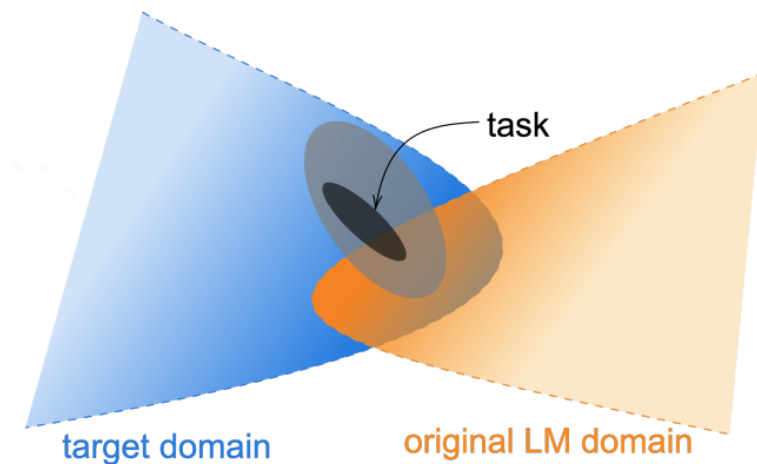


# Further Pretrain 简介

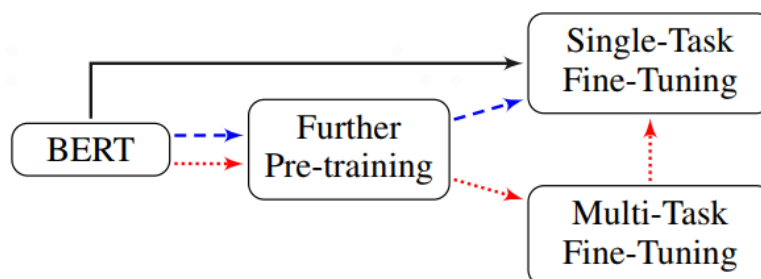
作者：纽约的自行车

日期：2022 年 10 月 10 号

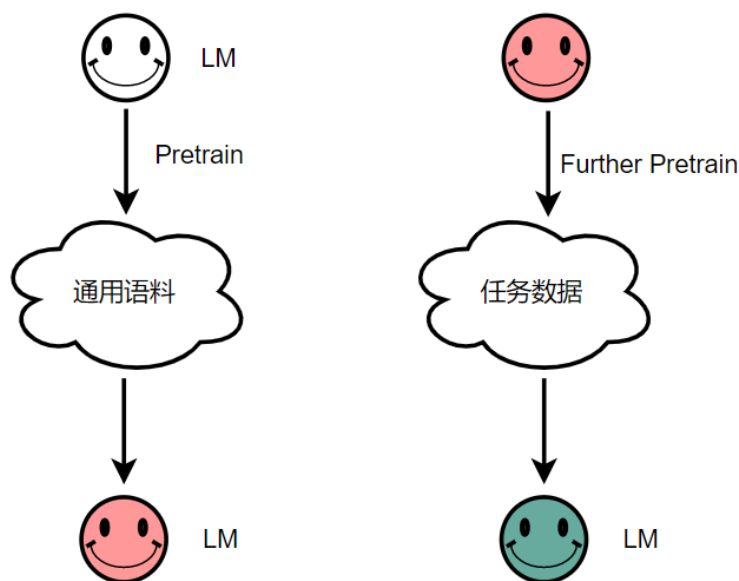
训练语言模型使用的通用语料通常与下游任务语料在领域上不同，如下图所示，黄色部分表示通用语料覆盖的领域，黑色部分表示下游任务数据覆盖的领域，灰色部分表示任务数据扩充所覆盖的领域，蓝色部分表示任务所在领域。下游任务数据和通用语料之间存在差异，那么是否有必要将语言模型迁移到下游任务领域上？



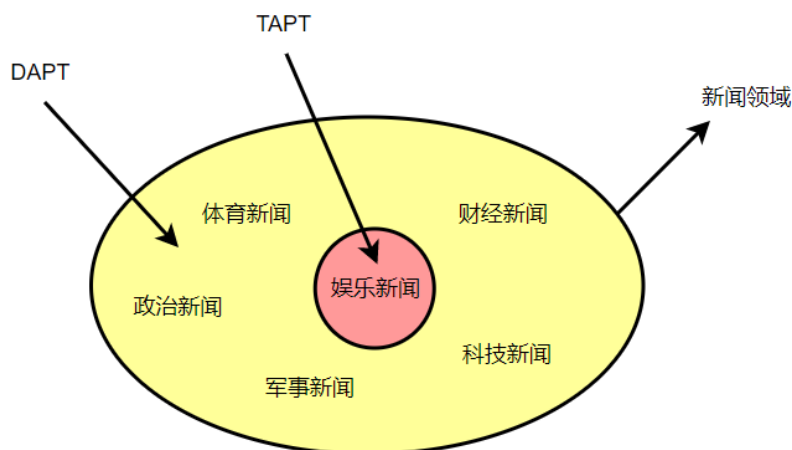
2019 年复旦大学邱锡鹏教授团队针对这个问题，尝试了 further pretrain [1]，并在文本分类任务上进行验证。相比用 Bert 直接在任务数据上 Finetune（下图黑色路径），作者增加了一个 Further Pretrain 步骤（下图蓝色路径），即在任务数据上继续做预训练，然后再 Finetune，使得流程由原来的两阶段变成三阶段。



