



微信扫一扫
关注该公众号



文 | 子龙

编 | 智商掉了一地

厉害了！作者将单一模型运用于三个不同领域的不同任务，结构简单且训练直观，还能有着出色的表现。

自Transformer横空出世，从NLP到CV，再到今天的多模态，无数基于Transformer的模型被应用于各类任务，似乎真的印证了当年文章的标题“Transformer is ALL you need”。然而，纯粹的NLP任务有BERT、RoBERTa，CV任务有ViT，多模态任务又有VLBERT、OSCAR，虽然都是基于Transformer的结构，但是仍然是针对不同任务设计不同模型，那么“万能”的Transformer能否构建出一个统合各类任务的模型，实现真的的一个模型解决所有问题呢？

今天文章的作者就关注到了当前各个模型的局限，提出了一个适用于NLP+CV+多模态的模型FLAVA，可运用于三种领域共计35个任务，且都有着出色的表现。

论文题目：

FLAVA: A Foundational Language And Vision Alignment Model

论文链接：

<https://arxiv.org/abs/2112.04482>

介绍

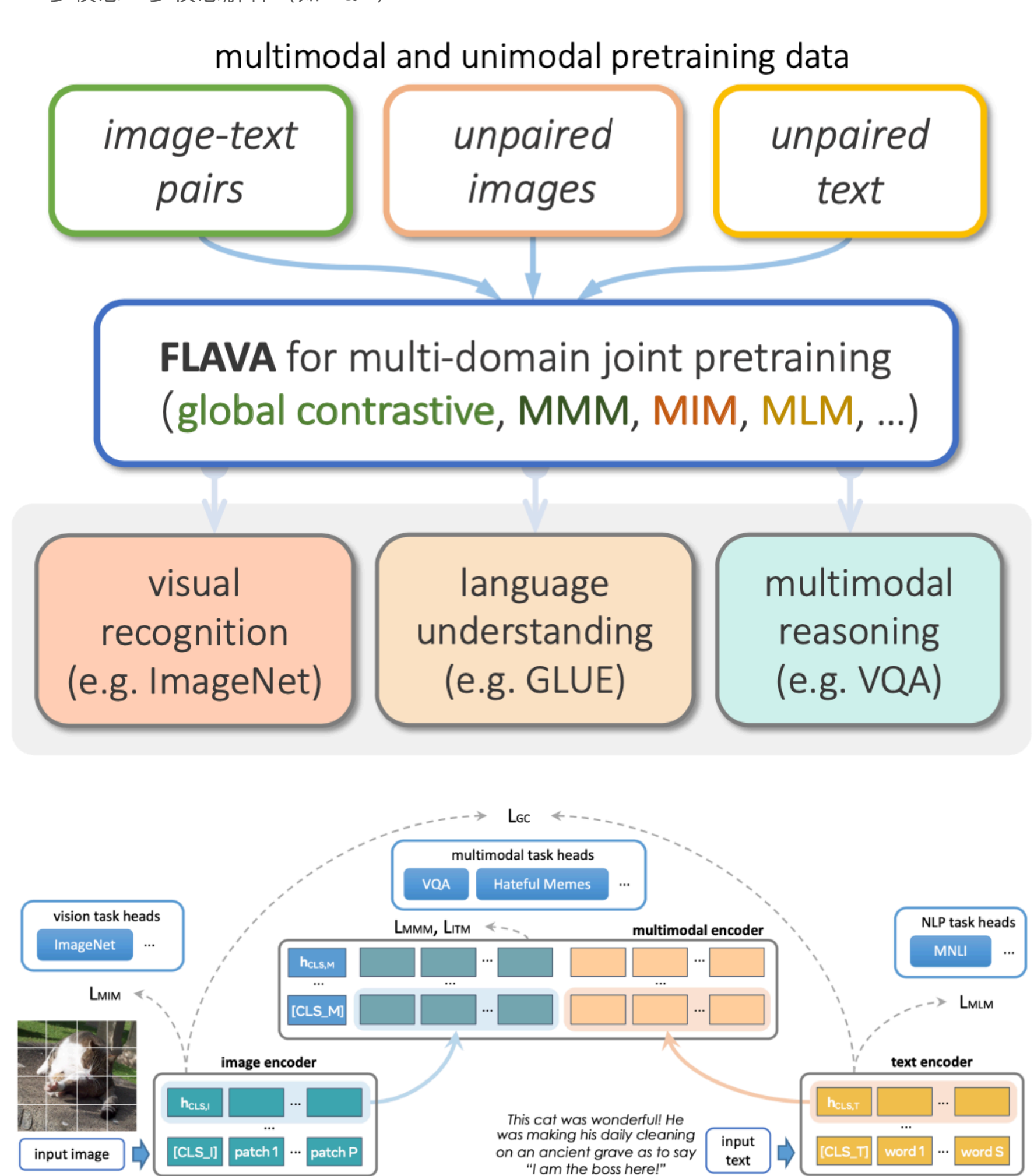
文章标题中，作者称模型为“Foundational”，他们不希望借助各种奇技淫巧的Tricks，而是通过尽可能简单的结构，配合直观的的训练手段，达到涵盖NLP、CV、多模态的目的。

