup>1</sup> Muhammad Maaz<sup>1</sup>Salman Khan<sup>1,2</sup> Fahad Shahbaz Khan<sup>1,3</sup><sup>1</sup>Mohamed bin Zayed University of AI <sup>2</sup>Australian National University <sup>3</sup>Linköping University

18

- 仅仅对语言模型做prompt tuning有局限性
- 本文设计了语言->视觉prompt的转换器，联合调整两个模态



(c) Performance comparison on base-to-novel generalization

# Multi-modal

19

- Vision language prompt coupling:  $\tilde{P}_k = \mathcal{F}_k(P_k)$ .

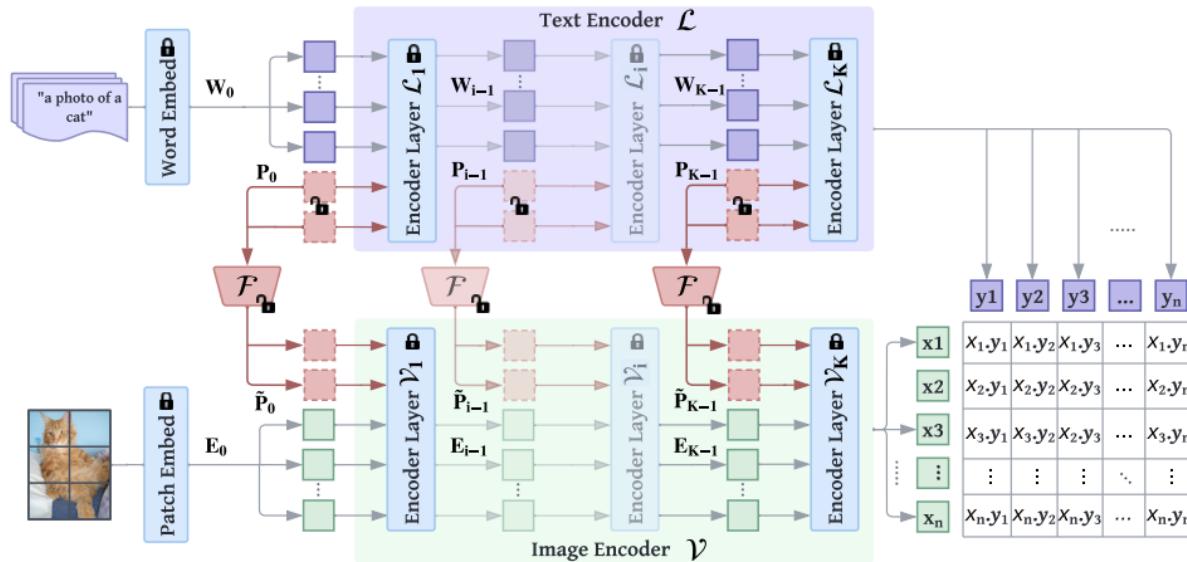


Figure 2. Overview of our proposed MaPLE (Multi-modal Prompt Learning) framework for prompt learning in V-L models. MaPLE tunes both vision and language branches where only the context prompts are learned, while the rest of the model is frozen. MaPLE conditions the vision prompts on language prompts via a V-L coupling function  $\mathcal{F}$  to induce mutual synergy between the two modalities. Our framework uses deep contextual prompting where separate context prompts are learned across multiple transformer blocks.



