

文 | Harris

刘鹏飞博士将近代NLP的研究划归为四种范式 [1] 并把预训练语言模型加持下的Prompt Learning看作是近代自然语言处理技术发展的“第四范式”。当我们使用新范式的方法的时候，能够意识到它带来的优异性可能是以某种“人力”牺牲为代价的。而如何让这种人力代价降到最低，往往就是新范式里需要解决的核心问题 [2]。Prompt Learning刚兴起之时，prompts大多多是人工设计的，为了减少人工，后来涌现出一系列用自动化方式获取prompts的研究工作。Soft prompts由于易于使用、易理解性和优异的性能在近一年来获得了广泛的关注。今年四月份，谷歌提出Prompt-tuning [3]，其为每个任务分配一个可训练的soft-prompt并保持预训练模型的参数不变，发现在使用较大的预训练模型时，Prompt-tuning可以媲美微调整个预训练模型（Model-tuning）的性能。虽然Prompt-tuning非常有效，其仍然面临如下三个问题：

- 在使用较小的预训练模型时，Prompt-tuning的表现仍旧与Model-tuning的表现有明显的差距
- 在Few-shot的场景下，Prompt-tuning的表现并不理想
- 相比于Model-tuning，Prompt-tuning的收敛速度较慢，需要训练更多的步数

幸运的是，这些问题都可以通过更好地初始化soft prompts解决！一个直观的想法是可以通过预训练模型词汇表里的一些词来初始化，比如对于分类任务就可以用类别标签对应的词来初始化。但是类别标签词数目有限而且并没有很多任务相关的信息。

一个更好的方法是：用在源任务上训练的soft prompts来初始化且标注任务的soft prompts！这个想法有点类似迁移学习的意思。迁移学习通俗来讲，就是用已有的知识来学习新的知识。核心是找到已有知识和新知识之间的相似性。迁移prompts也类似，如果源任务与目标任务越相似，那么迁移的效果可能就越好。针对prompts的可迁移性问题，谷歌和清华的研究员们进行了初步的探索，笔者接下来为大家一一解读。想要快速浏览 takeaway 干货的读者，可以直接移步文末总结部分。

相关论文：

- PPT: Pre-trained Prompt Tuning for Few-shot Learning
<https://arxiv.org/abs/2109.04332>
- SPoT: Better Frozen Model Adaptation through Soft Prompt Transfer
<https://arxiv.org/abs/2110.07904>
- On Transferability of Prompt Tuning for Natural Language Understanding
<https://arxiv.org/abs/2111.06719>

PPT: Pre-trained Prompt Tuning for Few-shot Learning

清华研究人员发现在Few-shot的场景下，即便使用非常大的预训练模型，相比于Model-tuning，Prompt-tuning的表现仍要差得多。受预训练模型的启发，他们想对prompts也进行预训练！他们关注的是分类任务，并把分类任务分成了三类：单句分类任务、句对分类任务和多项选择分类任务。之后，他们针对每一类都设计了一个自监督预训练任务，然后用预训练任务的prompts去初始化下游任务的prompts。此外，他们还提供了一个unified的版本，将所有的分类任务都看作是多项选择分类任务。

English Tasks									
		SST-2	SST-5	RACE-m	RACE-h	BoolQ	KTE	CB	
		Acc.	Acc.	Acc.	Acc.	Acc.	Acc.	Acc.	F1
FT (11B)	TS-Small	-	72.8 _{±1}	31.1 _{±4}	26.4 _{±6}	26.3 _{±5}	59.2 _{±6}	54.0 _{±7}	70.1 _{±6}
	TS-Base	-	74.6 _{±7}	28.8 _{±2}	27.2 _{±5}	26.7 _{±2}	61.9 _{±1}	56.1 _{±2}	70.4 _{±9}
	TS-Large	-	89.1 _{±2}	42.4 _{±2}	48.2 _{±8}	43.2 _{±7}	74.6 _{±9}	64.4 _{±4}	82.3 _{±2}
	TS-XL	-	89.6 _{±2}	38.4 _{±1}	55.0 _{±8}	50.9 _{±6}	77.2 _{±1}	62.3 _{±8}	81.9 _{±9}
PT (410K)	TS-XXL	-	91.4 _{±4}	40.6 _{±2}	62.9 _{±9}	54.8 _{±9}	80.8 _{±4}	64.1 _{±2}	86.5 _{±3}
	Vanilla PT	70.5 _{±5}	32.3 _{±3}	34.7 _{±2}	31.6 _{±3}	61.0 _{±3}	53.5 _{±5}	50.7 _{±1}	
	Hybrid PT	87.6 _{±6}	40.9 _{±2}	53.5 _{±2}	44.2 _{±4}	79.8 _{±5}	56.8 _{±6}	66.5 _{±2}	
	LM Adaptation	77.6 _{±9}	36.2 _{±2}	27.3 _{±2}	26.5 _{±4}	62.0 _{±3}	55.3 _{±9}	61.2 _{±7}	
	PPT	93.5 _{±3}	50.2 _{±2}	60.0 _{±2}	53.0 _{±2}	66.4 _{±7}	58.9 _{±7}	71.2 _{±2}	
	Hybrid PPT	93.8 _{±1}	50.1 _{±2}	62.5 _{±9}	52.2 _{±7}	82.0 _{±0}	59.8 _{±2}	73.2 _{±9}	
	Unified PPT	94.4 _{±3}	46.0 _{±3}	58.0 _{±9}	49.9 _{±3}	76.0 _{±7}	65.8 _{±3}	82.2 _{±4}	

