

别再Prompt了！谷歌提出tuning新方法，强力释放GPT-3潜力！

原创 Yimin_饭煲 夕小瑶的麦萌屋 2021-09-07 22:20



微信扫一扫
关注公众号

收录于话题
#麦萌屋@自然语言处理

69个 >



如果评选NLP圈的2020年度十大关键词，那么GPT-3（Language Models are Few shot Learners）一定榜上有名。

GPT-3庞大的参数量，优异的性能至今仍让圈内圈外人都津津乐道，而OpenAI发布的OpenAI API，更是为自然语言处理技术的大规模可扩展商业应用提供了一个极有前景的方向。不过，作为NLP研究者，我认为GPT-3对前沿研究的最大贡献是，展现了 Prompt-tuning 技术在通用任务（特别是零样本和小样本场景下）上的应用潜力。在GPT-3之前，Prompt-tuning 大多仅被用来探索语言模型中蕴藏的世界知识，而GPT-3之后，Prompt-tuning 就“登堂入室”，被用到了各种类型的NLP任务上（甚至还有多模态任务、代码分析任务），成为了近两年来发Paper的一个热点。不太熟悉 Prompt-tuning 的读者可以参考CMU最新发布在arxiv上的综述（Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing）。

麦萌屋之前也转载过该论文的中文解析《Fine-tune之后的NLP新范式：Prompt越来越火，CMU华人博士后出了篇综述文章》。

Prompt-tuning 和GPT-3互相成就，成为了NLP发展历史中不可忽略的里程碑。基于Prompt-tuning 让GPT-3处理各种类型的任务，固然取得了不错的表现，但仔细想想，GPT-3这样一个创造了各种奇迹的巨大模型，是否还有更大的零样本和小样本学习能力 尚未挖掘？Prompt-tuning 一定是利用GPT-3的最好方式吗？Google的研究员们不甘心止步于此，提出了 instruction-tuning，利用比GPT-3更少的参数量，在25个任务中的19个上显著超越GPT-3，告诉世界：GPT-3，还能更强！

论文题目：

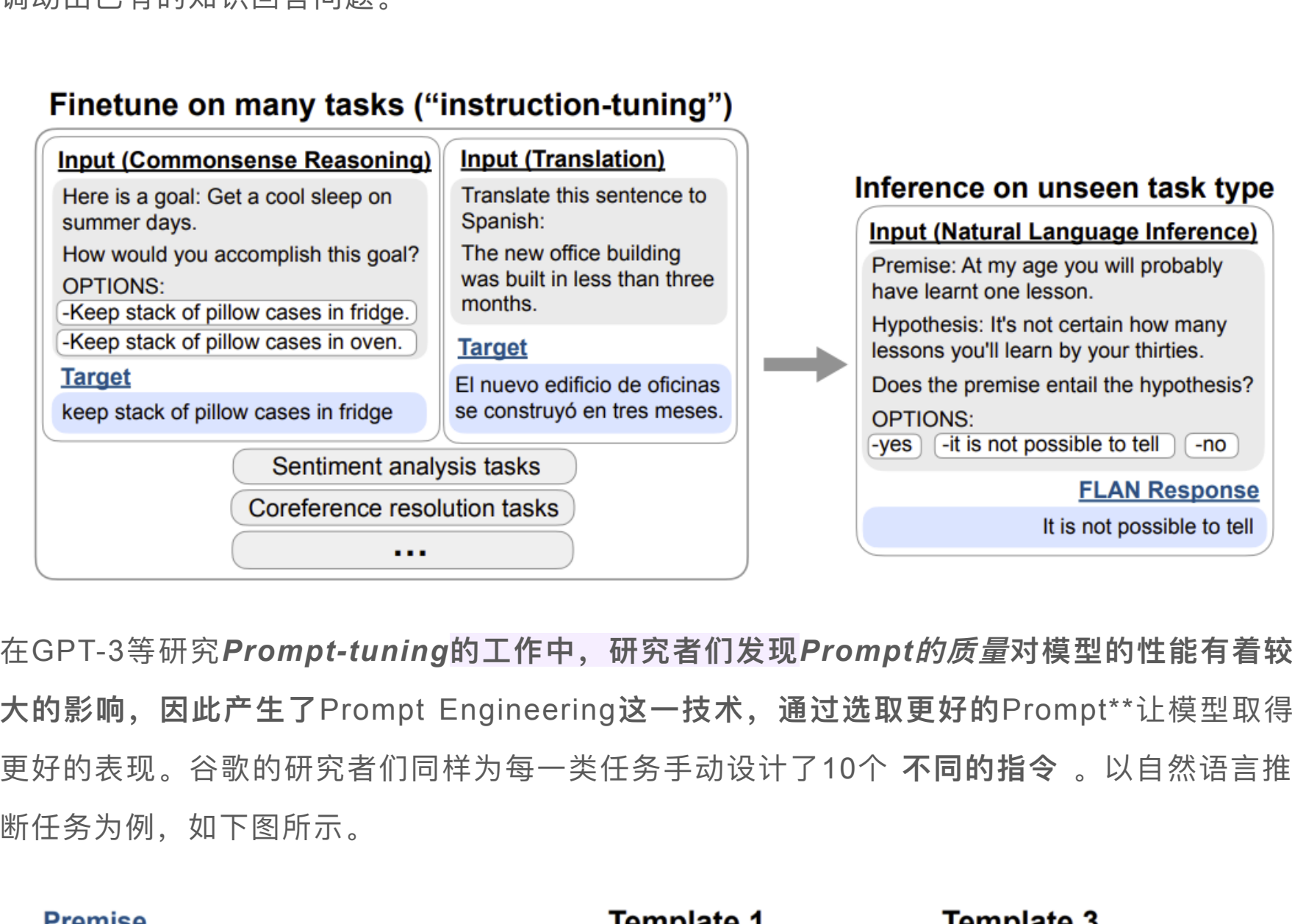
Finetuned Language Models Are Zero-Shot Learners