# Multi-modal

20

## □ 域泛化实验

(a) Average over 11 datasets				(b) ImageNet				(c) Caltech101			
	Base	Novel	HM		Base	Novel	HM		Base	Novel	HM
CLIP	69.34	74.22	71.70	CLIP	72.43	68.14	70.22	CLIP	96.84	94.00	95.40
CoOp	<b>82.69</b>	63.22	71.66	CoOp	76.47	67.88	71.92	CoOp	<b>98.00</b>	89.81	93.73
Co-CoOp	80.47	71.69	75.83	Co-CoOp	75.98	70.43	73.10	Co-CoOp	97.96	93.81	95.84
MaPLe	82.28	<b>75.14</b>	<b>78.55</b>	MaPLe	<b>76.66</b>	<b>70.54</b>	<b>73.47</b>	MaPLe	97.74	<b>94.36</b>	<b>96.02</b>
	+1.81	+3.45	+2.72		+0.68	+0.11	+0.37		-0.22	+0.55	+0.18
(d) OxfordPets				(e) StanfordCars				(f) Flowers102			
	Base	Novel	HM		Base	Novel	HM		Base	Novel	HM
CLIP	91.17	97.26	94.12	CLIP	63.37	<b>74.89</b>	68.65	CLIP	72.08	<b>77.80</b>	74.83
CoOp	93.67	95.29	94.47	CoOp	<b>78.12</b>	60.40	68.13	CoOp	<b>97.60</b>	59.67	74.06
Co-CoOp	95.20	97.69	96.43	Co-CoOp	70.49	73.59	72.01	Co-CoOp	94.87	71.75	81.71
MaPLe	<b>95.43</b>	<b>97.76</b>	<b>96.58</b>	MaPLe	72.94	74.00	<b>73.47</b>	MaPLe	95.92	72.46	<b>82.56</b>
	+0.23	+0.07	+0.15		+2.45	+0.41	+1.46		+1.05	+0.71	+0.85
(g) Food101				(h) FGVC Aircraft				(i) SUN397			
	Base	Novel	HM		Base	Novel	HM		Base	Novel	HM
CLIP	90.10	91.22	90.66	CLIP	27.19	<b>36.29</b>	31.09	CLIP	69.36	75.35	72.23
CoOp	88.33	82.26	85.19	CoOp	<b>40.44</b>	22.30	28.75	CoOp	80.60	65.89	72.51
Co-CoOp	90.70	91.29	90.99	Co-CoOp	33.41	23.71	27.74	Co-CoOp	79.74	76.86	78.27
MaPLe	<b>90.71</b>	<b>92.05</b>	<b>91.38</b>	MaPLe	37.44	35.61	<b>36.50</b>	MaPLe	<b>80.82</b>	<b>78.70</b>	<b>79.75</b>
	+0.01	+0.76	+0.39		+4.03	+11.90	+8.76		+1.08	+1.84	+1.48
(j) DTD				(k) EuroSAT				(l) UCF101			
	Base	Novel	HM		Base	Novel	HM		Base	Novel	HM
CLIP	53.24	<b>59.90</b>	56.37	CLIP	56.48	64.05	60.03	CLIP	70.53	77.50	73.85
CoOp	79.44	41.18	54.24	CoOp	92.19	54.74	68.69	CoOp	<b>84.69</b>	56.05	67.46
Co-CoOp	77.01	56.00	64.85	Co-CoOp	87.49	60.04	71.21	Co-CoOp	82.33	73.45	77.64
MaPLe	<b>80.36</b>	59.18	<b>68.16</b>	MaPLe	<b>94.07</b>	<b>73.23</b>	<b>82.35</b>	MaPLe	83.00	<b>78.66</b>	<b>80.77</b>
	+3.35	+3.18	+3.31		+6.58	+13.19	+11.14		+0.67	+5.21	+3.13

Table 3. Comparison with state-of-the-art methods on base-to-novel generalization. MaPLe learns multi-modal prompts and demonstrates strong generalization results over existing methods on 11 recognition datasets. Absolute gains over Co-CoOp are indicated in blue. ■



# 总结

21

- Prompt来自语言模型，本意是把不同任务转化为统一的语言描述模板输入模型，演化出各种变体。
- Prompt在语言模型效果更好，对视觉特征的改进较弱
- 因此，prompt tuning往往用于跨模态识别任务（CLIP系列）
  