Chinese Tasks									
		ChnSenti	Amazon	CCPM	C ³	LQCMC	CMNLI	OCNLI	
		Acc.	Acc.	Acc.	Acc.	Acc.	Acc.	Acc.	Acc.
FT (11B)	mTS-Small	-	76.1 _{±6}	29.9 _{±9}	31.9 _{±2}	29.6 _{±5}	52.4 _{±5}	36.5 _{±2}	34.9 _{±3}
	mTS-Base	-	78.2 _{±6}	36.4 _{±9}	40.4 _{±8}	29.4 _{±6}	50.9 _{±0}	36.3 _{±5}	35.4 _{±6}
	mTS-Large	-	79.1 _{±6}	31.0 _{±1}	46.0 _{±8}	29.5 _{±8}	52.1 _{±6}	35.8 _{±2}	35.2 _{±1}
	mTS-XL	-	82.7 _{±6}	35.5 _{±1}	68.3 _{±1}	29.7 _{±2}	52.9 _{±6}	36.8 _{±1}	35.6 _{±3}
	mTS-XXL	-	83.6 _{±5}	42.1 _{±8}	79.7 _{±1}	37.2 _{±3}	53.1 _{±0}	39.0 _{±4}	37.4 _{±2}
	CPM-2	-	86.1 _{±8}	42.5 _{±2}	81.8 _{±1}	38.4 _{±1}	58.8 _{±1}	40.7 _{±0}	38.5 _{±2}
PT (410K)	Vanilla PT	62.1 _{±1}	30.3 _{±4}	31.0 _{±7}	28.2 _{±4}	51.5 _{±4}	35.4 _{±5}	37.0 _{±3}	
	Hybrid PT	79.2 _{±0}	39.1 _{±4}	46.6 _{±10}	29.2 _{±5}	54.6 _{±3}	37.1 _{±6}	37.8 _{±4}	
	LM Adaptation	74.3 _{±2}	35.2 _{±4}	83.7 _{±28}	30.2 _{±5}	51.4 _{±9}	35.1 _{±3}	38.0 _{±1}	
	PPT	90.1 _{±8}	48.6 _{±1}	85.2 _{±8}	43.8 _{±2}	59.1 _{±6}	43.0 _{±8}	40.1 _{±4}	
	Hybrid PPT	89.5 _{±9}	48.8 _{±0}	83.5 _{±9}	46.0 _{±5}	67.3 _{±9}	41.3 _{±8}	38.7 _{±9}	
	Unified PPT	90.7 _{±3}	44.6 _{±1}	83.4 _{±9}	50.2 _{±6}	55.0 _{±4}	40.6 _{±4}	41.5 _{±8}	

从上图可以看到，在Few-shot的场景下，PPT比Prompt-tuning（图中Vanilla PT）的表现明显要好。相比于Model-tuning（图中FT），PPT在所有中文任务上都取得了更好的表现，并在英文任务上取得了类似的性能。上述结果充分显示出预训练prompts在Few-shot场景下的有效性。但是当训练数据增多时，其优势就会越来越小，见下图

