



学业繁重 先不聊了

文 | Yimin_饭煲

2021年, 多模态领域大概是人工智能研究者们关注者最多的一个领域了。随着各种模态数据 集的增长和算力的发展,研究者们开始不断地尝试在一个模型中融合来自各个模态的信息。

而在多模态领域的研究中,和视频相关的任务被认为是最复杂的。

一方面,高质量的视频数据集比图像数据集更加困难,因此数据集的数量和质量往往受限;另 一方面,视频数据集中含有文本、图像、语音等多个模态的信息,还要考虑时间线,融合起来 比单纯的图像-文本数据更加复杂。

在AI领域久负盛名的Allen研究所向这一复杂的问题发起了挑战,提出了MERLOT系列工作。

第一篇 MERLOT:Multimodal Neural Script Knowledge Models 发表于Neurips 2021, 使用了大量的视频数据进行自监督预训练,在12个视频问答任务上取得了SOTA;

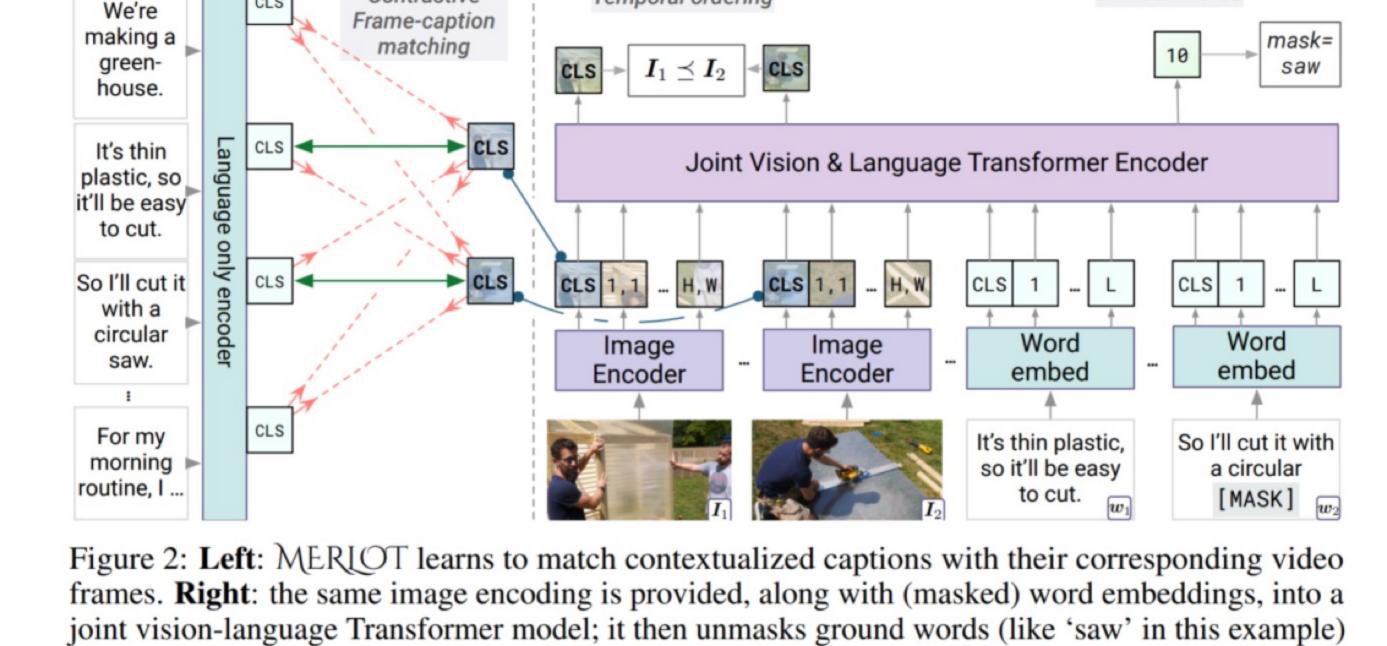
MERLOT Reserve: Neural Script Knowledge through Vision and

Unmask words

Language and Sound 则于今年年初刚刚发布,进一步深度融合了视频中的语音信息,在多 个任务上又取得了明显提升。下面,就让我们一起来学习这两篇十分Solid的工作吧~

论文链接: MERLOT: https://arxiv.org/pdf/2106.02636.pdf **MERLOT Reserve:** https://arxiv.org/pdf/2201.02639.pdf 主要方法

Temporal ordering



MERLOT这一工作使用了视觉编码器、语言编码器和联合编码器。作者们设计了三个任务来 进行优化。

and puts scrambled video frames into the correct order.