- 检测、分割、定位等任务，需要更灵活的设计



- 研究背景
- Prompt tuning
- Adapter tuning
- Parameter tuning
- 总结

# Adapter tuning

23

- 怎样把预训练的ViT backbone适应到密集预测任务？
- 本文根据检测任务需要的先验信息设计adapter

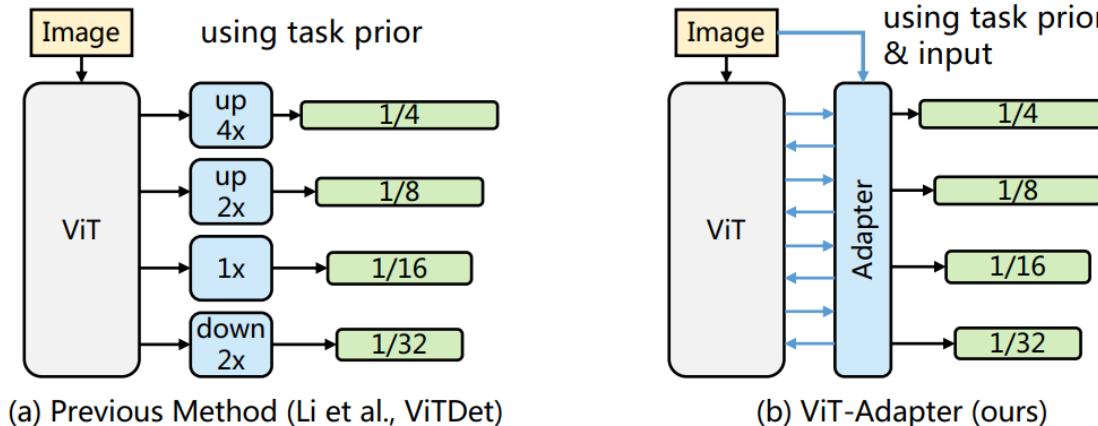


Figure 3: **Overview of ViT-Adapter and two related approaches.** Li et al. (2021b) and ViTDet (Li et al., 2022b) build simple feature pyramid to adapt plain ViT for object detection, which only consider task prior. Differently, our adapter utilizes both task prior and the input image.

