



研究生学习工作周报

| | |
|------|------------------|
| 院 系 | 人工智能学院 |
| 专 业 | 电子信息 |
| 姓 名 | 余依函 |
| 学 号 | 231226006052 |
| 导 师 | 周静 张俊驰 |
| 周报日期 | 2023 年 10 月 14 日 |

摘要

1. R2VOS 代码复现。
2. Position-guided Text Prompt for Vision-Language Pre-training 论文阅读

目录

| | |
|---------------------------------|-----------|
| 摘要 | I |
| 第一章 学习工作总结 | 1 |
| 1.1 Transformer 在数学建模中的应用 | 错误!未定义书签。 |
| 1.2 R2VOS | 错误!未定义书签。 |
| 1.3 工作内容 | 15 |
| 第二章 第六周总结及第七周学习计划 | 16 |
| 参考 | 17 |

第一章 学习工作总结

1.1 R2V0S 代码复现

本周完成了 R2V0S 的全部过程，可以通过给出的文本信息进行视频帧的目标检测，结果如下图所示。图 1.1 为结果示例，图 1.2 为部分错误调试。



图 1.1 结果示例

报错记录

```
(py38pt) (base) tbgtb:~/R2V05/R2V05-master$ python demo.py --with_box_refine --binary --freeze_text_encoder --output_dir=output/demo --resume=checkpoint.pth --backbone resnet50 --gpu 1 --use_cycle --mix_query --neg_cls --is_eval --use_cls --demo_exp 'a big track on the road' --demo_path 'demo/demo_examples'
Traceback (most recent call last):
  File "demo.py", line 15, in <module>
    import util.misc as utils
  File "/home/tb/R2V05/R2V05-master/util/misc.py", line 37, in <module>
    from torchvision.ops import new_empty_tensor
ImportError: cannot import name 'new_empty_tensor' from 'torchvision.ops' (/home/tb/anaconda3/envs/py38pt/lib/python3.8/site-packages/torchvision/ops/_init_.py)
(py38pt) (base) tbgtb:~/R2V05/R2V05-master$ python demo.py --with_box_refine --binary --freeze_text_encoder --output_dir=output/demo --resume=checkpoint.pth --backbone resnet50 --gpu 1 --use_cycle --mix_query --neg_cls --is_eval --use_cls --demo_exp 'a big track on the road' --demo_path 'demo/demo_examples'
Traceback (most recent call last):
  File "demo.py", line 16, in <module>
    from models import build_model
  File "/home/tb/R2V05/R2V05-master/models/_init_.py", line 1, in <module>
    from .referformer import build
  File "/home/tb/R2V05/R2V05-master/models/referformer.py", line 21, in <module>
    from .backbone import build_backbone, build_amm_backbone
  File "/home/tb/R2V05/R2V05-master/models/backbone.py", line 17, in <module>
    from .amm_resnet import amm_resnet50, amm_resnet101
  File "/home/tb/R2V05/R2V05-master/models/amm_resnet.py", line 4, in <module>
    from torchvision.models.utils import load_state_dict_from_url
ModuleNotFoundError: No module named 'torchvision.models.utils'
(py38pt) (base) tbgtb:~/R2V05/R2V05-master$ python demo.py --with_box_refine --binary --freeze_text_encoder --output_dir=output/demo --resume=checkpoint.pth --backbone resnet50 --gpu 1 --use_cycle --mix_query --neg_cls --is_eval --use_cls --demo_exp 'a big track on the road' --demo_path 'demo/demo_examples'
Traceback (most recent call last):
  File "demo.py", line 16, in <module>
    from models import build_model
  File "/home/tb/R2V05/R2V05-master/models/_init_.py", line 1, in <module>
    from .referformer import build
  File "/home/tb/R2V05/R2V05-master/models/referformer.py", line 23, in <module>
    from .segmentation import CrossModalFPHDecoder, VisionLanguageFusionModule
  File "/home/tb/R2V05/R2V05-master/models/segmentation.py", line 26, in <module>
    from sklearn.decomposition import PCA
ModuleNotFoundError: No module named 'sklearn'
```

错误截图

ERROR1: [ImportError: cannot import name 'new_empty_tensor' from 'torchvision.ops'](#)

报错位置

```
▼ misc.py
1 if float(torchvision.__version__[2:4]) < 7:
2     from torchvision.ops import new_empty_tensor
3     from torchvision.ops.misc import _output_size
```

解决方法

```
▼ misc.py
1 if float(torchvision.__version__[2:4]) < 7:
2     from torchvision.ops import new_empty_tensor
3     from torchvision.ops.misc import _output_size
```