第一个任务是Contrastive Frame-caption matching(标题-帧匹配),作者们使用视觉编码 器编码图片得到的[CLS]表示和文本编码器编码句子得到的[CLS]表示进行对比学习,使得图

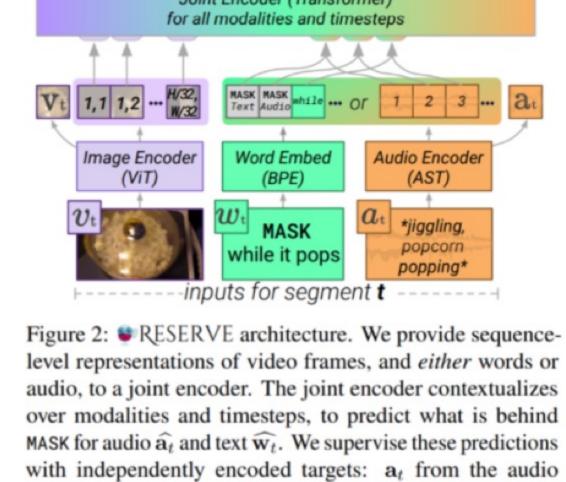
第二个任务是Masked Language Modeling,对模型的语言建模能力进行优化。

所有帧中随机选取i帧并进行打乱,将位置编码 (e.g. $[image_t]$)替换为随机且独特的位置编码 (e.g. $[image_{unk_0}]$). 这些随机的位置编码和原有的位置编码分别进行学习,可以让模型学到恢

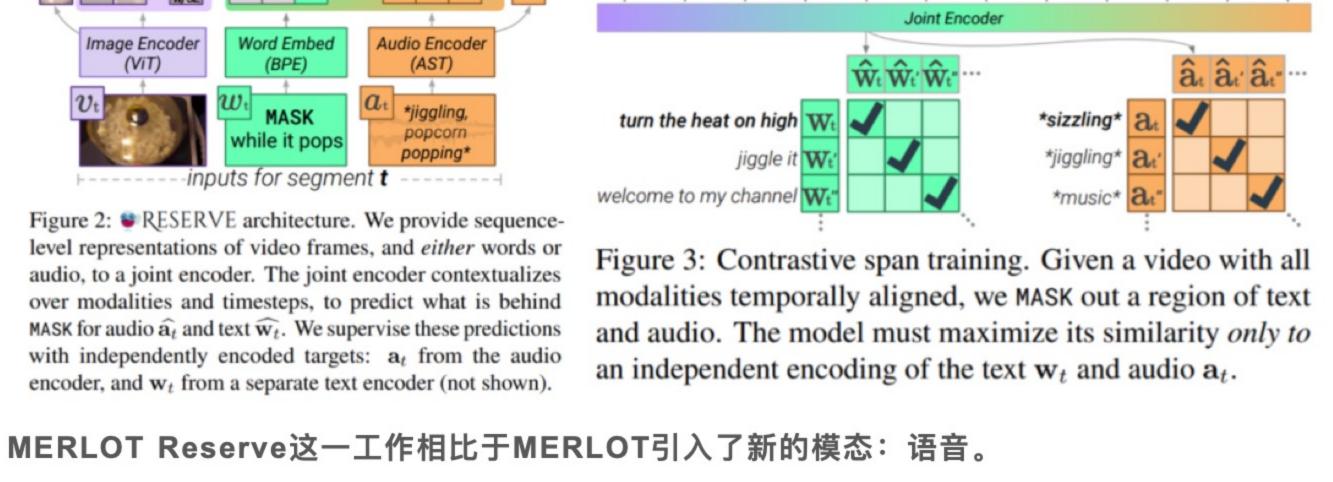
第三个任务是Temporal Reordering,在40%的情况下,随机选择一个整数i,从视频输入中的