# Adapter tuning

24

- CNN Stem, 增强底层纹理信息/空间归纳偏置
- 空间特征增强
- 多尺度特征

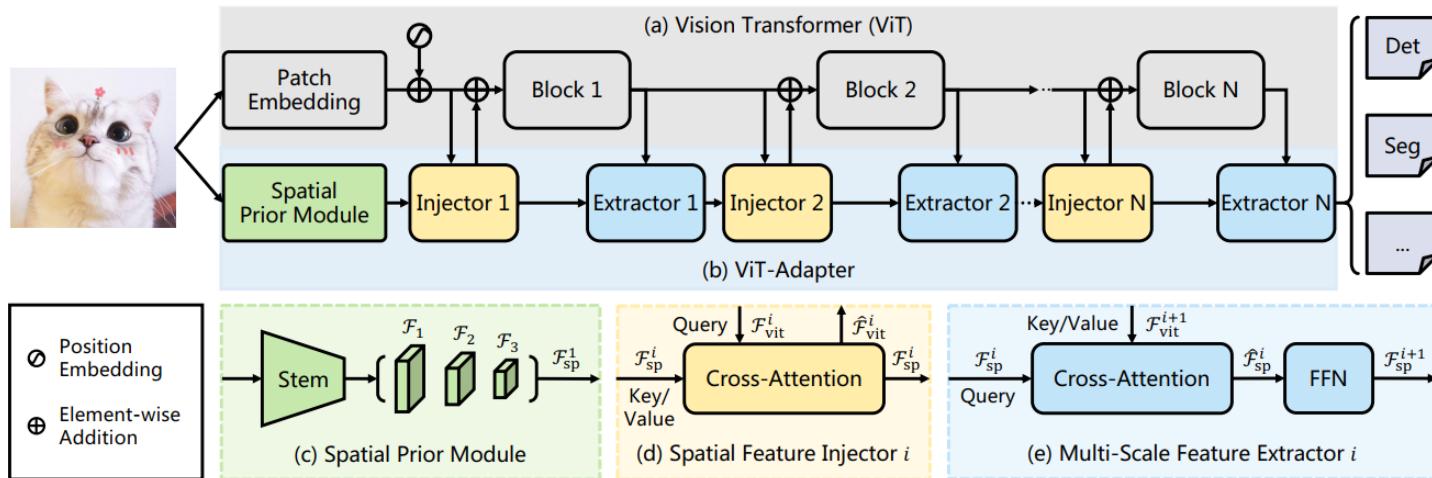
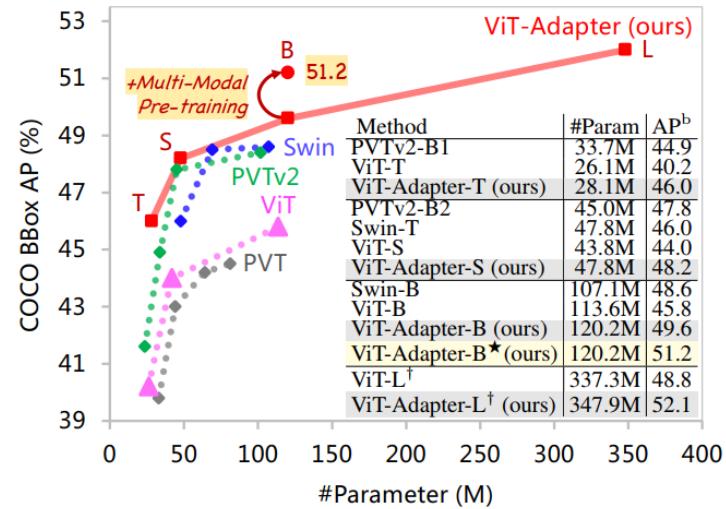


Figure 4: **Overall architecture of ViT-Adapter.** (a) The ViT, whose encoder layers are divided into  $N$  (usually  $N = 4$ ) equal blocks for feature interaction. (b) Our ViT-Adapter, which contains three key designs, including (c) a spatial prior module for modeling local spatial contexts from the input image, (d) a spatial feature injector for introducing spatial priors into ViT, and (e) a multi-scale feature extractor for reorganizing multi-scale features from the single-scale features of ViT.

# Adapter tuning

25

## □ 性能对比



Method	AP <sup>b</sup>	AP <sub>50</sub> <sup>b</sup>	AP <sub>75</sub> <sup>b</sup>	#P
Cascade Mask R-CNN 3×+MS schedule				
Swin-T (Liu et al., 2021b)	50.5	69.3	54.9	86M
Shuffle-T (Huang et al., 2021b)	50.8	69.6	55.1	86M
PVTv2-B2 (Wang et al., 2022a)	51.1	69.8	55.3	83M
Focal-T (Yang et al., 2021)	51.5	70.6	55.9	87M
ViT-S (Li et al., 2021b)	47.9	67.1	51.7	82M
ViT-Adapter-S (ours)	51.5	70.1	55.8	86M
ATSS 3×+MS schedule				
Swin-T (Liu et al., 2021b)	47.2	66.5	51.3	36M
Focal-T (Yang et al., 2021)	49.5	68.8	53.9	37M
PVTv2-B2 (Wang et al., 2022a)	49.9	69.1	54.1	33M
ViT-S (Li et al., 2021b)	45.2	64.8	49.0	32M
ViT-Adapter-S (ours)	49.6	68.5	54.0	36M
GFL 3×+MS schedule				
Swin-T (Liu et al., 2021b)	47.6	66.8	51.7	36M
PVTv2-B2 (Wang et al., 2022a)	50.2	69.4	54.7	33M
ViT-S (Li et al., 2021b)	46.0	65.5	49.7	32M
ViT-Adapter-S (ours)	50.0	69.1	54.3	36M