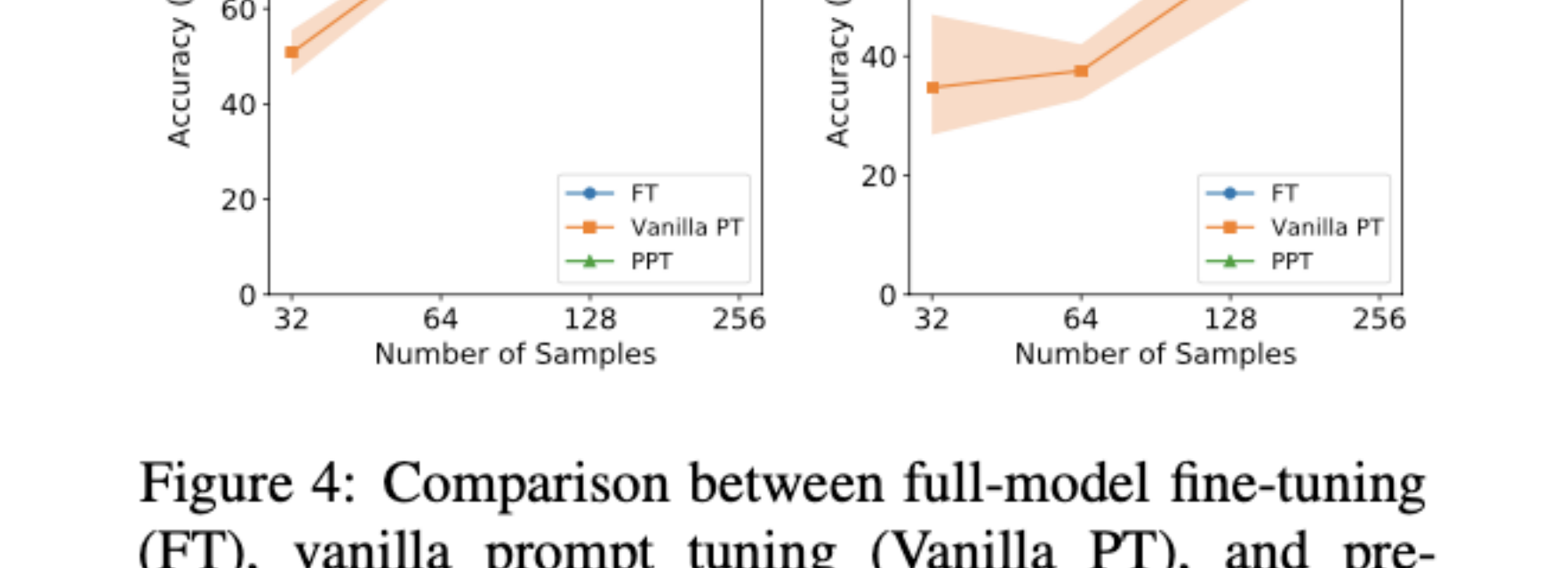
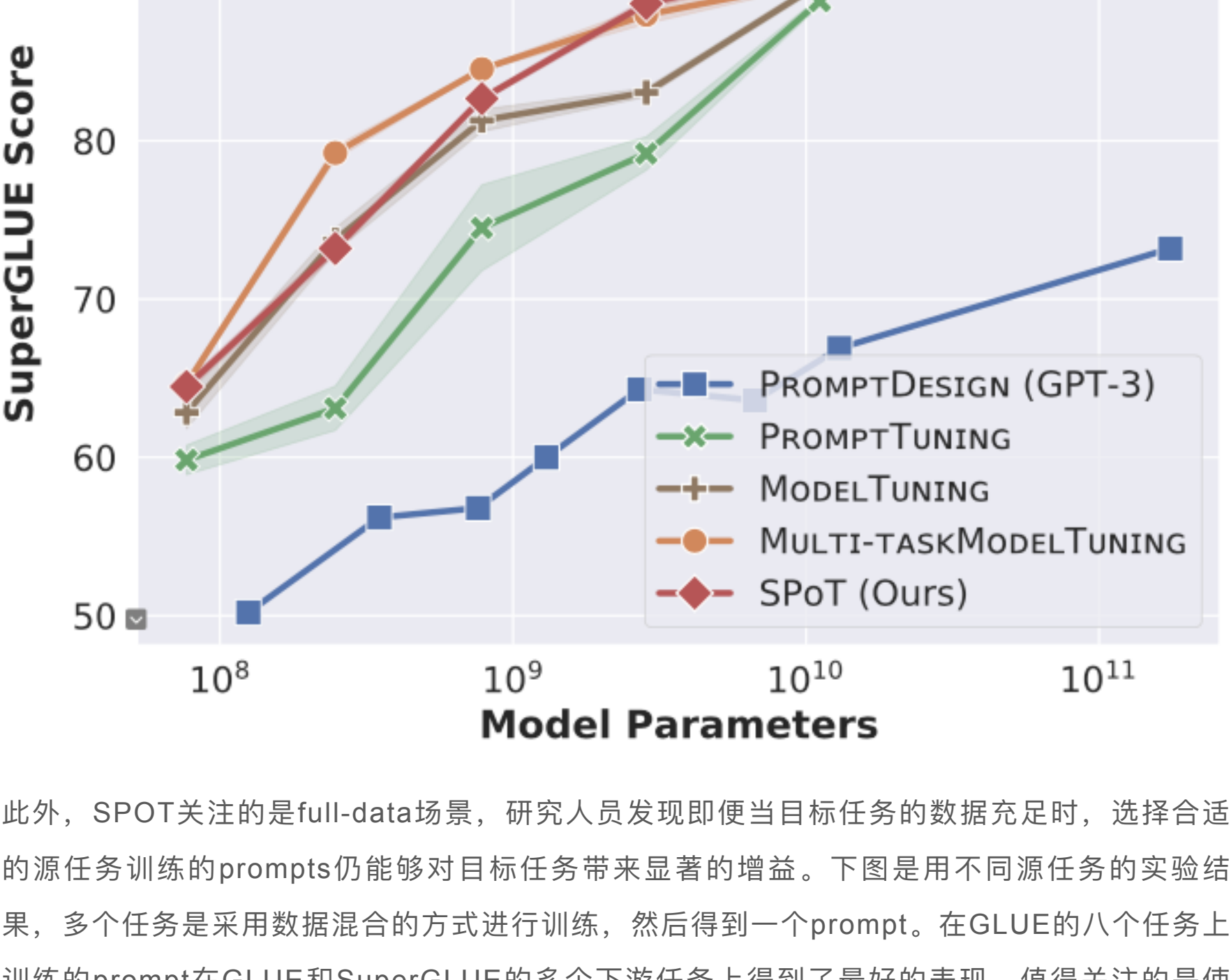