复被扰乱的帧顺序的能力。 这个任务的损失函数是针对一对视频帧 (t_i,t_j) 拼接隐状态 (h_{t_i},h_{t_j}) ,使用两层MLP分类器进行 二分类(t_i 和 t_j 的前后关系)。

Predict MASKed text and audio Joint Encoder (Transformer) for all modalities and timesteps third cup of MASK



像编码器具备较好的表示学习性能。



为了更好的融合来自三个模态的信息,**作者们提出了更通用,更统一的训练任务。**

数:

encoder, and \mathbf{w}_t from a separate text encoder (not shown).

对于每一个Batch的输入,只输入视频的帧和文本/语音中的一个(由于文本和语音的信息具有 重复性), 并且MASK文本/语音中的一部分。作者们提出了对比区域匹配(Contrastive Span

 $L_{mask->text} = rac{1}{|W|} \sum_{w_t \in W} (\log rac{\exp(\sigma \hat{w}_t w_t)}{\sum_{w \in W} \exp(\sigma \hat{w}_t w)})$ 其中 \hat{w}_t 为[MASK]位置的隐状态表示, w_t 为[MASK]掉的信息的隐状态表示,w为Batch中其他 样本(负样本)的隐状态表示。同理定义了 $L_{text->mask}$,定义:

Matching)这一任务, 给定匹配的视频帧/文本/语音数据, 以文本为例, 最小化交叉熵损失函

伏态表示。同理定义了
$$L_{text->mask}$$
,定义: $L_{text} = L_{text->mask} + L_{mask->text}$

同样的,可以定义 L_{audio} 和 L_{frame} 。定义总体的损失函数为:

 $L = L_{text} + L_{audio} + L_{frame}$

作者们还使用了一些技巧来提升得到的特征表示的质量,感兴趣的小伙伴可以去原文细读~。

数据集

对于大规模的预训练工作,除开训练方法之外另一个值得关注的部分就是使用的数据集了。

Conceptual ∪ COCO

₹ YT-Temporal-180M

作者们在12个视频问答数据集上开展了实验,大幅度刷新了SOTA。

Short

Short

Short

Short

Short

ERNIE-ViL-Large [123]

Villa-Large [39]

Villa-Base [39]

Situated Reasoning (STAR)

Model

Supervised SoTA

Just Ask (ZS)[122]

● RESERVE-B

RESERVE-L

RESERVE-L (+audio)

Random

(test acc; %)

ClipBERT [73]

Interaction Sequence Prediction Feasibility Overall

39.8 43.6 32.3 31.4 36.7

25.0 25.0 25.0 25.0 25.0

43.9 42.6 37.6 33.6 39.4

UNITER-Large [20]

Test

Test

Test

Test

Tasks

MSRVTT-QA

MSR-VTT-MC

TGIF-Transition

TGIF-Frame QA Test

TGIF-Action

HowTo100M

MERLOT收集的数据集为YT-Temporal-180M,从600万公开的YouTube视频中抽取得到。

作者们选取的数据集比起HowTo100M和VLOG等局限于特定领域的数据集范围更大,主题更

广。 后来的实验表明,如果仅使用HowTo100M这样的数据集进行训练,会降低模型在下游任务上

的性能。 **VCR** Dataset

58.9

66.3

75.2

MERIOT

43.1

90.9

94.0

96.2

69.5

41.5 [118]

88.2 [127]

82.8 [67]

87.8 [67]

60.3 [67]

	HowTo100M-sized YT-Temporal-180 YTT180M, raw ASR	M	72.8 72.8
tion data leads yet still below	video) data is important. Applying to poor results. Our model perform wour (more diverse) YT-Tempora e. Using raw ASR (vs. denoised AS	s be	otter on HowTo100M, 0M, even when con-

数据的有效性。 结果 作者们通过大量的实验证实了MERLOT和MERLOT Reserve的有效性。对于MERLOT模型,

Split Vid. Length ActBERT [127] ClipBERT_{8x2} [67]