Method	AP <sup>b</sup>	AP <sub>50</sub> <sup>b</sup>	AP <sub>75</sub> <sup>b</sup>	#P
ATSS 3×+MS schedule				
Swin-T (Liu et al., 2021b)	47.2	66.5	51.3	36M
Focal-T (Yang et al., 2021)	49.5	68.8	53.9	37M
PVTv2-B2 (Wang et al., 2022a)	49.9	69.1	54.1	33M
ViT-S (Li et al., 2021b)	45.2	64.8	49.0	32M
ViT-Adapter-S (ours)	49.6	68.5	54.0	36M
GFL 3×+MS schedule				
Swin-T (Liu et al., 2021b)	47.6	66.8	51.7	36M
PVTv2-B2 (Wang et al., 2022a)	50.2	69.4	54.7	33M
ViT-S (Li et al., 2021b)	46.0	65.5	49.7	32M
ViT-Adapter-S (ours)	50.0	69.1	54.3	36M

Table 2: **Object detection with different frameworks on COCO val2017.** For fair comparison, we initialize all ViT-S/B models with the regular ImageNet-1K pre-training (Touvron et al., 2021). “#P” denotes the number of parameters. “MS” means multi-scale training.



- 研究背景
- Prompt tuning
- Adapter tuning
- Parameter tuning
- 总结

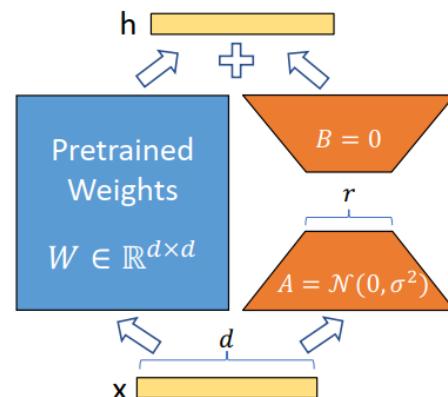
# Parameter tuning

Edward Hu\* Yelong Shen\* Phillip Wallis Zeyuan Allen-Zhu  
 Yuanzhi Li Shean Wang Lu Wang Weizhu Chen  
 Microsoft Corporation  
 {edwardhu, yeshe, phwallis, zeyuana,  
 yuanzhil, swang, luw, wzchen}@microsoft.com  
 yuanzhil@andrew.cmu.edu

27

- Adapter改变了模型结构，推理时带来更多负担。
- 将模型参数分解为  $W_0 + \Delta W$ ,  $W_0$  表示预训练参数,  $\Delta W$  表示微调后的变化量
- 训练阶段:  $h = W_0x + \Delta Wx = W_0x + BAx$  呈中的变化量
- 测试阶段: 将  $\Delta W$  加到预训练模型参数  $W_0$  中, 不改变结构, 不影响推断效率  

$$W = W_0 + \Delta W$$
- 低秩近似:  $\Delta W = BA$ ,  $B \in R^{d \times r}$ ,  $A \in R^{r \times d}$ ,  $r \ll d$



# Parameter tuning

28

Edward Hu\* Yelong Shen\* Phillip Wallis Zeyuan Allen-Zhu  
 Yuanzhi Li Shean Wang Lu Wang Weizhu Chen  
 Microsoft Corporation  
 {edwardhu, yeshe, phwallis, zeyuana,  
 yuanzhil, swang, luw, wzchen}@microsoft.com  
 yuanzhil@andrew.cmu.edu