Figure 4: Comparison between full-model fine-tuning (FT), vanilla prompt tuning (Vanilla PT), and pre-trained prompt tuning (PPT) when different numbers of training samples are available. For the small number of training samples, PPT is consistently the best. When the number grows, the performance of these methods becomes closer.

SPoT: Better Frozen Model Adaptation through Soft Prompt Transfer

不同于PPT需要人工设计预训练任务，SPoT用一个或多个源任务训练的prompts来初始化目标任务prompts。SPoT在预训练模型比较小时，仍可以达到与Model-tuning相近的性能，显著超越Prompt-tuning，见下图



此外，SPoT关注的是full-data场景，研究人员发现即便当目标任务的数据充足时，选择合适的源任务训练的prompts仍能够目标任务带来显著的增益。下图是用不同源任务的实验结果，多个任务是采用数据混合的方式进行训练，然后得到一个prompt。在GLUE的八个任务上训练的prompt在GLUE和SuperGLUE的多个下游任务上得到了最好的表现，值得关注的是使用单个任务如MNLI和SQUAD训练的prompt也带来了明显的提升。

Method	GLUE	SuperGLUE
BASELINE		
PROMPTTUNING	81.2 _{±0.4}	66.6 _{±0.2}
— longer tuning	78.4 _{±1.7}	63.1 _{±1.1}
SPoT with different source mixtures		
GLUE (8 tasks)	82.8 _{±0.2}	73.2 _{±0.3}
— longer tuning	82.0 _{±0.2}	70.7 _{±0.4}
C4	82.0 _{±0.0}	67.7 _{±0.3}
MNLI	82.5 _{±0.0}	72.6 _{±0.8}
SQuAD	82.2 _{±0.1}	72.0 _{±0.4}
SuperGLUE (8 tasks)	82.0 _{±0.1}	66.6 _{±0.2}
NLI (7 tasks)	82.6 _{±0.1}	71.4 _{±0.2}
Paraphrasing/similarity (4 tasks)	82.2 _{±0.1}	69.7 _{±0.5}
Sentiment (5 tasks)	81.1 _{±0.2}	68.6 _{±0.1}
MRQA (6 tasks)	81.8 _{±0.2}	68.4 _{±0.2}
RAINBOW (6 tasks)	80.3 _{±0.6}	64.0 _{±0.4}
Translation (3 tasks)	82.4 _{±0.2}	65.3 _{±0.1}
Summarization (9 tasks)	80.9 _{±0.3}	67.1 _{±1.0}
GEM (8 tasks)	81.9 _{±0.2}	70.5 _{±0.5}
All (C4 + 55 supervised tasks)	81.8 _{±0.2}	67.9 _{±0.9}