FLAVA基于三种不同的输入：

- 匹配的图片-文本
- 单独文本
- 单独图片

解决三个领域的问题：

- NLP：语言理解（如GLUE）
- CV：视觉识别（如ImageNet）
- 多模态：多模态解释（如VQA）



图片编码器(Image Encoder)

FLAVA直接借用既有模型ViT的结构，同时仿照ViT的处理方法，分割图片进行编码。在ViT输出的隐状态上，FLAVA利用单一模态数据集中的图片进行Masked Image Modeling。首先，利用dVAE将图片转化为类似词向量的token；再参照BEiT，对masked隐状态进行分类，即利用周围图片分块，预测masked的图片属于dVAE划分的哪一类，这样在图片上也可以像BERT那样做mask modeling。

文本编码器(Text Encoder)

FLAVA在文本部分多处理就相对简单，作者采取常见的Masked Language Modeling，对一部分masked token进行预测，和其他方法区别在于，FLAVA没有采用BERT之类纯文本语言模型的结构，而是和图片编码器一样，使用了ViT的结构，不过因为是不同的模态，自然采用了不同的模型参数。

多模态编码器(Multimodal Encoder)

在图片编码器和文本编码器之上，FLAVA添加了一层多模态编码器做模态融合，多模态编码器将前两者输出的隐状态作为输入，同样利用ViT的模型结构进行融合。

多模态预训练

在文本编码器和图片编码器中，FLAVA在单一模态上进行了预训练，在多模态预训练方面，FLAVA使用了三种多模态预训练任务：

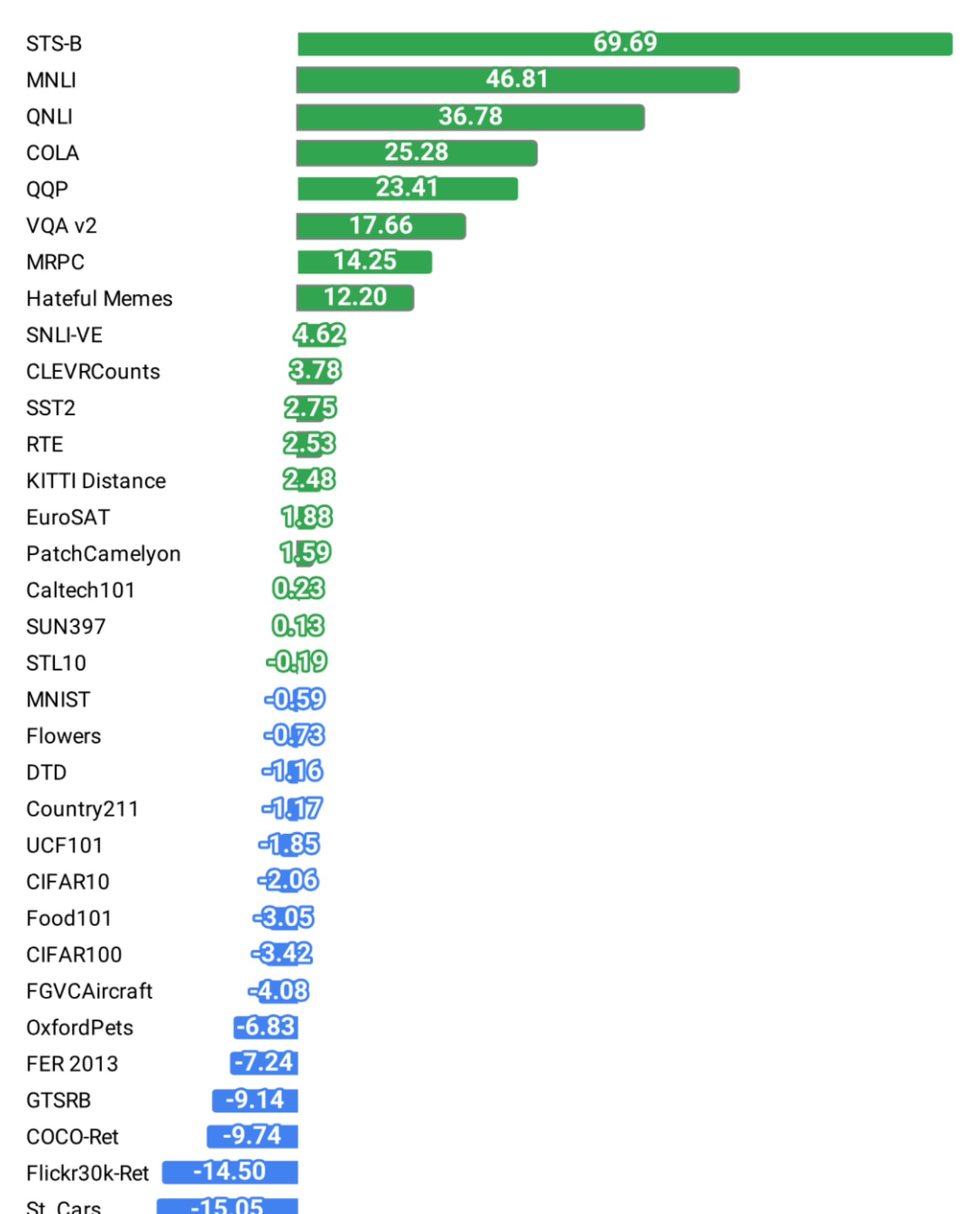
- 对比学习：FLAVA利用图片编码器和文本编码器的隐状态，增大相匹配的图片-文本对之间的余弦相似度，减小非匹配的图片-文本对之间的余弦相似度。
- Masked Multimodal Modeling：与图片编码器上的MIM类似，只不过改为利用多模态编码器的隐状态进行预测。
- 图片-文本匹配：与许多现有模型一样，FLAVA利用多模态编码器的[CLS]的隐状态，识别当前图片与文本是否匹配。