88.2

37.4

82.8

87.8

60.3

在MERLOT Reserve这一工作中,作者们扩充了数据集,提出了YT-Temporal-1B数据集,

包含2000万Youtube视频,进一步提升了数据集的多样性,而模型强大的性能也说明了扩充

	LSMDC-FiB QA	Test	Short	48.6		-	48.6 [127]	52.9	
	LSMDC-MC	Test	Short	-		-	73.5 [121]	81.7	
	ActivityNetQA	Test	Long	-		-	38.9 [118]	41.4	
	Drama-QA	Val	Long	-		-	81.0 [56]	81.4	
	TVQA	Test	Long			-	76.2 [56]	78.7	
	TVQA+	Test	Long	-		-	76.2 [56]	80.9	
	VLEP	Test	Long	-		-	67.5 [66]	68.4	
state of	the art methods	s in 12 c	downstre	art methods on videam tasks that invol	lve sl	hort and	long videos.		
能提升,	超过了许多使用	用了其(他监督信	信息的模型。					
					V	CR (tes	t)		
	_	Mod	el	Q-	→A	$QA{\rightarrow}R$	$Q{\rightarrow}AR$		

VilBERT [80] 73.3 74.6 54.8 B2T2 [4] 72.6 75.7 55.0 VisualBERT [76] 73.2 52.4 71.6

79.2

78.9

77.3

76.4

83.5

83.8

80.8

79.1

66.3

65.7

62.8

60.6

ased	MERLOT [126]	80.6	80.4	65.1
/ideo-based	● RESERVE-B	79.3	78.7	62.6
Vid	RESERVE-L	84.0	84.9	72.0
ar W sin	t leaderboard perform to the compare it with the ngle models, including the odels that utilize heaven (e.g. object detection)	rman larges ng im	ce on st sub age-ca	VCR. mitted aption apervi-
由于MERLOT Reserve	使用了大量的数据进行了自	自监督预	训练,	因此 在零样本学习上有着不
错的性能,在STAR数据	集上相比于有监督的SOTA	都有着	明显的	是升 。

39.8 40.5 35.5 36.0 38.0 CLIP (VIT-B/16) [91] 16.5 12.8 2.3 2.0 CLIP (RN50x16) [91] 13.4 14.5 2.1 39.9 41.7 36.5 **37.0** 38.7 2.3 44.4 40.1 38.1 35.0 39.4 17.9 15.6 2.7 26.1 42.6 41.1 37.4 32.2 38.3 15.6 19.3 4.5 26.7 RESERVE-B (+audio) **44.8 42.4 38.8** 36.2 **40.5** 20.9 17.5 3.7 29.1

Table 4: Zero shot results. On STAR, RESERVE obtains state-of-the-art results,

outperforming finetuned video models. It performs well on EPIC-Kitchens (verb and noun

EPIC-Kitchens

(val class-mean R@5; %)

Verb Noun Action

28.2 32.0 15.9

6.2 2.3 0.1

23.2 23.7 4.8

AVT+ [45]

LSMDC

(FiB test %)

Acc

52.9

0.1

31.0

MSR-VTT QA

(test acc %)

