



# How to Better Transfer the Pre-trained Language Model?

余阳 2022.03.23

## □ Effectiveness

- Finetuning
- Domain-adaptive finetuning

## □ Efficiency

- Prompt-based learning
- Parameter-efficient finetuning

## □ Conclusion

- Effectiveness
  - **Finetuning**
  - Domain-adaptive finetuning
- Efficiency
  - Prompt-based learning
  - Parameter-efficient finetuning
- Conclusion

- Pre-train -> Finetune (**NLP Paradigm #3**)
- Some tricks
  - "The lower layer of the PLM may contain more general information."
  - **Layer selection**: select the most effective layer(s) for the downstream task.
  - **Layer-wise decreasing learning rate**

$$\theta_t^l = \theta_{t-1}^l - \eta^l \cdot \nabla_{\theta^l} J(\theta)$$

$$\eta^{k-1} = \xi \cdot \eta^k$$

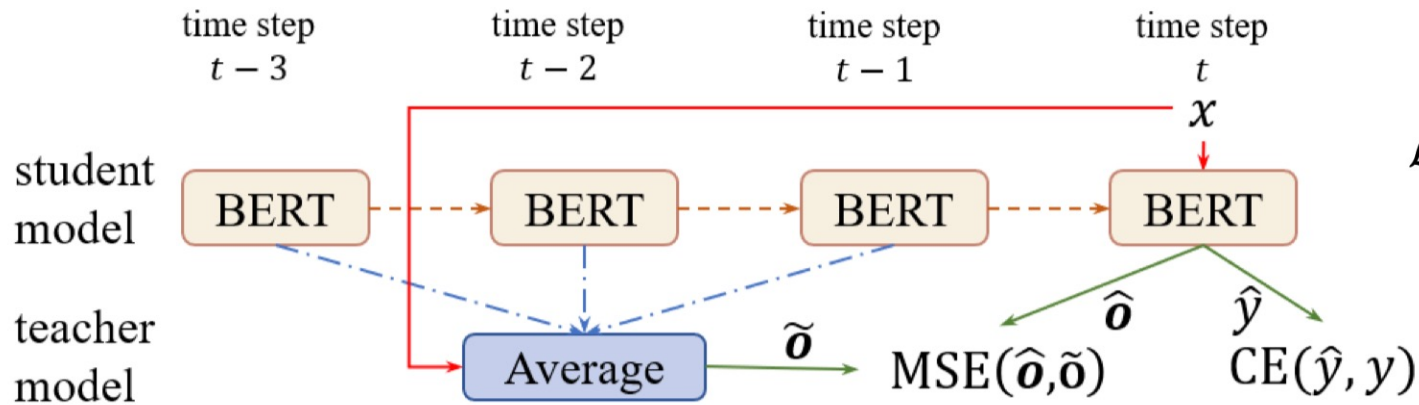
## □ Some tricks

□ **Self-ensemble & Self-distillation:** improve the stability of finetuning.

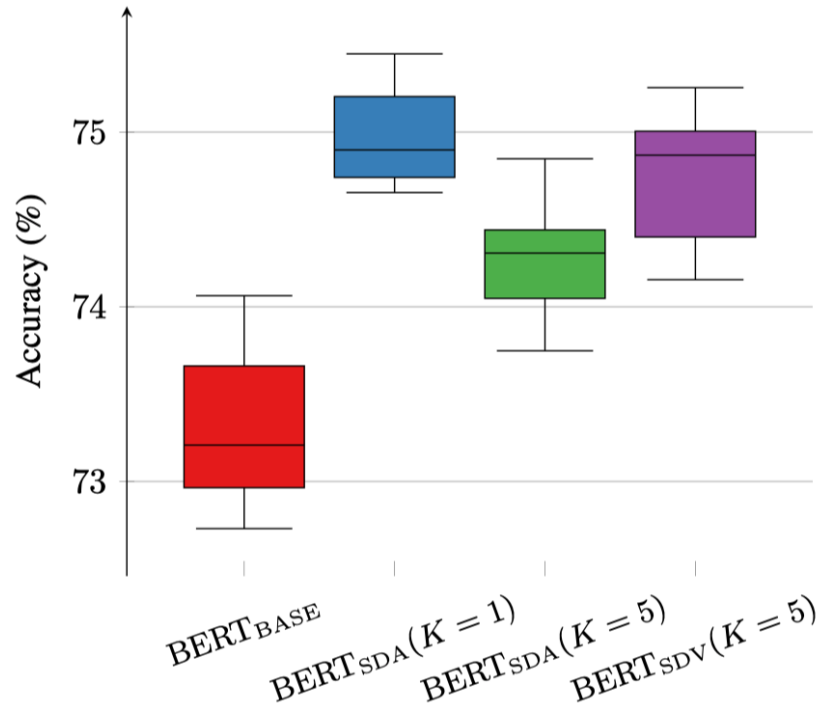
### □ Self-ensemble

$$\text{BERT}_{\text{SE}}(x; \bar{\theta}) = \text{BERT}\left(x; \frac{1}{T} \sum_{\tau=1}^T \theta_{\tau}\right)$$

### □ Self-distillation



- Some tricks
  - **Self-ensemble & Self-distillation:** improve the stability of finetuning.