效果

从上述模型细节可以看出，无论是模型结构，还是预训练任务，文本与图片之间高度对称，同时也设计也十分直观。接下来看看在35个任务上的表现。

Datasets	Eval method	FMD	PMD	PMD	PMD	(PMD+IN+CCs+News+BC)	PMD	400M [83]		
MNLI [11]	fine-tuning	–	73.23	70.99	76.82	–	78.06	80.33	32.85	33.52
CoLA [110]	fine-tuning	–	39.55	17.58	38.97	44.22	50.65	11.02	24.65	11.02
MRPC [29]	fine-tuning	–	73.24	76.31	79.14	78.91	84.16	68.74	69.91	68.74
QQP [69]	fine-tuning	–	86.68	85.94	88.49	98.61	98.74	59.17	65.33	65.33
SST-2 [97]	fine-tuning	–	87.96	86.47	89.33	90.14	90.94	83.49	88.19	83.49
QNLI [88]	fine-tuning	–	82.32	71.85	84.77	86.40	87.31	49.46	50.54	50.54
RTE [7, 25, 36, 40]	fine-tuning	–	50.54	51.99	51.99	54.87	57.76	15.07	15.23	15.23
STS-B [1]	fine-tuning	–	78.89	57.27	84.29	83.21	85.67	13.70	15.98	15.98
NLP Avg.	–	–	71.55	64.80	74.22	75.55	78.19	46.44	50.50	50.50
Visual										
ImageNet [90]	linear eval	41.79	–	74.09	74.34	73.49	75.54	72.95	80.20	80.20
Food101 [11]	linear eval	53.30	–	87.77	87.53	87.39	88.51	85.49	91.56	91.56
CIFAR100 [53]	linear eval	76.30	–	93.44	92.37	92.63	92.87	91.25	94.93	94.93
CIFAR100 [58]	linear eval	55.57	–	78.37	78.01	76.49	77.68	74.40	81.10	81.10
Cars [56]	linear eval	14.71	–	72.12	72.07	66.81	70.87	62.84	85.92	85.92
Aircraft [74]	linear eval	13.83	–	49.74	48.90	44.73	47.31	40.02	51.40	51.40
DTD [20]	linear eval	55.53	–	76.86	76.91	75.80	77.29	73.40	78.46	78.46
Pets [79]	linear eval	34.48	–	84.98	84.93	82.77	84.82	79.61	91.66	91.66
Caltech101 [32]	linear eval	67.36	–	94.91	95.32	94.95	95.74	93.76	95.51	95.51
Flowers102 [76]	linear eval	67.23	–	96.36	96.39	95.58	96.70	94.94	97.12	97.12
MNIST [60]	linear eval	96.40	–	98.39	98.58	98.70	98.82	97.38	99.01	99.01
STL10 [21]	linear eval	80.12	–	98.06	98.31	98.32	98.89	97.29	99.09	99.09
EuroSAT [41]	linear eval	95.48	–	97.00	96.98	97.04	97.26	95.70	95.38	95.38
GTSRB [109]	linear eval	63.14	–	79.92	77.93	77.71	79.46	76.34	88.61	88.61
KITTI [35]	linear eval	86.03	–	87.83	88.84	88.70	89.04	84.89	86.56	86.56
PCAM [106]	linear eval	85.10	–	85.02	85.51	85.72	85.31	83.99	83.72	83.72
UCF101 [86]	linear eval	46.34	–	82.69	82.90	81.42	83.32	77.85	85.17	85.17
CLEVR [52]	linear eval	61.51	–	79.35	81.66	80.62	79.66	73.64	75.89	75.89
FER 2013 [38]	linear eval	50.98	–	59.96	60.87	58.99	61.12	57.04	68.36	68.36
SUN397 [113]	linear eval	52.45	–	81.27	81.41	81.05	82.17	79.96	82.05	82.05
SST [83]	linear eval	57.77	–	56.67	59.25	56.40	57.11	56.84	74.68	74.68
Country211 [83]	linear eval	8.87	–	27.27	26.75	27.01	28.92	25.12	30.10	30.10
Vision Avg.	–	57.46	–	79.14	79.35	78.29	79.44	76.12	82.57	82.57
Visual										
VQw2 [39]	fine-tuning	–	67.13	71.69	71.29	71.29	72.49	59.81	54.83	54.83
SUN397 [114]	fine-tuning	–	73.27	78.36	78.14	78.89	73.53	74.27		
Hateful Memes [53]	fine-tuning	–	55.58	70.72	72.45	76.09	56.59	63.93		
Flickr30K [81] TR R@1	zero-shot	–	68.30	69.30	64.50	67.70	60.90	82.20		
Flickr30K [81] TR R@5	zero-shot	–	91.50	92.90	90.30	94.00	88.90	96.60		
Flickr30K [81] IR R@1	zero-shot	–	60.56	63.16	60.04	65.22	56.48	62.08		
Flickr30K [81] IR R@5	zero-shot	–	86.68	87.70	86.46	89.38	83.60	85.68		
COCO [66] TR R@1	zero-shot	–	43.08	43.48	39.88	42.74	37.12	52.48		
COCO [66] TR R@5	zero-shot	–	75.82	76.76	72.84	76.76	69.48	76.68		
COCO [66] IR R@1	zero-shot	–	37.59	38.46	34.95	38.38	33.29	33.07		
COCO [66] IR R@5	zero-shot	–	67.28	67.68	64.63	67.47	62.47	58.37		
Multimodal Avg.	–	–	–	66.25	69.11	67.32	69.92	62.02	67.29	67.29
Macro Avg.	–	19.15	23.85	70.06	74.23	73.72	75.85	61.52	66.78	66.78