为了探究源任务到目标任务的prompt可迁移性受什么因素影响，作者用16个源任务和10个目标任务构造了160个任务对。对于每一对任务，用源任务训练得到的prompt去初始化目标任务prompt，他们发现任务之间的相似度是一个重要的影响因素。作者假设不同任务prompts之间的相似度可以反映任务之间的相似度。这样的话，给定一个目标任务就可以去检索相似的源任务，用这些源任务来帮助模型执行目标任务。有如下几种方式

- Best of top-k: 依次使用top-k源任务的prompts去初始化目标任务prompt，然后选择表现最好的那个
- Top-k weighted average: 用top-k源任务的prompts的线性加权去初始化目标任务prompt，权重就是各个任务与目标任务的prompts的相似度
- Top-k multi-task mixture: 混合top-k源任务的数据进行训练得到一个prompt去初始化目标任务prompt

几种不同策略的实验结果见下图

Method	Change		Avg. score
	Abs.	Rel.	
<hr/>			
BASELINE	-	-	74.7 _{0.7}
<hr/>			
BRUTE-FORCE SEARCH ($k = 48$)			
<hr/>			
ORACLE	6.0 _{0.5}	26.5 _{1.1}	80.7 _{0.0}
<hr/>			
COSINE SIMILARITY OF AVERAGE TOKENS			
<hr/>			
BEST OF TOP- k			
$k = 1$	1.5 _{0.5}	11.7 _{1.1}	76.2 _{0.1}
$k = 3$	2.7 _{0.6}	16.6 _{1.1}	77.4 _{0.3}
$k = 6$	3.8 _{0.1}	20.0 _{1.1}	78.5 _{0.5}
$k = 9$	4.5 _{0.4}	22.2 _{1.1}	79.2 _{0.1}
$k = 12$	5.0 _{0.9}	23.6 _{2.2}	79.7 _{0.4}
$k = 15$	5.4 _{0.8}	24.9 _{1.8}	80.1 _{0.3}
<hr/>			
PER-TOKEN AVERAGE COSINE SIMILARITY			
<hr/>			
BEST OF TOP- k			
$k = 1$	2.0 _{0.4}	12.1 _{1.1}	76.7 _{0.7}
$k = 3$	2.9 _{0.6}	17.0 _{0.6}	77.5 _{0.4}
$k = 6$	4.5 _{0.5}	22.1 _{1.2}	79.2 _{0.1}
$k = 9$	4.6 _{0.5}	22.6 _{0.9}	79.5 _{0.2}
$k = 12$	5.0 _{0.6}	23.5 _{1.4}	79.6 _{0.1}
$k = 15$	5.3 _{0.9}	24.5 _{2.2}	80.0 _{0.4}
<hr/>			
TOP- k WEIGHTED AVERAGE			
best $k = 3$	1.9 _{0.5}	11.5 _{2.7}	76.6 _{0.1}
<hr/>			
TOP- k MULTI-TASK MIXTURE			
best $k = 12$	3.1 _{0.5}	15.3 _{2.8}	77.8 _{0.1}
<hr/>			

可以看到使用不同的策略相比于baseline即随机初始化都能够得到显著的提升。

On Transferability of Prompt Tuning for Natural Language Understanding

相比于之前的两篇工作，这个工作定义了更多的衡量prompts之间的相似度的指标。作者还提

可以看到使用不同的策略相比于baseline即随机初始化都能够得到明显的提升。

On Transferability of Prompt Tuning for Natural Language Understanding

相比于之前的两篇工作，这个工作定义了更多的衡量prompts之间的相似度的指标。作者还探究了prompts的zero-shot transfer的能力，即在源任务上训练的prompts直接用到目标任务上而不用目标任务的数据finetune，选择合适的源任务可以使用随机prompts得到明显更好的表现。见下图