图 1.2 错误调试图

1.2 Position-guided Text Prompt for Vision-Language Pre-training 论文阅读

本周对 Position-guided Text Prompt for Vision-Language Pre-training 论文阅读进行阅读，此文章主要研究的是提高 VLP 的效率，本文提出了一个新颖的观念，目标检测器并不是成功检测的必要因素，如何利用位置信息对于 VLP 模型至关重要。可以尝试部署到 R2VOS 的预训练之中。如图 1.3 是部分论文阅读笔记。

Position-guided Text Prompt for Vision-Language Pre-training

💡 Meta Data

| | |
|-----------|--|
| Title | Position-guided Text Prompt for Vision-Language Pre-training |
| Journal | |
| Authors | Alex Jinpeng Wang; Pan Zhou; Mike Zheng Shou; Shuicheng Yan |
| Pub. date | 2023-06-07 |
| 期刊标签 | |
| DOI | 10.48550/arXiv.2212.09737 |
| 附件 | Wang et al 2023 Position-guided Text Prompt for Vision-Language Pre-training.pdf |

📖 研究背景 & 基础 & 目的

背景：

缺乏视觉的定位的模块、提升使用 VLP 训练的跨模态模型的视觉基础能力

提出了位置引导文本提示 PTP 增强跨模态模型的视觉 ground。

PTP 放弃了目标检测器直接进行推理

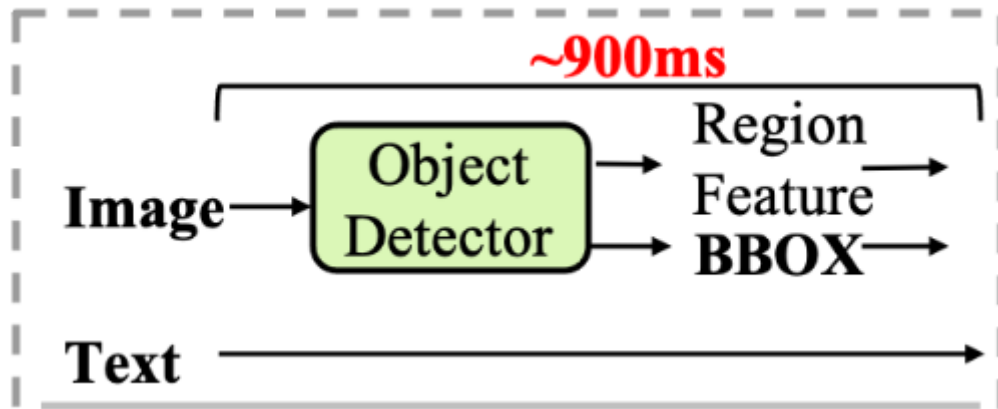
基础：

通用的跨模态模型首先以自监督的方式对大规模图像描述数据进行预训练，以获得足够的数据以获得更好的泛化能力。

然后对下游任务进行微调以进行适应。

“region features” (Wang 等, 2023, p. 1)

为了对位置信息进行建模，传统的 VLP 模型 采用在 1600 个类视觉基因组上预训练的 fast-rcnn 来提取显着特征区域特征和边界框。

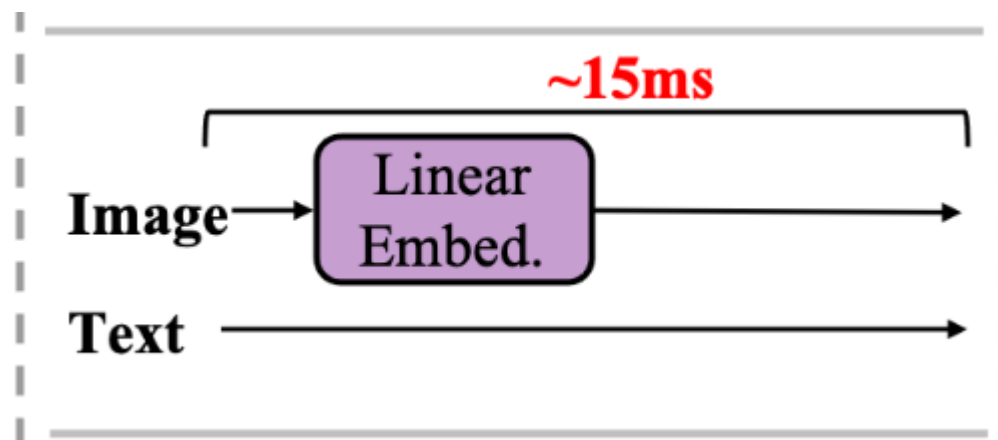


“To model the position information, traditional VLP models [3,23,46] (the top of Fig. 1 (a)) employ a faster-rcnn [34] pre-trained on the 1600 classes Visual Genome [17] to extract salient region features and bounding boxes.” (Wang 等, 2023, p. 1)

