

# Learning to Prompt for Open-Vocabulary Object Detection with Vision-Language Model

CVPR 2022

报告人：徐静远



# 目录

2

1

作者介绍

2

研究背景

3

研究方法

4

实验效果

5

总结



# 目录

3

1

作者介绍

2

研究背景

3

研究方法

4

实验效果

5

总结



# 作者介绍

4

Yu Du<sup>1</sup> Fangyun Wei<sup>2†</sup> Ziheng Zhang<sup>1</sup> Miaojing Shi<sup>3†</sup> Yue Gao<sup>2</sup> Guoqi Li<sup>1</sup>

<sup>1</sup>Tsinghua University <sup>2</sup>Microsoft Research Asia <sup>3</sup>King's College London

{duyu20, zh-zhang17}@mails.tsinghua.edu.cn liguoqi@mail.tsinghua.edu.cn

{fawe, yuegao}@microsoft.com miaojing.shi@kcl.ac.uk



Fangyun Wei

Microsoft Research Asia

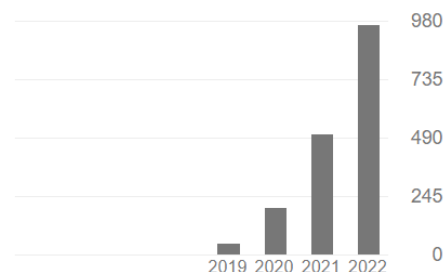
在 microsoft.com 的电子邮件经过验证

[Computer Vision](#) [Deep Learning](#) [Machine Learning](#)



引用次数

	总计	2017 年至今
引用	1739	1736
h 指数	16	16
i10 指数	20	20



开放获取的出版物数量

[查看全部](#)

1 篇文章

5 篇文章

无法查看的文章

可查看的文章

根据资助方的强制性开放获取政策

标题

引用次数

年份

[GCNet: Non-local networks meet squeeze-excitation networks and beyond](#)

1090

2019

Y Cao, J Xu, S Lin, F Wei, H Hu

CVF International Conference on Computer Vision Workshop (ICCVW), 1971-1980

[End-to-End Semi-Supervised Object Detection with Soft Teacher](#)

142

2021

M Xu, Z Zhang, H Hu, J Wang, L Wang, F Wei, X Bai, Z Liu

ICCV 2021

[Point-set anchors for object detection, instance segmentation and pose estimation](#)

75

2020

F Wei, X Sun, H Li, J Wang, S Lin

ECCV 2020

[Aligning Pretraining for Detection via Object-Level Contrastive Learning](#)

42

2021

F Wei, Y Gao, Z Wu, H Hu, S Lin

NeurIPS 2021 (Spotlight)

[RelationNet++: Bridging Visual Representations for Object Detection via Transformer Decoder](#)

41

2020

C Chi, F Wei, H Hu

NeurIPS 2020 (Spotlight)

# 目录

5

1

作者介绍

2

**研究背景**

3

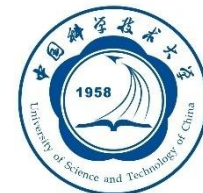
研究方法

4

实验效果

5

总结



## 研究背景一：目标检测

6

## 常用的数据集范式

► 以COCO数据集为例 [1]

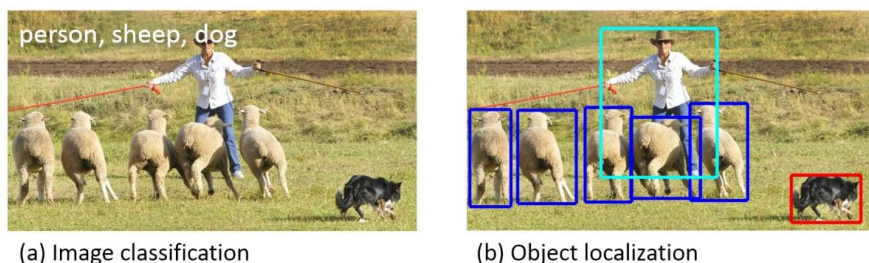


图1: COCO数据集对应的目标检测任务



图2: COCO数据集包含的80类

## □ 常用目标检测方法

➤ 以Faster-rcnn 为例 [2]