图中下划线表示最优结果，加粗表示在公开数据集上训练的最优结果。



从各个任务平均上看，FLAVA能够取得整体上的最优结果，多模态任务平均比CLIP高出2个百分点左右，整体平均比CLIP高出10个百分点左右。从具体任务上看，在不少任务上都取得了十分显著的提高，如STS-B数据集提高了69.69，MNLI数据集提高了46.81。


小结

不同于现有模型，FLAVA最大的特点，也可以说是创新点，在于作者实现了将单一模型运用于三个不同领域的不同任务，而且都有着不错的效果，虽然FLAVA并没有奇迹般在所有任务上都达到SOTA，但是整体性能上并不弱于现有模型，同时有着更广阔的运用场景，模型设计也没有各种奇技淫巧，这对未来研究通用模型有着很大的启发。



萌屋作者：子龙(Ryan)
本科毕业于北大计算机系，曾混迹于商汤和MSRA，现在是宅在UCSD(Social-Dead)的在读PhD，主要关注多模态中的NLP和数据 mining，也在探索更多有意思的Topic，原本只是公众号的吃瓜群众，被各种有趣的推送吸引就土子赋船，希望借此沾沾小屋的灵气，paper++，早日成为有猫的程序员！

作品推荐：
1 [别再用搞纯文本了！多模态文档理解更被时代需要！](#)
2 [Transformer哪家强？Google爸爸择优良！](#)
3 [预训练语言真的是世界模型？](#)



后台回复关键词【**入群**】
加入卖萌屋NLP/IR/Rec 求职讨论群
后台回复关键词【**顶会**】
获取ACL、CIKM等各大顶会论文集！



FOLLOW ME



STAR ME

喜欢此内容的人还喜欢

如何在自动驾驶的视觉感知中检测极端情况？
3D视觉工坊

致敬ConcConv！英特尔提出即插即用的“万金油”动态卷积ODConv
AI科技评论

模型大十倍，性能提升几倍？谷歌研究员进行了一番研究
墨创AI