是什么+在哪里

没有上下文信息推理速度较慢

“end-to-end” (Wang 等, 2023, p. 2)



速度很快但无法很好的学习目标及其关系

目的

缓解端到端模型的位置缺失问题，并同时保持下游任务的快速推理时间

“In this work, we aims to ease the position missing problem for these end-to-end models, and keep fast inference time for downstream tasks at the same time.” (Wang 等, 2023, p. 2)

在图像和文本中添加基于位置的共同标识，将视觉 grounding 重新表达为一个填空题。从而简化对目标信息的学习。

在图像数据中实现自然语言的表达

生成标签块

将图像划分为 $N \times N$ 的块，并识别每个 block 中的目标。

生成文本提示

将查询文本放入基于位置的文本查询模板中。

“To ground natural language expressions in image data, PTP contains two components: (1) block tag generation to divide image into $N \times N$ blocks and to identify object in each block, and (2) text prompt generation that puts the query text into a position-based text query template.” (Wang 等, 2023, p. 2)

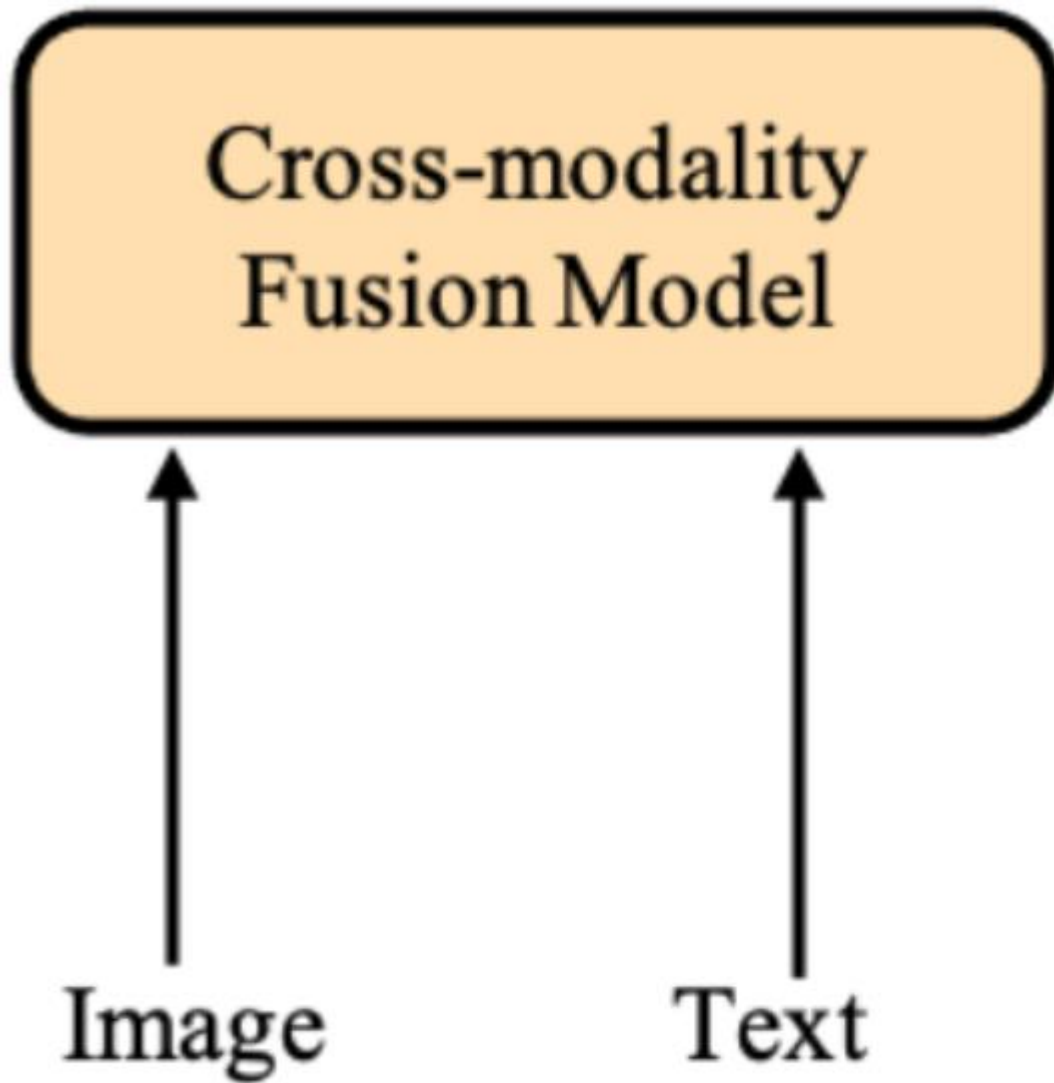
研究内容

Vision-language Pre-training Models

三个模型

单目标模型、多目标模型、多目标+融合编码模型（自己理解的翻译）。

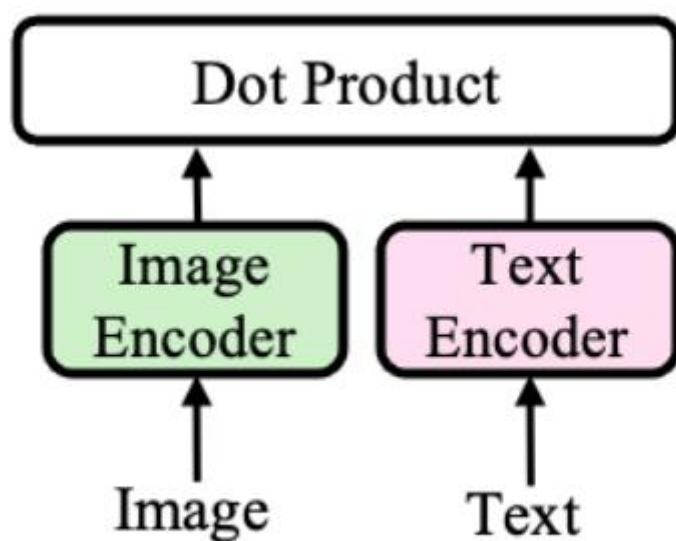
“One-stream Model” (Wang 等, 2023, p. 2)



(a). One-stream Model

对图像和文本输入进行处理。

“Dual-stream Model” (Wang 等, 2023, p. 2)

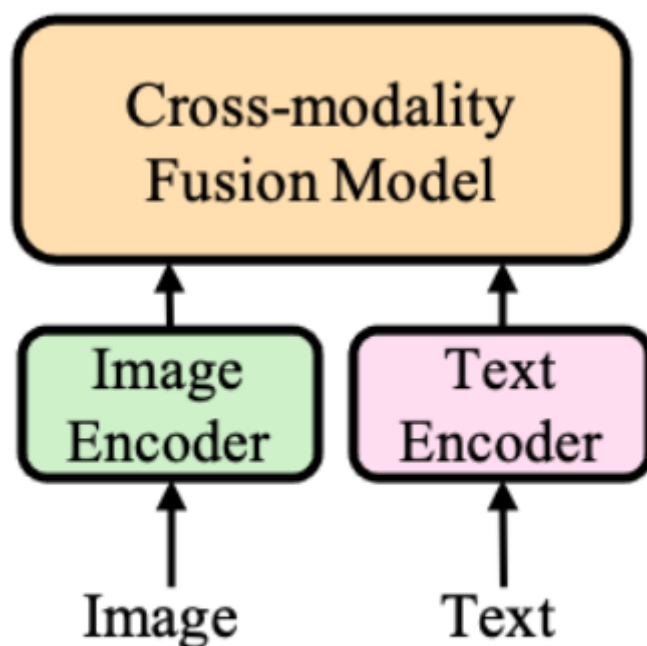


(b). Dual-stream Model

每种模态使用独立但同样 expensive transformer 编码器。

这两种模式在输入层没有串联，而在浅层，汇集的图像矢量和文本矢量之间存在交互。

“Dual-stream with Fusion Model” (Wang 等, 2023, p. 2)



(c). Dual-stream + Fusion Encoder

以上两种模式的融合

“Prompt Learning for Computer Vision” (Wang 等, 2023, p. 2)

多模态提示，它为 VLPT 模型提供了多模态提示调整，在一些视觉语言任务中取得了可喜的成果。

“Learn Position Information in VLP” (Wang 等, 2023, p. 2)

grounding 对于多模态任务至关重要，为了将此能力引入 VLP 将区域特征和边界框向量进行连接。但由于在下游任务中目标提取非常耗时。目前的研究为解决这个问题都是根据特定框架设计的，难以扩展。