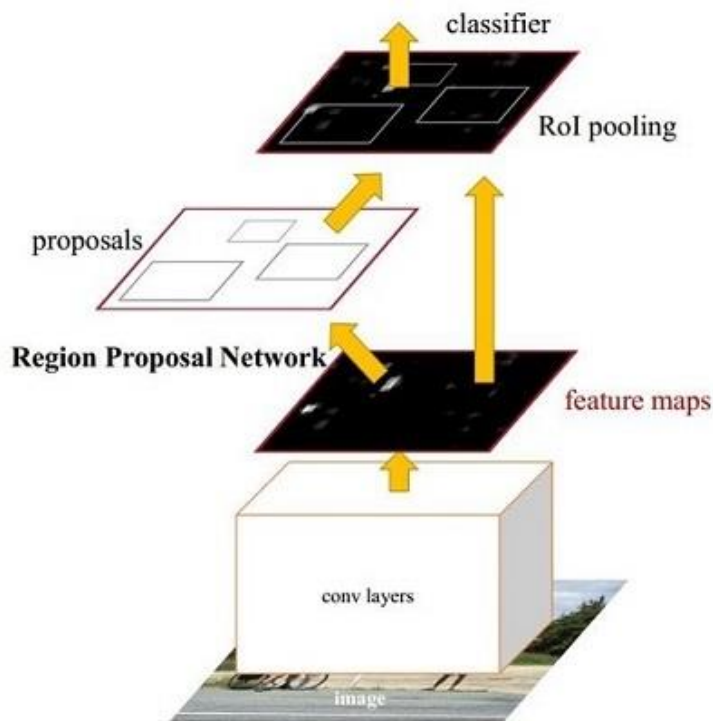


图3: Faster-rcnn方法示意图

[1]. Lin et al. Microsoft COCO: Common Objects in Context. ECCV2014

[2]. Ren et al. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. NIPS2015

# 研究背景二：视觉语言模型

7

- 可用于开放类的视觉语言模型
  - 以CLIP为例[1]