top1 top5

0.1 0.5

3.0 11.9

2.3 9.7

2.9 8.8

3.7 10.8

4.4 11.5

4.0 12.0

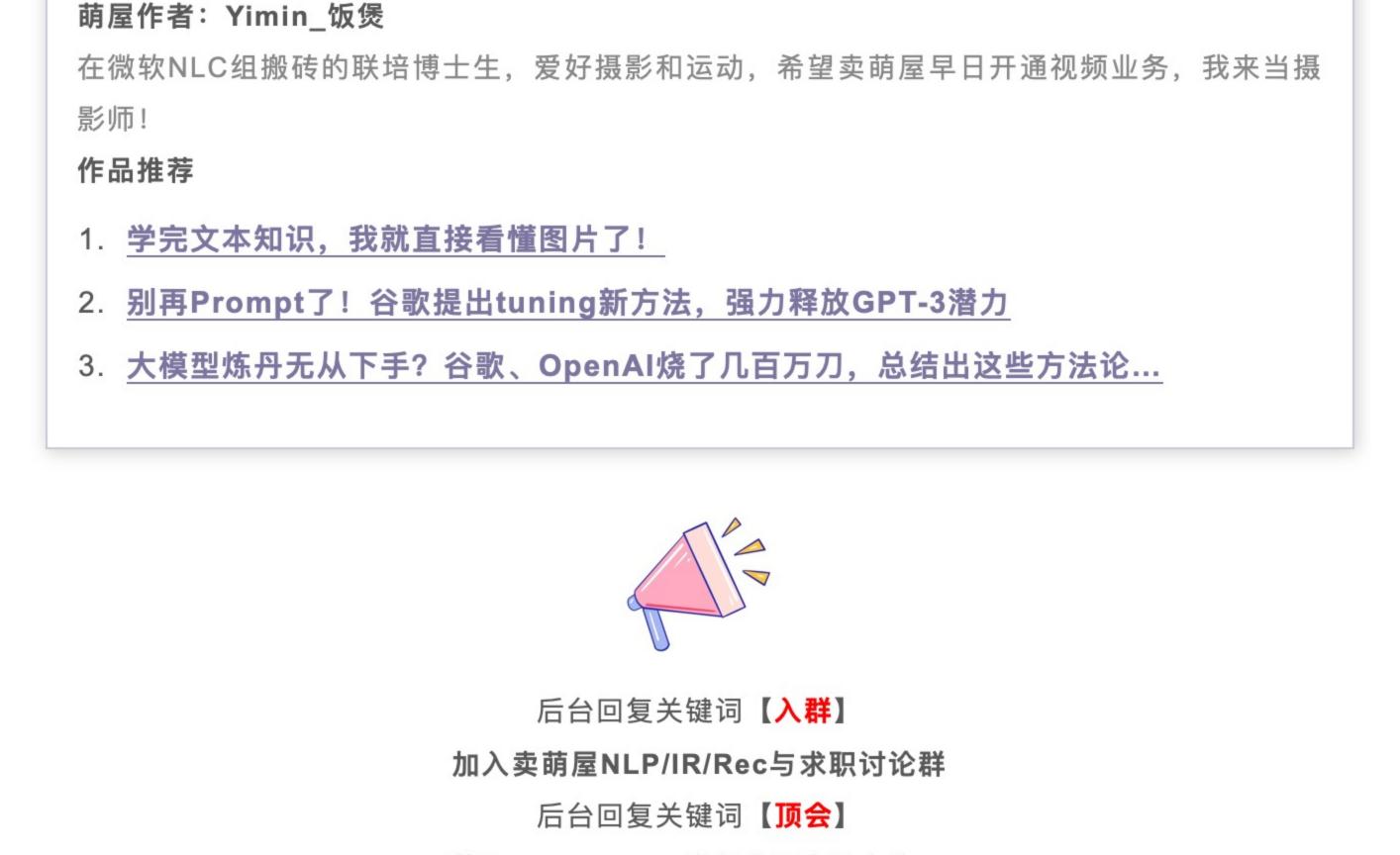
5.8 13.6

43.1

MERLOT [126]

forecasting), along with LSMDC, despite their long-tail distributions. On MSR-VTT QA, it outperforms past work on weakly-supervised video QA. Further, it outperforms CLIP (that cannot handle dynamic situations), and benefits from audio when given.
结语
多模态技术的发展和商业化,也许将会带来人工智能应用的新一轮爆发。以往的多模态应用面
临着模型架构复杂、缺少数据、缺少算力等一系列问题,而随着Transformer结构一统天下,
互联网上各模态数据的井喷式增长,计算资源越发普及,这些问题都在慢慢得到解决。
MERLOT系列工作刷新了我们认知中视频理解领域的上限,向我们展示了视频、语音、文本多

模态高效融合的一种可能性。未来,让我们一起努力朝着多模态领域的"BERT"模型进发吧!





喜欢此内容的人还喜欢 Nat. Mach. Intell. | MolCLR:一个用于分子表征学习的自监督框架 DrugAl 《Datawhale强化学习教程》出版了!

Datawhale