本文目标提出了一个学习位置信息的通用框架，可以插入现有框架。

“Position-guided Text Prompt” (Wang 等, 2023, p. 3)

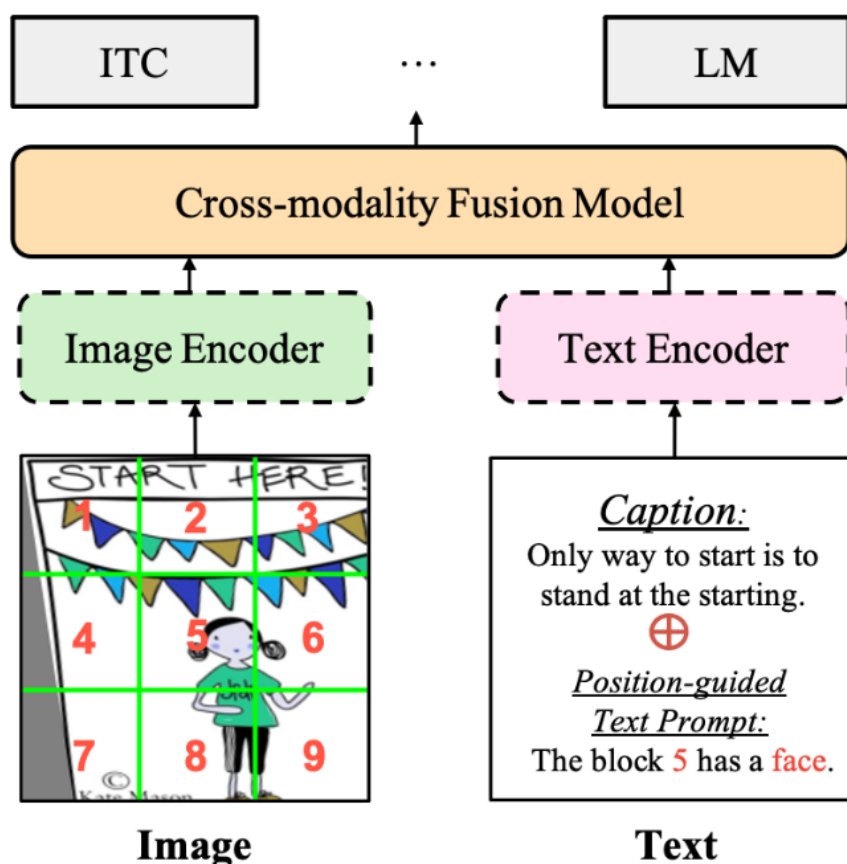
以 VILT、CLIP、BLIP 为例介绍如何将 PTP 和 VLP 框架结合来增强其视觉基础能力。

“PTP Paradigm” (Wang 等, 2023, p. 3)

PTP 帮助跨模态感知物体，并将物体与相关文本对其。

PTP 不同于传统的视觉语言配准方法，PTP 将物体特征和边界框作为输入

- 1.生成区块标签：将输入的图像分成若干块，并识别每个区块中的物体。
- 2.生成文本提示：根据 1 中的物体位置信息，将视觉 ground 任务转换为填空问题。



“Block Tag Generation” (Wang 等, 2023, p. 3)

对每一对图像文本平均分成 $N \times N$ 个区块

1.目标检测:

采用 VinVL 中的 **Faster-rcnn** 提取每幅图像中的所有目标。

选出 **预测置信度最高的前 k** 个目标，用 $\mathcal{O} = \left\{ o_i \right\}_{i=1}^K$ 表示 其中 $o_i = (z_i, q_i)$ 表示具有 4 维区域的位置向量 z 和目标类别 q

2.CLIP 模型:

采用区域监测而非 **haevy** 目标检测，PTP 可以通过 CLIP(ViT-B) 模型生成分块的目标监督

1.提取短语:

提取整个文本语料库中频率出现最高的 M 个短语或关键词（默认为 3000 个)使用 **NLTK** 被视为词汇记为 V 。

提取文本特征通过 CLIP 编码器嵌入所有 M 个短语或关键词。

获取每个区块中的图像嵌入 h 并计算文本特征的相似度。

相似度得分最高的关键字和短语为特定块的最终目标标签。

以下是每个块的目标标签索引计算公式

$$I = \operatorname{argmax}_{y \in [1, \dots, M]} \left(\frac{\exp(h^T e_y)}{\sum_{w \in V} \exp(h^T e_w)} \right), \quad (1)$$

$\operatorname{argmax}()$ 返回最大值坐标

$\exp()$ 以 e 为底的指数函数

h 是所选块的视觉特征嵌入

3. 优点

与目标检测器相比

产生更多样化的目标标签

块标签的生成速度快比 Faster-RCNN 块 40 倍

“Text Prompt Generation” (Wang 等, 2023, p. 3) 文本提示生成

设置一个简单的文本提示

“The block [P] has a [O].”