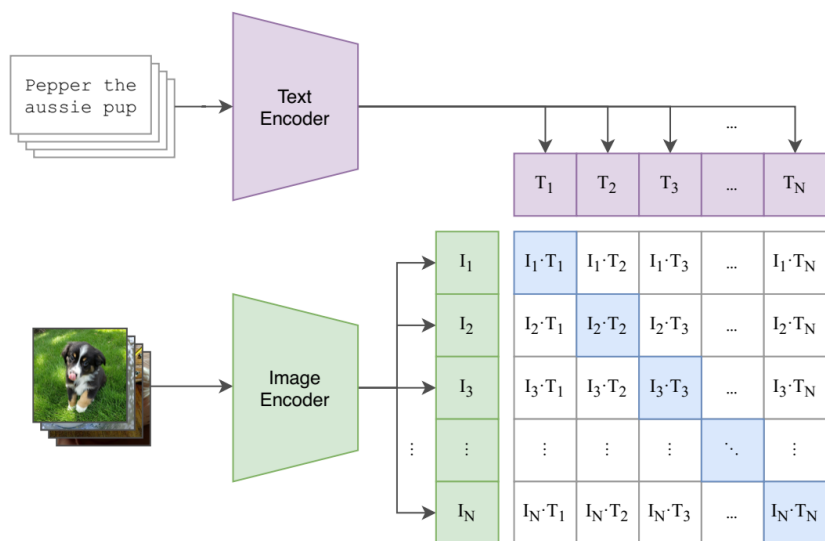


图1：CLIP模型的训练方式，4亿训练对

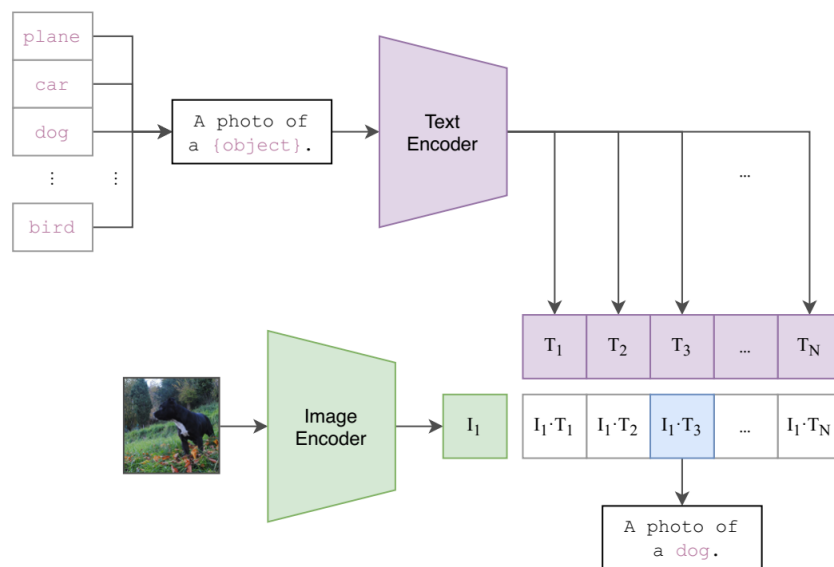


图2：CLIP模型的推理方式，可用于开放类

# 研究背景三：开放类目标检测

8

- 全监督场景下，如何检测更多类别？
  - 增加新类的标注和数据
- 缺点：
  - 类别标签固定
  - 增加收集数据，成本高昂

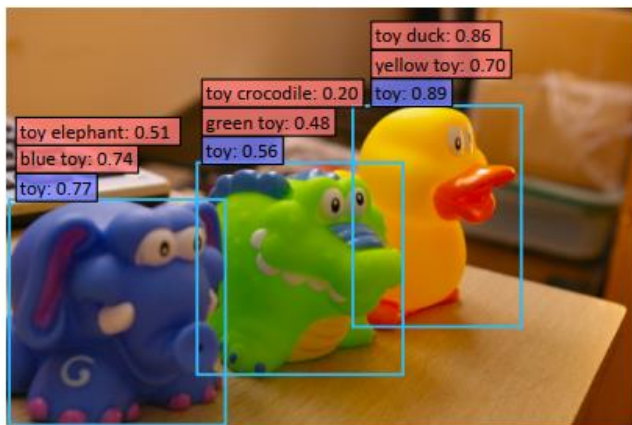
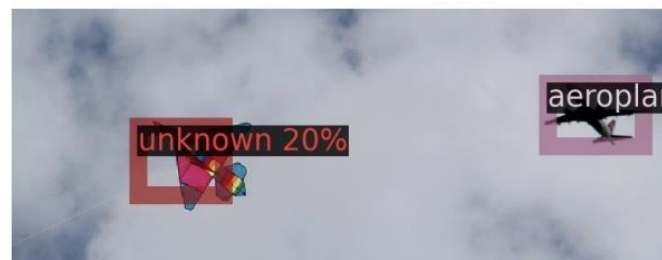
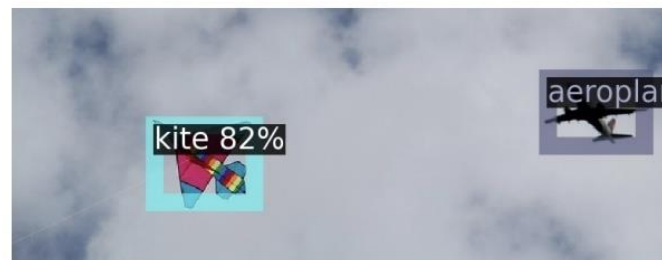


图1：不固定标签新类的目标检测



(a)



(b)