论文链接：

<https://arxiv.org/abs/2109.01652>

方法：FLAN

谷歌的研究员们将自己的方法取名为FLAN（Finetuned LANGUAGE Models are zero-shot Learners），相比于GPT-3（LANGUAGE Models are zero-shot Learners），区别在于Finetune。FLAN的核心思想是，当面对给定的任务A时，首先将模型在大量的其他不同类型的任务上进行微调，微调的方式是将任务的指令与数据进行拼接（可以理解作为一种Prompt），随后给出任务A的指令，直接进行推断。具体示例可见下图。



在GPT-3等研究Prompt-tuning的工作中，研究者们发现Prompt的质量对模型的性能有着较大的影响，因此产生了 Prompt Engineering 这一技术，通过选取更好的Prompt“让模型取得更好的表现。谷歌的研究者们同样为每一类任务手动设计了10个不同的指令，以自然语言推断任务为例，如下图所示。



实验结果

实验设置

作者们选取了12类共计62个常见的自然语言处理和生成任务开展实验：

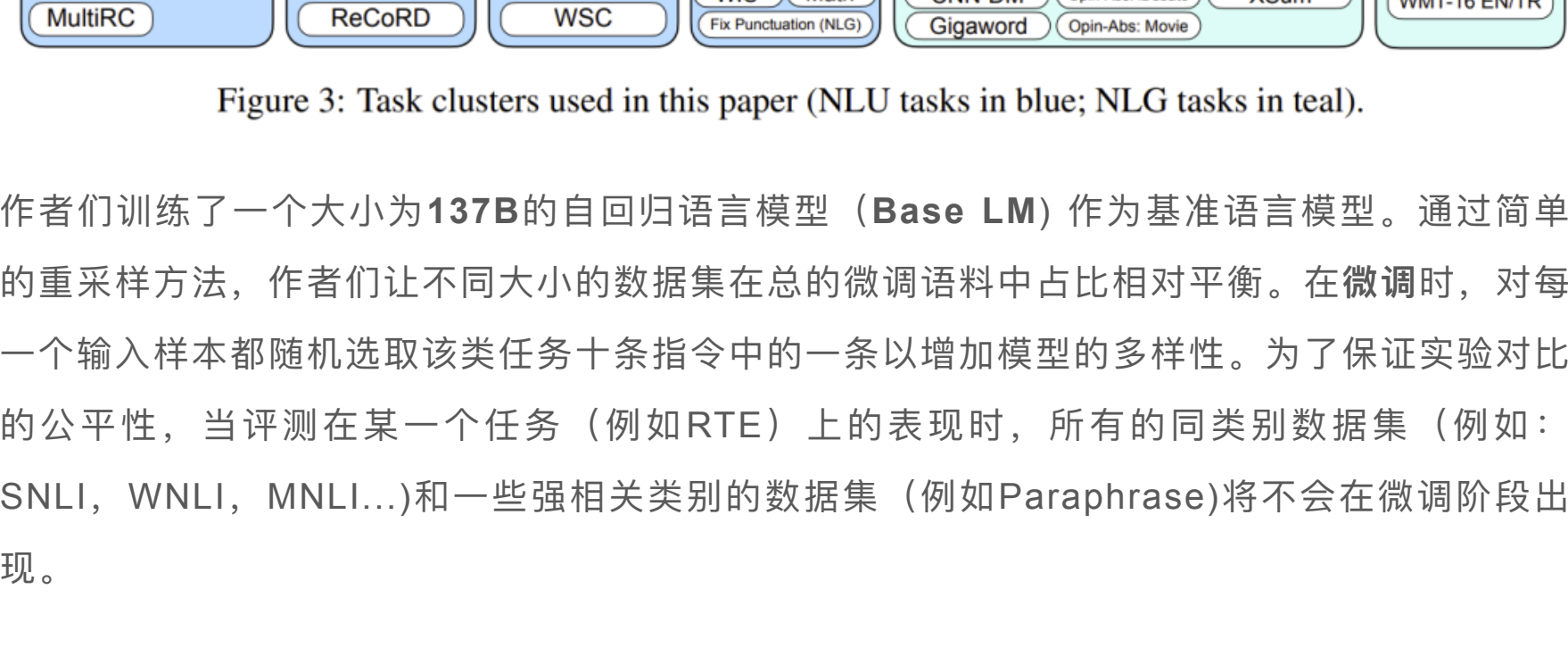


Figure 3: Task clusters used in this paper (NLU tasks in blue; NLG tasks in teal).

作者们训练了一个大小为137B的自回归语言模型（Base LM）作为基准语言模型。通过简单的重采样方法，作者们让不同大小的数据集在总的微调语料中占比相对平衡。在微调时，对每一个输入样本都随机选取该类任务十条指令中的一条以增加模型的多样性。为了保证实验对比的公平性，当评测在某一个任务（例如RTE）上的表现时，所有的同类别数据集（例如：SNLI, WNLI, MNL...）和一些强相关类别的数据集（例如Paraphrase）将不会在微调阶段出现。

作者们使用了T5-11B和GPT-3作为基线模型。对于FLAN方法，作者同时给出了在目标任务上选择随机指令（no prompt engineering）和在目标任务验证集上最优指令（best dev template）

结果

作者们发现，FLAN这一方法在与指令更相关的任务上表现更好（例如自然语言推断，问答），而在与常识更相关的任务上表现较为普通。