P 是所选块的索引表示块的位置 O 表示块 P 的标签。每一个 P 可能会包含多个 O 所以每次随机选择一个 O

“In this way, each sentence in our PTP incorporates fine-grained object position and language into a model, and thus provides a new way to align the objects and pertinent text.” (Wang 等, 2023, p. 4)

通过这种方式，我们的 PTP 中的每个句子都将细粒度的对象位置和语言合并到模型中，从而提供了一种对齐对象和相关文本的新方法。

“Pre-training with PTP” (Wang 等, 2023, p. 4) 使用 PTP 进行预训练

在这项工作中，我们将我们的 PTP 集成到主流 VLP 框架中，产生了 PTP-VILT、PTP-

CLIP 和 PTP-BLIP。

1.集成到现有任务中:

提示文本和原始标题被简单的填充到一起。 $x=[w,q]$ 其中 w 是文本, q 是本文方法生成的文本提示。

使用传统目标检测的端到端的训练方式训练 VLP 模型。

PTP-BLIP 中采用 LM loss、ITM、ITC loss

TC (图像文本对比学习), ITM (图像文本匹配), MLM (Masked Language Modeling, 有时会扩展到 MIM), LM (Language Modeling, 大部分可以看作是 captioning)

PTP-ViLT 中采用 ITM、MLM loss

PTP-CLIP 采用 ITC loss

2.新的 pretext task:

$$\mathcal{L}_{BLIP}(\theta) = -\mathbb{E}_{\lambda \sim D} \left[\sum_{\mathcal{L}}^{\mathcal{L}=J} \log B^{\theta}(\lambda^{\mathcal{L}} | \lambda^{<\mathcal{L}}) \right]$$

上式是目标预测的损失函数

D 为预测数据

y_1, \dots, y_{t-1} 是文本提示 q 的训练标记序列

t 为时间步

$p(t)=p(*|y_1, \dots, y_{t-1})$ 是需要通过设计来进行预测的概率分布

θ 是模型的可训练参数。

PTP 不需要修改任何基础网络, 并且可以应用于任何没有附加功能的 VLP 模型。PTP 旨在从原始像素图像中学习位置信息。只有在预训练阶段才需物体位置信息, 在下游的任务中采用正常的端到端的方式评估模型, 无需目标的信息, 从而摆脱繁重的目标特征提取以提高效率。

“Experiments” (Wang 等, 2023, p. 4)

“Experimental Settings” (Wang 等, 2023, p. 4)

1.“Datasets” (Wang 等, 2023, p. 4)

使用 COCO、VG、SBU、CC3M 组成 4M 设置。

Datasets. As in earlier studies [23, 46], we begin by using a 4M setup made up of four popular pre-training datasets (COCO [24], VG [17], SBU [29] and CC3M [35]). Following recent work [19], we also explore 14M setting, which includes additional CC12M [6] (actually only 10M image urls available) dataset besides 4M datasets. We refer readers to supplementary material for more dataset details.

2. “Training Settings” (Wang 等, 2023, p. 4)

PyTorch+8*NVIDIA A100 GPU

图像增强，探索 RandAugment 并使用除颜色反转之外的所有原始策略，因为颜色信息很重要。

RandAugment 是通过网格搜索两个参数(N, M)来寻找最佳DA的方法，其中 N 为 DA 的数量(变换)，M 为 Augmentation 的程度。这是为图像识别而设计的，所以在视频的情况下，必须对每一帧进行应用。

以图像相同的方式增强边界框以进行仿射变换。

在预训练期间随机采集分辨率为 224×224 的图像，并将图像分辨率提高到 384×384 进行微调。

3. “Baselines” (Wang 等, 2023, p. 4)

ViTB/16 作为基础视觉编码器并使用相同的数据集。

“Main Results” (Wang 等, 2023, p. 4)