图2：新类检测需要重新收集数据示意



# 本研究baseline

9

## □ 直观方法 ViLD

- 将新类和已知类名称交给文本编码器获得嵌入
- 将候选框内的图像交给视觉编码器获得嵌入
- 计算文本和图像嵌入的距离进行识别

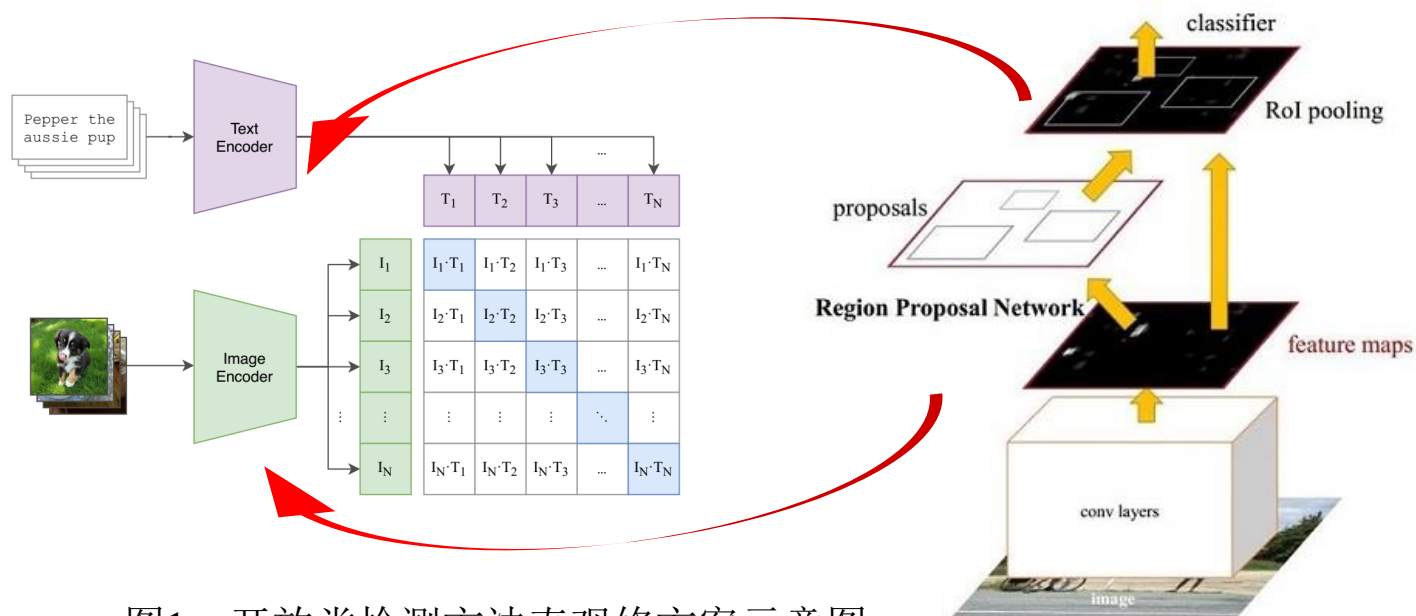


图1：开放类检测方法直观修方案示意图

# 本研究baseline

10

## ViLD架构

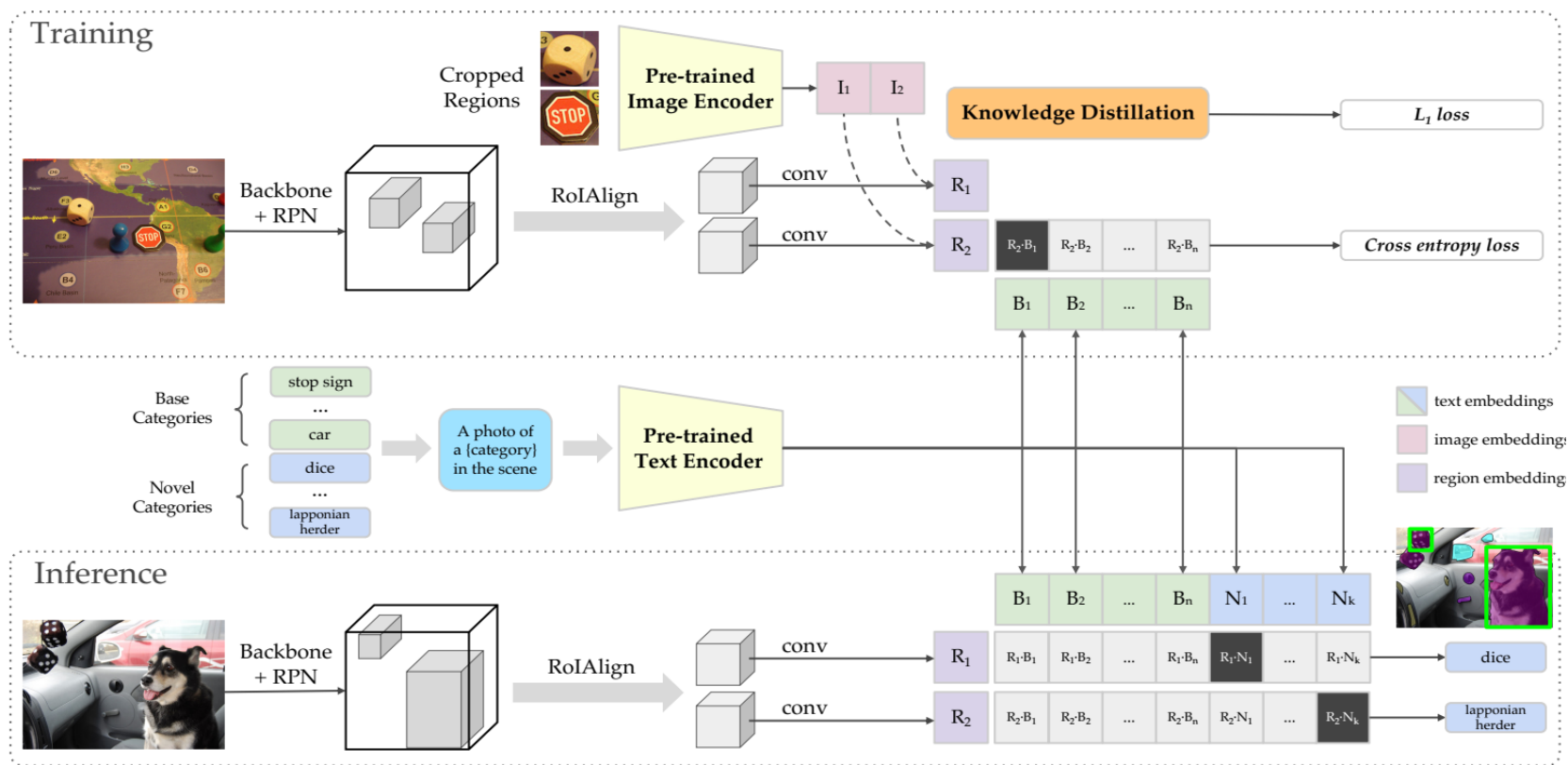


图1: ViLD方法整体框架, 上半部分是训练阶段, 下半部分是推理阶段, 黄色预训练模型表示是固定的。

# CoOp方案



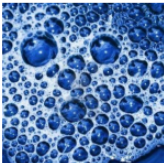

Caltech101	Prompt	Accuracy	Flowers102	Prompt	Accuracy
	a [CLASS].	82.68		a photo of a [CLASS].	60.86
	a photo of [CLASS].	80.81		a flower photo of a [CLASS].	65.81
	a photo of a [CLASS].	86.29		a photo of a [CLASS], a type of flower.	66.14
	$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	<b>91.83</b>		$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	<b>94.51</b>
	(a)			(b)	
Describable Textures (DTD)	Prompt	Accuracy	EuroSAT	Prompt	Accuracy
	a photo of a [CLASS].	39.83		a photo of a [CLASS].	24.17
	a photo of a [CLASS] texture.	40.25		a satellite photo of [CLASS].	37.46
	[CLASS] texture.	42.32		a centered satellite photo of [CLASS].	37.56
	$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	<b>63.58</b>		$[V]_1 [V]_2 \dots [V]_M$ [CLASS].	<b>83.53</b>
	(c)			(d)	