在自然语言推断任务和问答任务上，FLAN在零样本场景下就已经超过了小样本GPT-3的效果，在许多任务上甚至与有监督模型达到了相当的表现。

	NATURAL LANGUAGE INFERENCE				
	ANLI-R1 acc.	ANLI-R2 acc.	ANLI-R3 acc.	CB acc.	RTE acc.
Supervised model	57.4 [±]	48.3 [±]	43.5 [±]	96.8 [±]	92.5 [±]
Base LM 137B zero-shot	39.6	39.9	39.3	42.9	73.3
- few-shot	39.0	37.5	40.7	44.8	70.8
GPT-3 175B zero-shot	34.6	35.4	34.5	46.4	58.9
- few-shot	36.8	34.0	40.2	82.1	70.4
FLAN 137B zero-shot					
- no prompt engineering	47.7 [±] _{std=1.4}	43.9 [±] _{std=1.5}	47.0 [±] _{std=1.5}	64.1 [±] _{std=1.7}	78.3 [±] _{std=1.9}
- best dev template	46.4 [±] _{std=1.6}	44.4 [±] _{std=1.0}	48.5 [±] _{std=1.3}	83.9 [±] _{std=1.8}	84.1 [±] _{std=1.9}

Table 1: Results on natural language inference. For FLAN, we report both the average of up to ten templates (proxying the expected performance without prompt engineering), as well as the test set performance of the template that had the highest performance on the dev set. The triangle [±] indicates improvement over few-shot GPT-3. The up-arrow [±] indicates improvement only over zero-shot GPT-3. *T5-11B, *BERT-large.