## 性能对比

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB <sub>base</sub> (FT)*	125.0M	<b>87.6</b>	94.8	90.2	<b>63.6</b>	92.8	<b>91.9</b>	78.7	91.2	86.4
RoB <sub>base</sub> (BitFit)*	0.1M	84.7	93.7	<b>92.7</b>	62.0	91.8	84.0	81.5	90.8	85.2
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.3M	87.1 <sub>±.0</sub>	94.2 <sub>±.1</sub>	88.5 <sub>±1.1</sub>	60.8 <sub>±.4</sub>	93.1 <sub>±.1</sub>	90.2 <sub>±.0</sub>	71.5 <sub>±2.7</sub>	89.7 <sub>±.3</sub>	84.4
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.9M	87.3 <sub>±.1</sub>	94.7 <sub>±.3</sub>	88.4 <sub>±.1</sub>	62.6 <sub>±.9</sub>	93.0 <sub>±.2</sub>	90.6 <sub>±.0</sub>	75.9 <sub>±2.2</sub>	90.3 <sub>±.1</sub>	85.4
RoB <sub>base</sub> (LoRA)	0.3M	87.5 <sub>±.3</sub>	<b>95.1</b> <sub>±.2</sub>	89.7 <sub>±.7</sub>	63.4 <sub>±1.2</sub>	<b>93.3</b> <sub>±.3</sub>	90.8 <sub>±.1</sub>	<b>86.6</b> <sub>±.7</sub>	<b>91.5</b> <sub>±.2</sub>	<b>87.2</b>
RoB <sub>large</sub> (FT)*	355.0M	90.2	<b>96.4</b>	<b>90.9</b>	68.0	94.7	<b>92.2</b>	86.6	92.4	88.9
RoB <sub>large</sub> (LoRA)	0.8M	<b>90.6</b> <sub>±.2</sub>	96.2 <sub>±.5</sub>	<b>90.9</b> <sub>±1.2</sub>	<b>68.2</b> <sub>±1.9</sub>	<b>94.9</b> <sub>±.3</sub>	91.6 <sub>±.1</sub>	<b>87.4</b> <sub>±2.5</sub>	<b>92.6</b> <sub>±.2</sub>	<b>89.0</b>
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	3.0M	90.2 <sub>±.3</sub>	96.1 <sub>±.3</sub>	90.2 <sub>±.7</sub>	<b>68.3</b> <sub>±1.0</sub>	<b>94.8</b> <sub>±.2</sub>	<b>91.9</b> <sub>±.1</sub>	83.8 <sub>±2.9</sub>	92.1 <sub>±.7</sub>	88.4
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	0.8M	<b>90.5</b> <sub>±.3</sub>	<b>96.6</b> <sub>±.2</sub>	89.7 <sub>±1.2</sub>	67.8 <sub>±2.5</sub>	<b>94.8</b> <sub>±.3</sub>	91.7 <sub>±.2</sub>	80.1 <sub>±2.9</sub>	91.9 <sub>±.4</sub>	87.9
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	6.0M	89.9 <sub>±.5</sub>	96.2 <sub>±.3</sub>	88.7 <sub>±2.9</sub>	66.5 <sub>±4.4</sub>	94.7 <sub>±.2</sub>	92.1 <sub>±.1</sub>	83.4 <sub>±1.1</sub>	91.0 <sub>±1.7</sub>	87.8
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	0.8M	90.3 <sub>±.3</sub>	96.3 <sub>±.5</sub>	87.7 <sub>±1.7</sub>	66.3 <sub>±2.0</sub>	94.7 <sub>±.2</sub>	91.5 <sub>±.1</sub>	72.9 <sub>±2.9</sub>	91.5 <sub>±.5</sub>	86.4
RoB <sub>large</sub> (LoRA)†	0.8M	<b>90.6</b> <sub>±.2</sub>	96.2 <sub>±.5</sub>	<b>90.2</b> <sub>±1.0</sub>	68.2 <sub>±1.9</sub>	<b>94.8</b> <sub>±.3</sub>	91.6 <sub>±.2</sub>	<b>85.2</b> <sub>±1.1</sub>	<b>92.3</b> <sub>±.5</sub>	<b>88.6</b>
DeBERTa <sub>XXL</sub> (FT)*	1500.0M	91.8	<b>97.2</b>	92.0	72.0	<b>96.0</b>	92.7	93.9	92.9	91.1
DeBERTa <sub>XXL</sub> (LoRA)	4.7M	<b>91.9</b> <sub>±.2</sub>	96.9 <sub>±.2</sub>	<b>92.6</b> <sub>±.6</sub>	<b>72.4</b> <sub>±1.1</sub>	<b>96.0</b> <sub>±.1</sub>	<b>92.9</b> <sub>±.1</sub>	<b>94.9</b> <sub>±.4</sub>	<b>93.0</b> <sub>±.2</sub>	<b>91.3</b>