图1: CoOp的效果图

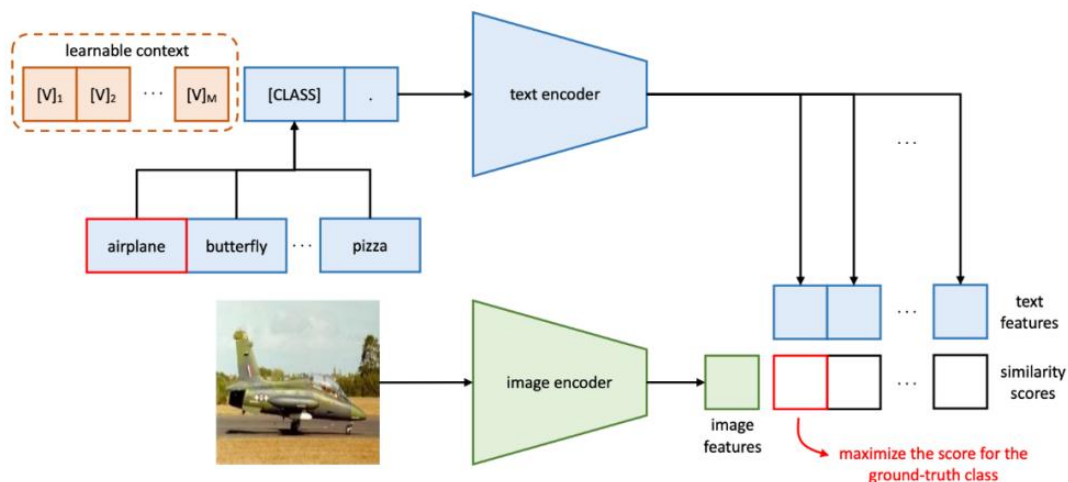


图2: CoOp的框架图

# 目录

12

1

作者介绍

2

研究背景

3

**研究方法**

4

实验效果

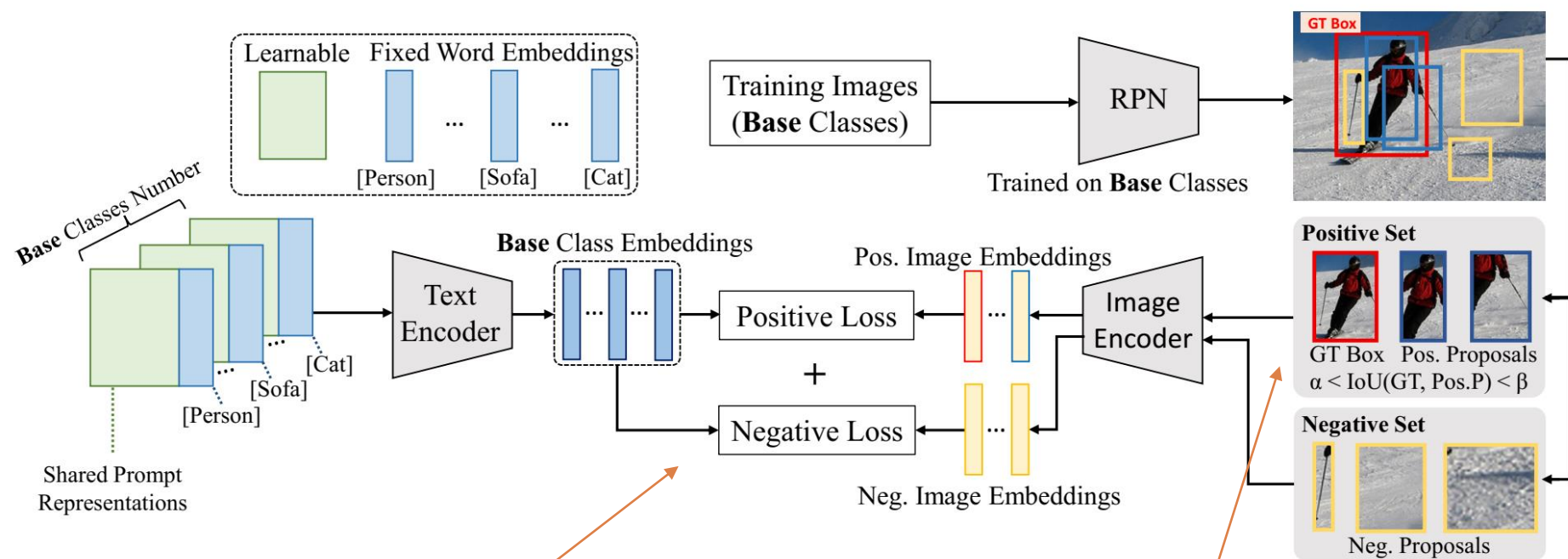
5

总结



# 研究方法

13



背景解释方案