	READING COMPREHENSION			OPEN-DOMAIN QA		
	BoolQ acc.	MultiRC F1 acc.	OBQA acc.	ARC-c acc.	ARC-e acc.	TriviaQA EM acc.
Supervised model	91.2 [±]	88.2 [±]	85.4 [±]	92.6 [±]	81.1 [±]	36.6 [±]
Base LM 137B zero-shot	81.0	60.0	41.8	76.4	42.0	3.2
- few-shot	79.7	59.6	40.6	80.9	49.4	22.1
GPT-3 175B zero-shot	80.0	72.9	57.6	68.8	51.4	14.6
- few-shot	77.5	74.8	65.4	70.1	51.5	29.9
FLAN 137B zero-shot						
- no prompt engineering	80.2 [±] _{std=1.7}	74.5 [±] _{std=1.4}	77.4 [±] _{std=1.0}	79.5 [±] _{std=1.6}	61.7 [±] _{std=1.7}	18.6 [±] _{std=1.3}
- best dev template	82.9 [±] _{std=1.4}	77.5 [±] _{std=1.7}	78.4 [±] _{std=1.0}	79.6 [±] _{std=1.7}	63.1 [±] _{std=1.6}	20.7 [±] _{std=1.3}

Table 2: Results on reading comprehension and open-domain question answering. For FLAN, we report both the average of up to ten templates (proxying the expected performance without prompt engineering), as well as the test set performance of the template that had the highest performance on the dev set. The triangle [±] indicates improvement over few-shot GPT-3. The up-arrow [±] indicates improvement only over zero-shot GPT-3. *T5-11B, *BERT-large.

	COMMONSENSE REASONING				COREFERENCE		
	CoPA acc.	HellaSwag acc.	PiQA acc.	StoryCloze acc.	ReCoRD acc.	WSC273 acc.	Winogrande acc.
Supervised model	94.8 [±]	47.3 [±]	66.8 [±]	89.2 [±]	93.4 [±]	72.2 [±]	93.8 [±]
Base LM 137B zero-shot	90.0	57.0	80.3	79.5	87.8	81.0	68.3
- few-shot	89.6	58.8	80.2	83.7	87.6	81.5	68.4
GPT-3 175B zero-shot	91.0	78.9	81.0	83.2	90.2	88.3	70.2
- few-shot	92.0	79.3	82.3	87.7	89.0	88.6	77.7
FLAN 137B zero-shot							
- no prompt engineering	90.6 [±] _{std=1.0}	56.4 [±] _{std=0.5}	80.9 [±] _{std=0.8}	92.2 [±] _{std=1.5}	67.8 [±] _{std=1.3}	80.8 [±] _{std=1.7}	67.3 [±] _{std=1.3}
- best dev template	91.0 [±] _{std=1.0}	56.7 [±] _{std=0.5}	80.5 [±] _{std=0.8}	93.4 [±] _{std=1.7}	72.5 [±] _{std=1.3}	-	71.2 [±] _{std=1.0}