- Effectiveness
  - Finetuning
  - **Domain-adaptive finetuning**
- Efficiency
  - Prompt-based learning
  - Parameter-efficient finetuning
- Conclusion

# Domain-adaptive Finetuning

- ❑ The PLM is usually pre-trained on general corpora.
  - ❑ BookCorpus
  - ❑ Wikipedia
  - ❑ CCNews, OpenWebText, CommonCrawl, etc.
- ❑ Narrow the **data distribution gap** between the pre-training data and the downstream task data.

PT	100.0	54.1	34.5	27.3	19.2
News	54.1	100.0	40.0	24.9	17.3
Reviews	34.5	40.0	100.0	18.3	12.7
BioMed	27.3	24.9	18.3	100.0	21.4
CS	19.2	17.3	12.7	21.4	100.0
	PT	News	Reviews	BioMed	CS

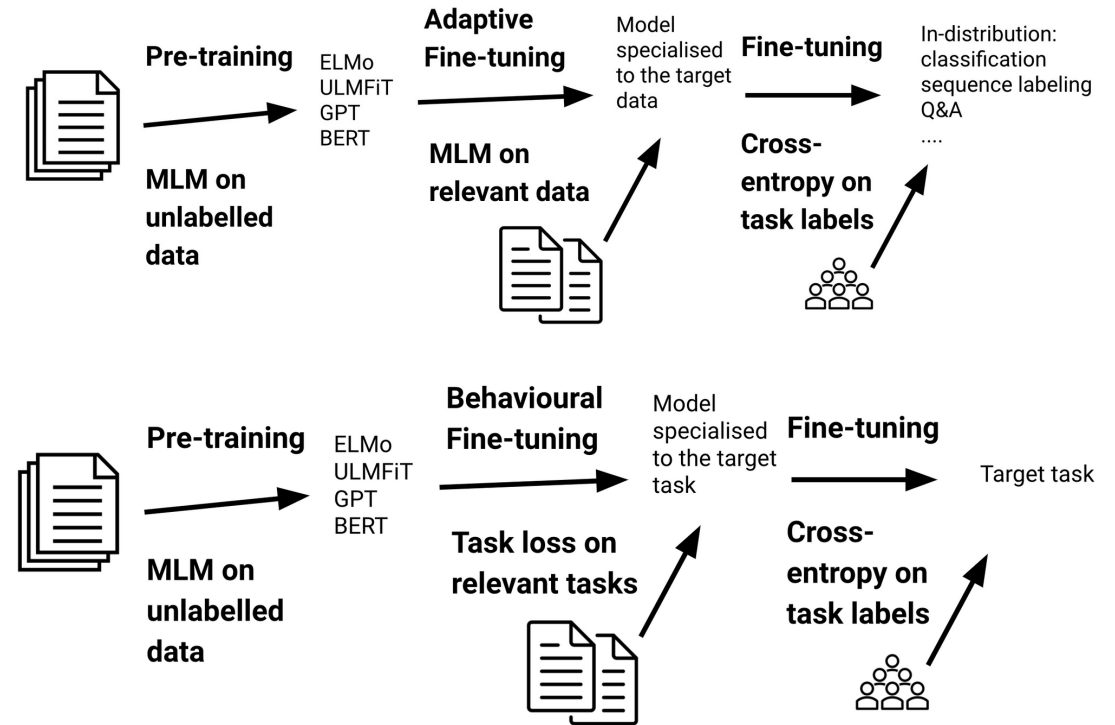
Figure 2: Vocabulary overlap (%) between domains. PT denotes a sample from sources similar to ROBERTA’s pretraining corpus. Vocabularies for each domain are created by considering the top 10K most frequent words (excluding stopwords) in documents sampled from each domain.