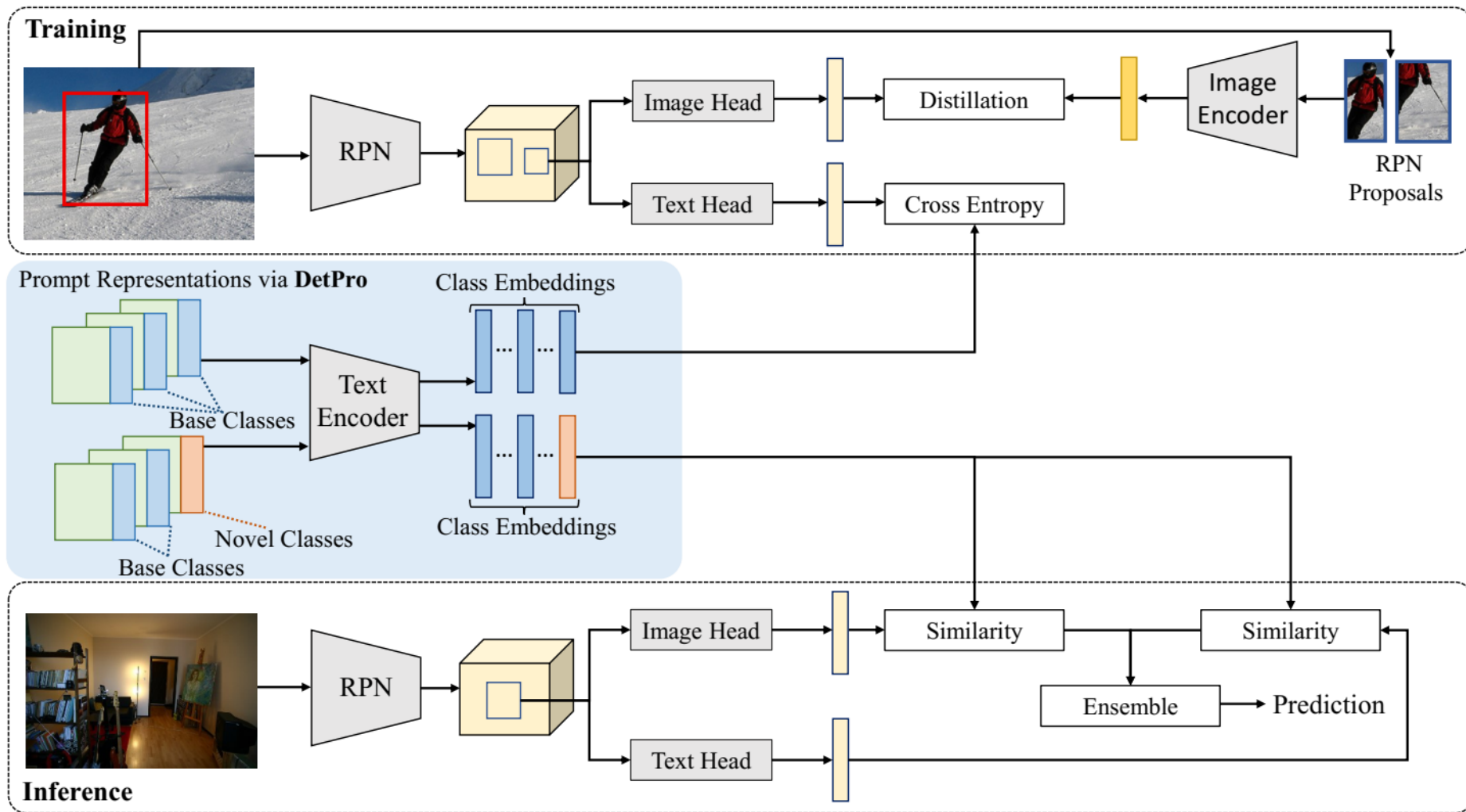
图1: DetPro的框架图

上下文打分方案



# 研究方法

14



# 目录

15

1

作者介绍

2

研究背景

3

研究方法

4

**实验效果**

5

总结





# 实验效果

16

## □ LVIS v1 [1]

- 866个基类 (frequent & common), 337个新类 (rare)