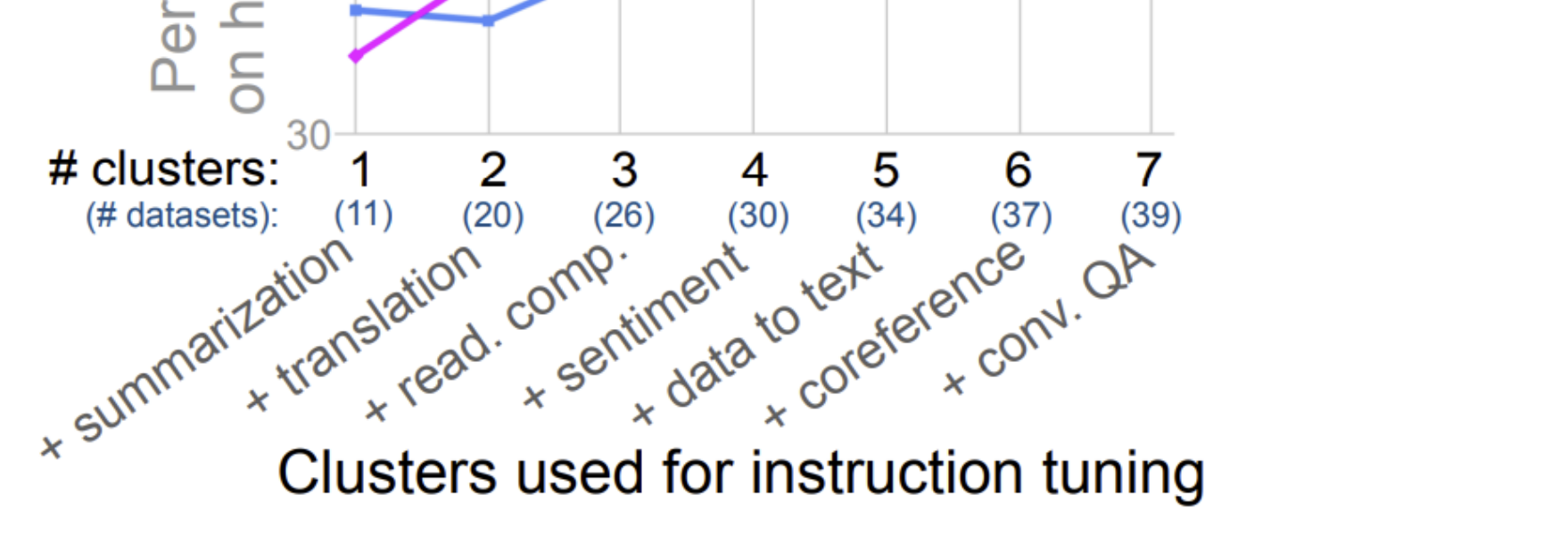
Table 3: Results (accuracy in %) for commonsense reasoning and coreference resolution. For FLAN, we report both the average of up to ten templates (proxying the expected performance without prompt engineering), as well as the test set performance of the template that had the highest performance on the dev set. *T5-11B, *BERT-large. The triangle [±] indicates improvement over few-shot GPT-3. The up-arrow [±] indicates improvement only over zero-shot GPT-3.