Source Task	Target Task															
	IMDB	SST-2	laptop	restaurant	Movie	Tweet	MNLI	QNLI	SNLI	deontology	justice	QQP	MRPC	random	prompt	
IMDB	90	86	45	60	82	32	37	56	34	50	50	40	68			
SST-2	74	94	66	74	71	46	36	52	34	50	50	57	53			
laptop	80	87	77	79	69	53	34	50	34	50	50	54	44			
restaurant	68	86	72	81	67	57	33	52	34	50	50	45	65			
Movie	86	74	32	29	79	33	33	51	33	51	51	37	68			
Tweet	78	86	73	77	73	74	39	53	36	50	50	49	52			
MNLI	57	58	53	66	52	40	81	76	77	50	50	50	66			
QNLI	50	51	51	61	50	41	48	90	56	50	50	37	68			
SNLI	58	55	54	66	55	41	74	70	88	50	50	60	60			
deontology	51	51	3	1	52	41	33	50	33	73	59	37	64			
justice	63	52	2	1	57	42	33	51	33	62	70	39	57			
QQP	58	51	26	37	52	46	34	44	30	50	50	87	70			
MRPC	50	50	2	1	50	41	33	50	33	50	50	67	84			
random	51	51	2	1	50	40	33	50	33	50	50	37	68			

此外，使用合适的源任务，目标任务Prompt-tuning的收敛速度明显加快，见下图

Task Type	SA					NLI					EJ				
Task	IMDB	SST-2	laptop	restaurant	Movie	Tweet	MNLI	QNLI	SNLI	deontology	justice	QQP	MRPC	PT	PI
Labels	2	2	4	4	2	3	3	2	3	2	2	2	2	2	2
Accuracy (PT) (%)	89.9	93.8	77.3	80.7	79.2	74.5	80.6	90.5	88.5	72.9	70.0	86.9	83.9		
Accuracy (TPTTask) (%)	90.0	93.9	76.6	83.5	80.2	74.2	83.3	90.6	88.1	76.6	70.1	87.5	82.6		
Convergence Time (%)	90.6	65.3	77.3	28.9	41.7	52.3	46.5	94.2	94.1	75.0	34.1	133.0	57.7		
Comparable result Time (%)	53.1	54.5	-	3.3	1.5	-	2.3	94.2		12.4	2.2	107.0			

总结与展望

基于之前三篇研究工作，我们可以得到如下结论：

- 好的初始化对Prompt-tuning十分重要，可以解决Prompt-tuning的三大问题。Prompt transfer是初始化目标任务prompts的一种非常有效的方式。
- 源任务与目标任务之间的相似性是迁移有效性的一个重要因素。可以利用任务的prompts之间的相似度来估计任务的相似度。值得注意的是，这些实验结果都只是显示了一种趋势，更高的相似度并不能保证更好的迁移。
- 通过为多个源任务训练prompts可以得到一个prompts库，对于某个目标任务就可以检索一个或者是多个源任务的prompts去初始化目标任务prompt。目标任务prompt可以通过训练一定步数得到（不一定要等到收敛），检索可以通过prompts的相似度去进行排序，也可以通过源任务的prompts在目标任务上zero-shot transfer的表现去排序。如果使用多个源任务的prompts，有多种方式，例如逐个使用选最好或是加权平均。更复杂的方式可能获得更好的性能。
- PPT在目标任务训练数据较多时提升有限，而SPoT仍能有明显的提升。这可能是因为预训练任务与下游任务之间仍旧有显著的差距，更合适的源任务可能是更好的选择。

基于此，我们有如下展望：

- 如果可以将一个任务分解成多个子任务，那么是否可以训练多个子任务的prompts然后用这些子任务的prompts去得到原任务的prompt？比如任务型对话系统有多个模块，并不是每个数据都有所有模块的标注，通过这种方式就可以利用几乎所有现存的数据集，每个prompt只需要编码对应子任务的信息，然后将多个子任务的prompt组合起来就可以为原任务提供很好的初始化。
- 更进一步，是否可以设计任务相关的prompts，综合考虑原任务与子任务的特点与联系，使得子任务上训练的prompt直接放在原任务上复用？
- 是否可以结合Prompt-tuning和Model-tuning在目标任务上得到更好的表现？结合第二点展望，能否将此推广到多任务学习，多个任务之间设计相关的prompts然后一起训练？

笔者认为Prompt-tuning是一个十分有前景的技术，如何使其更加有效、适用更多场景还需要进一步探索。这就仰仗各位研究者（包括正在读文章的你）的努力啦~~