PTP 在每个视觉下游任务的部署。

“Image-Text Retrieval” (Wang 等, 2023, p. 4) 图文检索

本文在 COCO 和 Flickr30K 基准上评估图像到文本检索 (TR) 和文本到图像检索 (IR) 的 PTP。

ViLT 是视觉和语言模型中最简单的架构，因为它使用 Transformer 模块代替单独的深度视觉嵌入器来提取和处理视觉特征。这种设计本质上显著改善了运行时间和参数效率。这是我们首次在不使用区域特征或深度卷积视觉嵌入器的情况下，在视觉和语言任务上取得了胜任的性能。

BLIP 是一种多模态 Transformer 模型，主要针对以往的视觉语言训练 (Vision-Language Pre-training, VLP) 框架的两个常见问题：大多数现有的预训练模型仅在基于理解的任务或者基于生成的任务方面表现出色，

很少有可以兼顾的模型。

在 MSCOCO 上对于 Image->Text Recall@1 提高了 13.8%

本文只在预训练阶段使用目标检测器。这表明目标检测器并不是成功的秘诀，如何利用位置信息对于 VLP 模型至关重要。

“Image Captioning” (Wang 等, 2023, p. 5) 图像字幕

数据集: No-Caps、COCO。

均使用在 COCO 上微调并具有 LM 损失的模型进行评估。

CIDEr 是 BLEU 和向量空间模型的结合。它把每个句子看成文档，然后计算 [TF-IDF 向量](#)（只不过 term 是 n-gram 而不是单词）的余弦夹角，据此得到候选句子和参考句子的相似度，同样是不同长度的 n-gram 相似度取平均得到最终结果。优点是不同的 n-gram 随着 TF-IDF 的不同而有不同的权重，因为整个语料里更常见的 n-gram 包含了更小的信息量。

“Visual Question Answering” (Wang 等, 2023, p. 5) 视觉问答

与 ViLT 相比，PTP 提高了 1.8%

“Visual Reasoning” (Wang 等, 2023, p. 6) 视觉推理

自然语言视觉推理 (NLVR2) 任务是一个二元分类任务，给定两个图像的三元组和一个自然语言问题。虽然不如 SimVLM 性能但是此方法也接近 VinVLLarge 模型，采用更大的模型并使用来自强对象检测器的对象特征而不是原始像素图像作为输入。

“Video-Language Tasks” (Wang 等, 2023, p. 6) 视频语言任务

在这个实验中分析了我们的方法对视频语言任务的泛化能力。

“Ablation & Design Choices” (Wang 等, 2023, p. 6)

“The Variations of Architecture.” (Wang 等, 2023, p. 6) 结构变体

使用三种不同类型的 baseline 进行实验: ViLT、CLIP 和 BLIP，以探索 PTP 的影响。比较这些 baseline 实验的结果，我们发现 PTP 极大地提高了 i2t 和 t2i 性能。这表明 PTP 具有良好的通用性。由于本文在下游任务中不使用目标检测器或提示，因此计算成本与基线模型保持一致，但比基于目标特征的 VinVL 快 20 倍。

“Text Prompt vs. Additional Pretext Task” (Wang 等, 2023, p. 7) 文本提示与附加借口任务

提示比借口要好得多，特别是对于 COCO 字幕 CIDEr (127.2 vs 123.5)。在这项工作中，由于其效率，本文使用提示作为默认值。

“Other Types of Text Prompt” (Wang 等, 2023, p. 7)

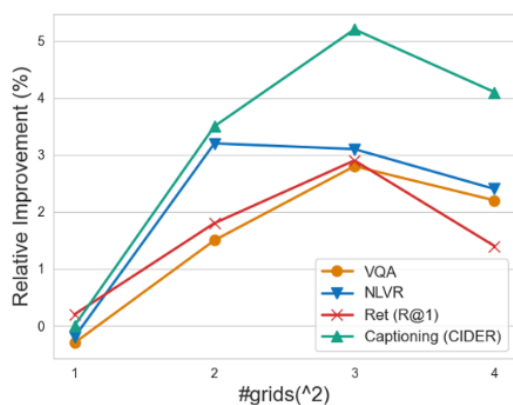
“The Importance of Position in Text Prompt” (Wang 等, 2023, p. 7)