在翻译任务上，零样本场景下的FLAN明显优于GPT-3，但相比于小样本GPT-3的表现仍然有差距。

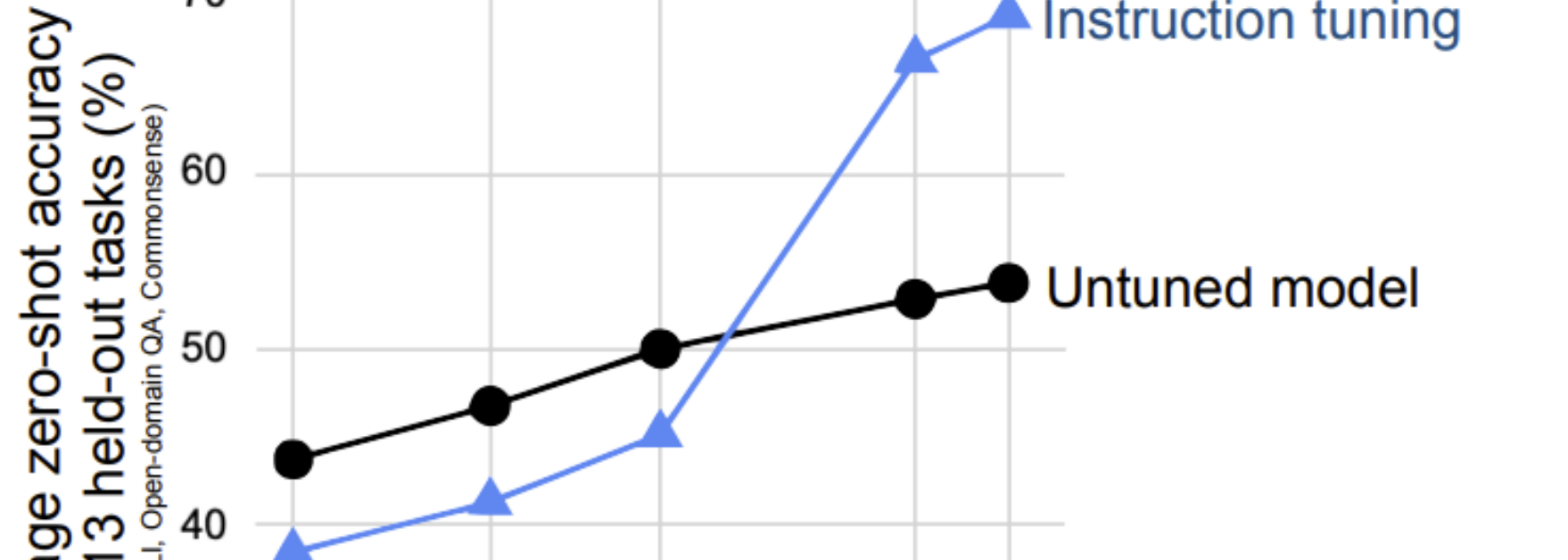
	TRANSLATION					
	French		German		Romanian	
	En→Fr BLEU	Fr→En BLEU	En→De BLEU	De→En BLEU	En→Ro BLEU	Ro→En BLEU
Supervised model	45.6 [±]	35.0 [±]	41.2 [±]	38.6 [±]	38.5 [±]	39.9 [±]
Base LM 137B zero-shot	11.2	7.2	7.7	20.8	3.5	9.7
- few-shot	31.5	34.7	26.7	36.8	22.9	37.5
GPT-3 175B zero-shot	25.2	21.2	24.6	27.2	14.1	19.9
- few-shot	32.6	39.2	29.7	40.6	21.0	39.5
FLAN 137B zero-shot						
- no prompt engineering	32.0 [±] _{std=0.3}	35.6 [±] _{std=1.4}	24.2 [±] _{std=0.7}	39.4 [±] _{std=1.2}	16.9 [±] _{std=1.3}	36.1 [±] _{std=1.0}
- best dev template	34.0 [±] _{std=1.4}	36.5 [±] _{std=1.5}	27.0 [±] _{std=1.2}	39.8 [±] _{std=1.6}	18.4 [±] _{std=1.3}	36.7 [±] _{std=1.7}

Table 4: Translation results (BLEU) for WMT 14 En/Fr and WMT 16 En/De and En/Ro. For FLAN, we report both the average of up to ten templates (proxying the expected performance without prompt engineering), as well as the test set performance of the template that had the highest performance on the dev set. *EduNov et al. (2018), *Durrani et al. (2014), *Wang et al. (2019b), *Sennrich et al. (2016), *Liu et al. (2020). The triangle [±] indicates improvement over few-shot GPT-3. The up-arrow [±] indicates improvement only over zero-shot GPT-3.

作者们还研究了增加指令微调阶段任务的数目对FLAN模型效果的影响。结果表明，随着指令微调任务数目的增加，模型在各种任务上都能够取得更好的表现。



作者们同时研究了模型大小对FLAN模型效果的影响，一个有趣的现象是，当模型的大小较小时，指令微调反而会模型的表现变差。作者认为的原因时，当模型较小时，在大量的任务上做指令微调会“填满”模型的容量，损害模型的泛化能力，使得模型在新任务上表现较差。



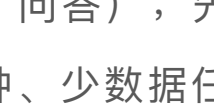
结语

熟悉NLP相关领域的同行们，也许认为这篇文章又是一篇“A+B”的工作（A= Prompt-tuning, B= Multi-task Learning）。基于Prompt的工作正值大热的时期（你敢相信九月的第一周就有四篇和Prompt有关的NLP论文挂Arxiv吗~），而通过在不同种类的微调任务上多任务学习提升性能也并不新颖，例如早期Microsoft的工作MT-DNN，Facebook的工作MUPPET。