Table 2: RoBERTa<sub>base</sub>, RoBERTa<sub>large</sub>, and DeBERTa<sub>XXL</sub> with different adaptation methods on the GLUE benchmark. We report the overall (matched and mismatched) accuracy for MNLI, Matthew’s correlation for CoLA, Pearson correlation for STS-B, and accuracy for other tasks. Higher is better for all metrics. \* indicates numbers published in prior works. † indicates runs configured in a setup similar to Houldby et al. (2019) for a fair comparison.



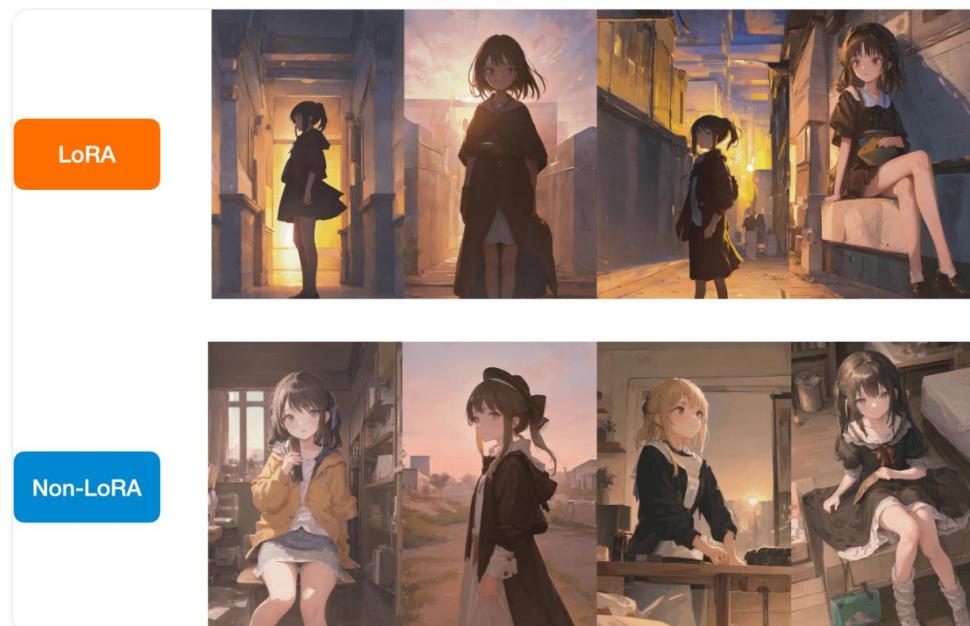
# Parameter tuning

29

- Parameter tuning方法不改变模型结构，更加普适。Lora最近也被广泛用于Diffusion model生成过程的微调。

prompt = "masterpiece, best quality, 1girl, at dusk"

Below is a comparison between the LoRA and the non-LoRA results:





- 研究背景
- Prompt tuning
- Adapter tuning
- Parameter tuning
- 总结



# 总结

31

- 大语言模型的发展整体上领先视觉模型，可以借鉴其经验
- 随着模型增大，full fine tuning会减少，由prompt, adapter等方式代替。
- 视觉模型不同于语言模型，需要更多先验知识。  
核心在于根据下游任务的特点设计相应的可微调参数



# Thanks!