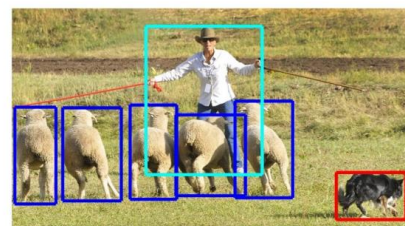


## □ COCO [2]

- 48个基类, 17个新类, 移除不包含在WordNet的15类



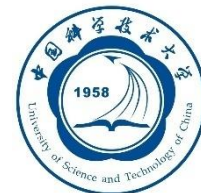
(a) Image classification



(b) Object localization

[1]. Gupta et al. Lvis: A dataset for large vocabulary instance segmentation. CVPR2019

[2]. Lin et al. Microsoft COCO: Common Objects in Context. ECCV2014



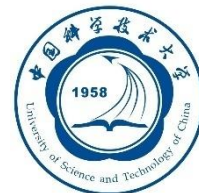


# 实验效果

17

## □ LVIS数据集

Method	Epoch	Detection				Instance segmentation			
		$AP_r$	$AP_c$	$AP_f$	AP	$AP_r$	$AP_c$	$AP_f$	AP
Supervised (base)	20	0.0	26.1	34.0	24.7	0.0	24.7	29.8	22.4
Supervised (base+novel)	20	15.5	25.5	33.6	27.0	16.4	24.6	30.6	25.5
ViLD (base) [7]	460	16.7	26.5	34.2	27.8	16.6	24.6	30.3	25.5
ViLD* (base) [7]	20	17.4	27.5	31.9	27.5	16.8	25.6	28.5	25.2
DetPro (base)	20	<b>20.8</b>	27.8	32.4	28.4	<b>19.8</b>	25.6	28.9	25.9



# 实验效果

18

## □ VOC, COCO, Object365 数据集

Method	Pascal VOC		COCO						Objects365					
	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>
Supervised	78.5	49.0	46.5	67.6	50.9	27.1	67.6	77.7	25.6	38.6	28.0	16.0	28.1	36.7
ViLD* [7]	73.9	<b>57.9</b>	34.1	52.3	36.5	21.6	38.9	46.1	11.5	17.8	12.3	4.2	11.1	17.8
DetPro	<b>74.6</b>	<b>57.9</b>	<b>34.9</b>	<b>53.8</b>	<b>37.4</b>	<b>22.5</b>	<b>39.6</b>	<b>46.3</b>	<b>12.1</b>	<b>18.8</b>	<b>12.9</b>	<b>4.5</b>	<b>11.5</b>	<b>18.6</b>

Table 2. We evaluate the LVIS-trained model on Pascal VOC test set, COCO validation set and Object365 validation set.



# 实验效果

19

## □ Ablation Study

Background proposals	$AP_r$	$AP_c$	$AP_f$	AP
10%	<b>19.1</b>	25.4	28.2	<b>25.4</b>
30%	18.3	<b>25.6</b>	<b>28.4</b>	<b>25.4</b>
50%	17.8	<b>25.6</b>	<b>28.4</b>	<b>25.4</b>
100%	17.6	25.1	28.2	25.0

Table 4. Ablation on number of background proposals involved in DetPro training.

GT	FG	BG	$AP_r$	$AP_c$	$AP_f$	AP
✓			15.3	<b>25.4</b>	27.9	24.6
✓	✓		16.9	25.1	27.7	24.7
✓		✓	17.7	25.3	<b>28.2</b>	25.1
✓	✓	✓	<b>19.1</b>	<b>25.4</b>	<b>28.2</b>	<b>25.4</b>

Table 5. Ablation study on the involvement of different training data. 'GT': ground-truth; 'FG': foreground; 'BG': background.



# 实验效果

20

## □ Ablation Study

Length	$AP_r$	$AP_c$	$AP_f$	AP
4	18.7	24.9	28.2	25.1
<b>8</b>	<b>19.1</b>	<b>25.6</b>	<b>28.3</b>	25.2
16	17.7	<b>25.6</b>	<b>28.3</b>	<b>25.3</b>

Table 7. Ablation study on context lengths.

Position	$AP_r$	$AP_c$	$AP_f$	AP
Front	16.4	24.5	<b>28.3</b>	24.6
Middle	18.0	25.1	<b>28.3</b>	25.1
<b>End</b>	<b>19.1</b>	<b>25.4</b>	28.2	<b>25.4</b>

Table 8. Ablation study of inserting class token into different positions of prompt representation.