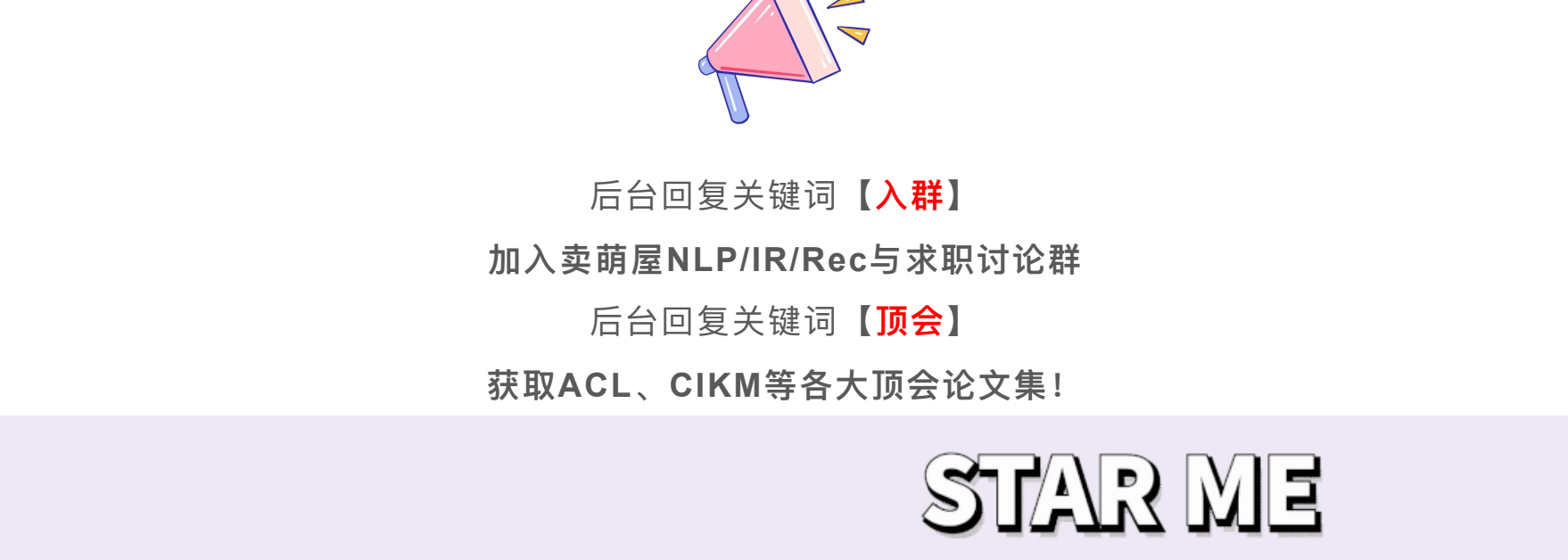
不过，笔者认为，这样的A+B，或许是未来通用自然语言处理模型的一个可能的解决方案。

首先通过大量的无标记语料训练千亿参数级别的大规模自回归预训练模型，第二步，通过设计指令（Instruction Tuning）的方式让这样的模型能够对理解和生成任务进行微调。在微调的过程中可以采用类似于课程学习的方式，先学习底层的任务（如命名实体识别，序列语义标注），再学习上层的任务（如逻辑推理，问答）；先学习资源丰富的任务（如英语/大数据任务），再学习资源较少的任务（如小语种、少数族任务），并利用适配器(Adapter)保留模型中任务专用的部分。最后，给出指令让模型面对新数据、新任务进行推理。

这样通用性更强的工作应该不会太远，也许资源丰富的大厂们已经在搞了呢~



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



喜欢此内容的人还喜欢

若被制裁，中国AI会雪崩吗？

夕小瑶的麦萌屋



21 Finetuned Language Models Are Zero-Shot Learners
<https://arxiv.org/abs/2109.01652>

22 Multi-Task Deep Neural Networks for Natural Language Understanding
<https://arxiv.org/abs/1919.14411>

23 Muppet: Massive Multi-Task Representations with Pre-Finetuning
<https://arxiv.org/abs/2101.11038>