在这个实验中检查了提示 PTP 获取各种粒度信息（例如没有位置信息）的效果。

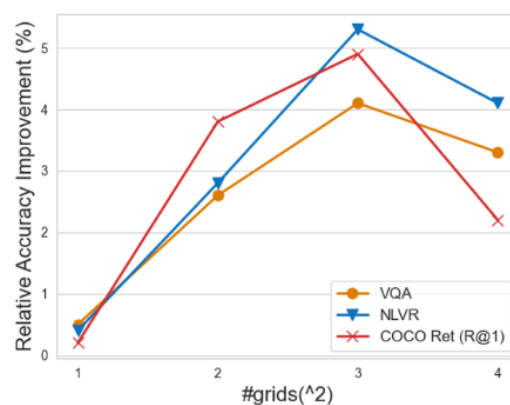
每一个信息都很重要，如果缺失，下游的性能会变的越来越差。

“Number of Blocks” (Wang 等, 2023, p. 7)

建议使用 3×3 块，因为它具有准确性。



(a). PTP-BLIP



(b). PTP-ViLT

“Is Object Detector Necessary?” (Wang 等, 2023, p. 8)

除了基于 ResNext152 的强大目标检测器之外，我们还使用了一个较小的 Faster-rcnn 网络，该网络利用 ResNet101 作为主干。

使用更强的检测器会带来更好的结果，但同时也会带来巨大的计算成本。

结论:剪辑模型是 PTP 中目标检测器的一个很好的替代方案。

“Visualization” (Wang 等, 2023, p. 8)

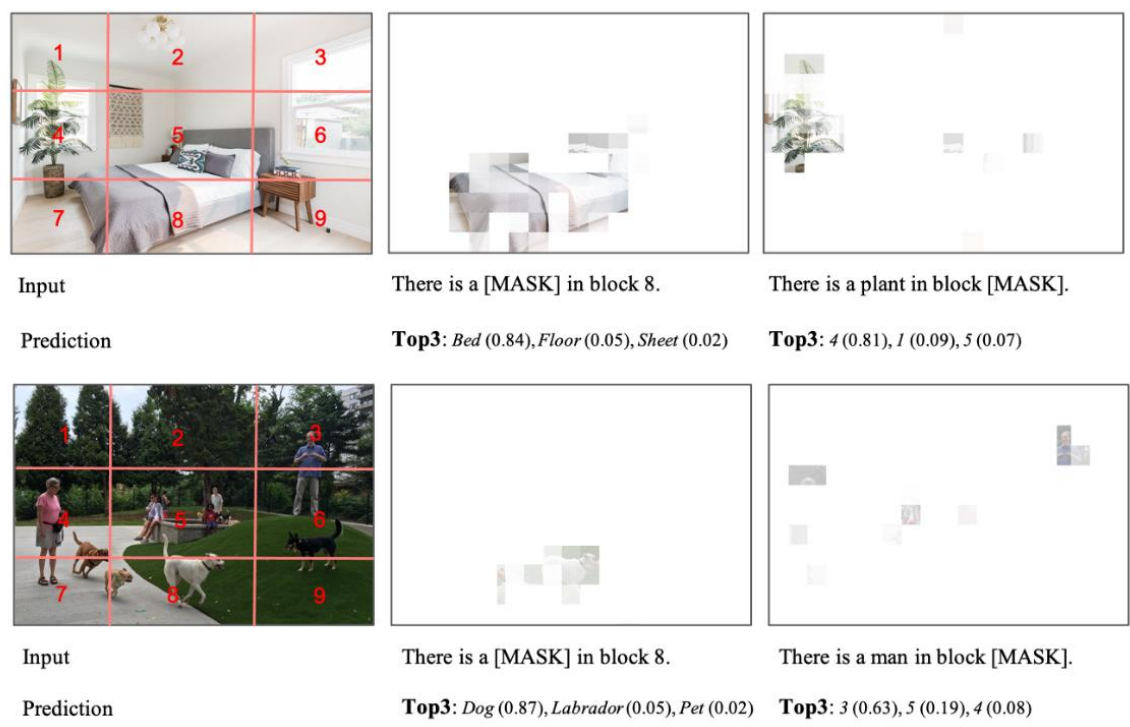


图 1.3 论文阅读笔记

1.3 工作内容

对领导留言板进行爬取。

第二章 第七周总结及第六周学习计划

| 序号 | 上周任务 | 完成情况及备注 | 本周任务 |
|-----|---|---------|----------------------|
| 1 | R2VOS 代码复现 | 完成 | R2VOS 与 PTP 是否可以结合使用 |
| 2 | Position-guided Text Prompt for Vision-Language Pre-training 论文阅读 | 完成 | CVPR 论文阅读一篇（目前还没找好） |
| 3 | | | |
| 备注： | | | |

参考

1. <https://github.com/lxa9867/R2VOS.git>
2. <https://arxiv.org/abs/1706.03762v7>
3. <https://arxiv.org/abs/1706.03762v7>