## 实验效果

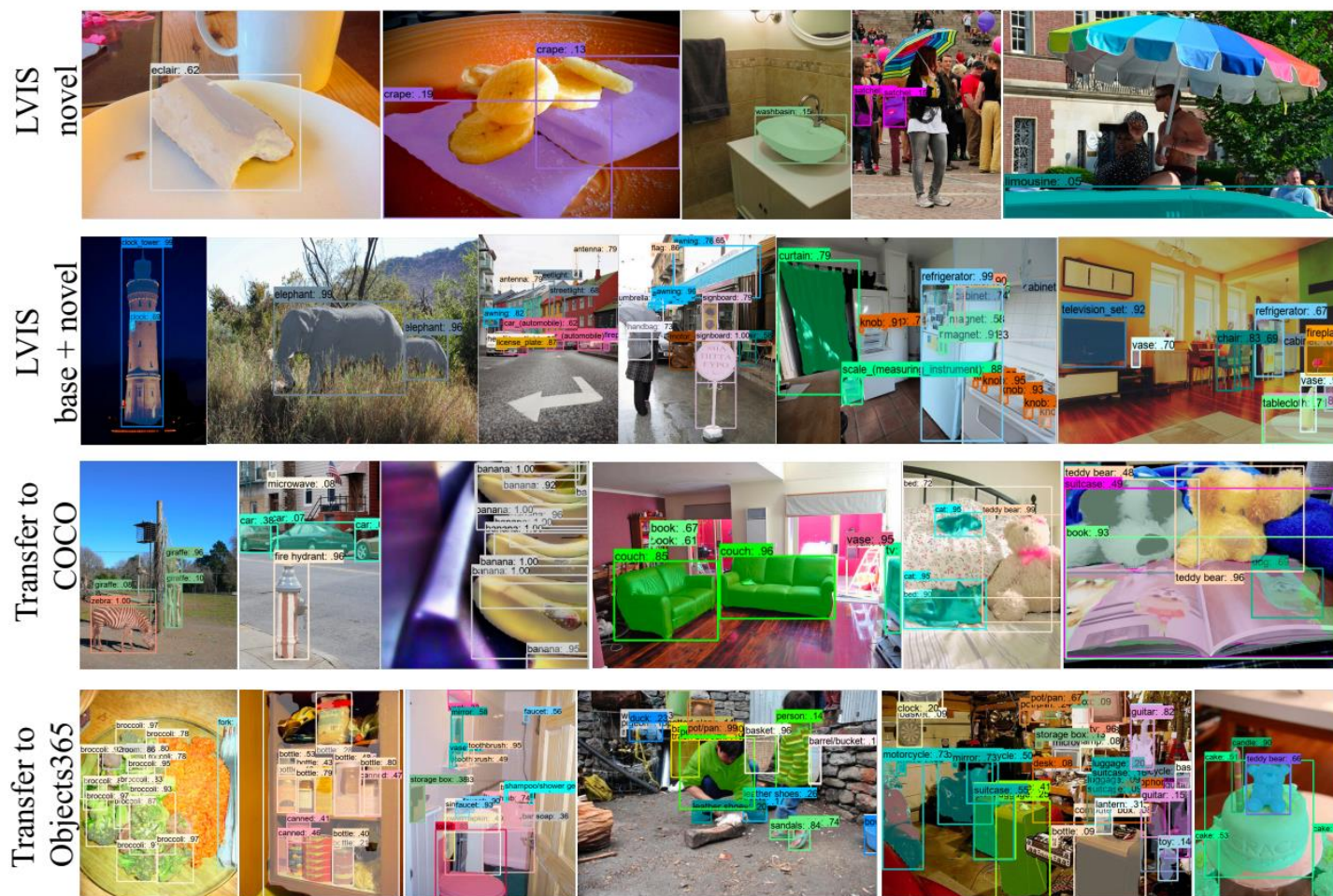


图1: 在LVIS, COCO, Object365数据集上的可视化实验。

# 目录

22

1

作者介绍

2

研究背景

3

研究方法

4

实验效果

5

总结



# 总结

23

## □ 总结

- 本文对于prompt learning出发，针对两个细节改进：
  - 基于ViLT方法处理背景类使用统一表征的问题
  - 基于ViLT方法处理前景框过于粗糙的问题
- 提升RPN对开放词汇目标的效果会是一个可能的改进方向

Table 1: **Training with only base categories achieves comparable average recall (AR) for novel categories on LVIS.** We compare RPN trained with base only vs. base+novel categories and report the bounding box AR.

Supervision	$AR_r@100$	$AR_r@300$	$AR_r@1000$
base	39.3	48.3	55.6
base + novel	41.1	50.9	57.0