什么是 further pretrain? further pretrain 是 pretrain 的延续，further pretrain 与 pretrain 的区别是换了个数据集，并且 further pretrain 使用的模型是在大规模语料上训练后的预训练模型，further pretrain 的训练任务可以和 pretrain 的训练任务相同，也可以不同，最常用的 further pretrain 训练任务是 mlm。



2020 年 Suchin Gururanga [2] 在 ACL 会议上发表了一篇主题相同的论文，论文中作者将 Further Pretrain 分为 Domain-Adaptive Pretraining (DAPT) 和 Task-Adaptive Pretraining (TAPT)，并做了更详细的实验论证。DAPT (Domain-adaptive pretraining) 是在任务数据所属领域的语料上进行 further pretrain，该方法适用于任务数据量非常少的场景。TAPT (Task-adaptive pretraining) 是直接在任务数据上进行 further pretrain，适用于任务数据量较大的场景。下图展示了 DAPT 和 TAPT 的区别。假设下游任务：娱乐新闻情感分类。DAPT 是指在新闻领域语料上做 further pretrain，但是需要排除娱乐新闻(下游任务数据)。TAPT 是指在下游任务数据上做 further pretrain。



作者选取低资源数据集和高资源数据集，下表中第二列左上角有十字标记的表示高资源。八个数据集分成四个领域，每个领域两个数据集，全部都是文本分类任务。

下表是实验结果，DATP+TAPT 是三阶段方式，使用 roberta 模型，现在相关语料上做 DAPT，然后在任务数据上做 TAPT，最后 finetune。虽然最终效果相比 DAPT 和 TAPT 确有提升，但是训练耗时长。

实验结论：

1. 在大数据集上 TAPT 平均提升 1.25%，在小数据集上 TAPT 平均提升 2.73%；

Domain	Task	Label Type	Train (Lab.)
BIO MED	CHEMPROT	relation classification	4169
	<sup>†</sup> RCT	abstract sent. roles	18040
CS	ACL-ARC	citation intent	1688
	SciERC	relation classification	3219
NEWS	HYPERPARTISAN	partisanship	515
	<sup>†</sup> AGNEWS	topic	115000
REVIEWS	<sup>†</sup> HELPFULNESS	review helpfulness	115251
	<sup>†</sup> IMDB	review sentiment	20000

Domain	Task	ROBERTA	Additional Pretraining Phases		
			DAPT	TAPT	DAPT + TAPT
BIO MED	CHEMPROT	81.9 <sub>1.0</sub>	84.2 <sub>0.2</sub>	82.6 <sub>0.4</sub>	<b>84.4</b> <sub>0.4</sub>
	<sup>†</sup> RCT	87.2 <sub>0.1</sub>	87.6 <sub>0.1</sub>	87.7 <sub>0.1</sub>	<b>87.8</b> <sub>0.1</sub>
CS	ACL-ARC	63.0 <sub>5.8</sub>	75.4 <sub>2.5</sub>	67.4 <sub>1.8</sub>	<b>75.6</b> <sub>3.8</sub>
	SciERC	77.3 <sub>1.9</sub>	80.8 <sub>1.5</sub>	79.3 <sub>1.5</sub>	<b>81.3</b> <sub>1.8</sub>
NEWS	HYPERPARTISAN	86.6 <sub>0.9</sub>	88.2 <sub>5.9</sub>	<b>90.4</b> <sub>5.2</sub>	90.0 <sub>6.6</sub>
	<sup>†</sup> AGNEWS	93.9 <sub>0.2</sub>	93.9 <sub>0.2</sub>	94.5 <sub>0.1</sub>	<b>94.6</b> <sub>0.1</sub>
REVIEWS	<sup>†</sup> HELPFULNESS	65.1 <sub>3.4</sub>	66.5 <sub>1.4</sub>	68.5 <sub>1.9</sub>	<b>68.7</b> <sub>1.8</sub>
	<sup>†</sup> IMDB	95.0 <sub>0.2</sub>	95.4 <sub>0.1</sub>	95.5 <sub>0.1</sub>	<b>95.6</b> <sub>0.1</sub>

2. 在大数据集上 DAPT 平均提升 0.55%，在小数据集上 DAPT 平均提升 4.95%；
3. 在大数据集上 TAPT 优于 DAPT，在小数据集上 DAPT 优于 TAPT；
4. 在大数据集上 DAPT+TAPT 相比 TAPT 有略微提升；
5. 在 roberta 的训练语料领域 news 上，TAPT 也有明显提升；

## 参考文献

- [1] How to Fine-Tune BERT for Text Classification? [J]. 2019.
- [2] Don't Stop Pretraining: Adapt Language Models to Domains and Tasks[J]. 2020.