Domain	Pretraining Corpus	# Tokens	Size
BIOMED	2.68M full-text papers from S2ORC (Lo et al., 2020)	7.55B	47GB
CS	2.22M full-text papers from S2ORC (Lo et al., 2020)	8.10B	48GB
NEWS	11.90M articles from REALNEWS (Zellers et al., 2019)	6.66B	39GB
REVIEWS	24.75M AMAZON reviews (He and McAuley, 2016)	2.11B	11GB



# Domain-adaptive Finetuning

- ❑ Self-supervised task on downstream domain corpora
  - ❑ Domain-adaptive pre-training (DAPT) / Task-adaptive pre-training (TAPT)
  - ❑ Contrastive Learning, etc.
- ❑ Supervised relevant task
  - ❑ Sequential
  - ❑ Parallel (i.e., multi-task finetuning)



Domain	Task	ROBERTA	Additional Pretraining Phases		
			DAPT	TAPT	DAPT + TAPT
BIOMED	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>
	†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>
	†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	†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>
	†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>

- Effectiveness
  - Finetuning
  - Domain-adaptive finetuning
- Efficiency
  - **Prompt-based learning**
  - Parameter-efficient finetuning
- Conclusion

## □ Motivation

- The PLM need to be finetuned for every new downstream task.
- Finetuning is costly for extremely large PLMs.
  - E.g., T5 (11B), GPT-3 (175B)
- Finetuning requires thousands to hundred of thousands task-specific examples.
  - Human can perform a new language task with a few examples or simple instructions.

## □ Solution

- In-context learning (prompt + demonstration) -> prompt-based learning

# Prompt-based Learning



## The three settings we explore for in-context learning

### Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

### One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← example
3 cheese => ..... ← prompt
```

### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ← demonstration
4 plush girafe => girafe peluche ← examples
5 cheese => ..... ← prompt
```

## Traditional fine-tuning (not used for GPT-3)

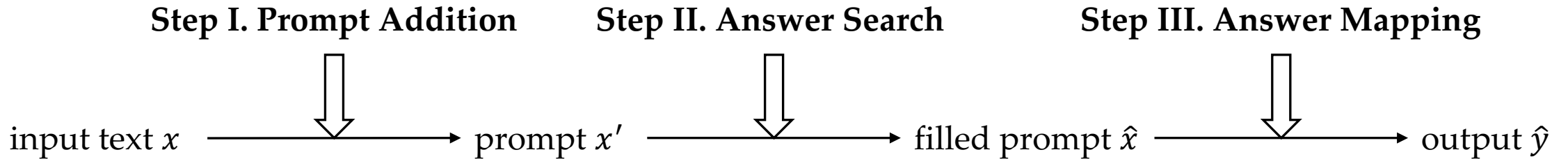
### Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



# Prompt-based Learning

- **Definition:** Prompt is the technique of making better use of the knowledge from the PLM by adding additional texts to the input.
- **Basic framework**



## □ Prompt addition

- Design a **template** with two slots: input slot [X] and answer slot [Z].
- Fill slot [X] with the input text  $x$ .

### □ E.g., sentiment analysis (**prefix prompt**)

template = "[X] The movie is [Z]."

input text  $x$  = "I love this movie."

prompt  $x'$  = "I love this movie. The movie is [Z]."

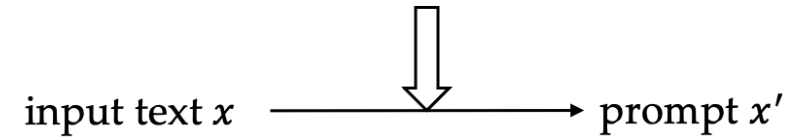
### □ E.g., named entity recognition (**cloze prompt**)

template = "[X1] [X2] is a [Z] entity."

input text [X1] = "Mike went to Paris." [X2] = "Paris"

prompt  $x'$  = "Mike went to Paris. Paris is a [Z] entity."

### Step I. Prompt Addition



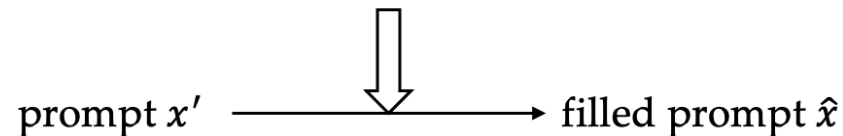
## □ Answer search

- Search for text  $\hat{z} \in \mathcal{Z}$  that maximize the score of a PLM  $P(\cdot; \theta)$ .

$$\hat{z} = \underset{z \in \mathcal{Z}}{\text{search}} P(f_{\text{fill}}(\mathbf{x}', z); \theta)$$

- $\mathcal{Z}$  can be the entire vocabulary set, or a small subset specific to the target task.
- The search function can be implemented as *argmax* or *sampling*.
- E.g.,  $\mathcal{Z} = \{\text{"excellent"}, \text{"good"}, \text{"OK"}, \text{"bad"}, \text{"horrible"}\}$  for sentiment analysis.

### Step II. Answer Search



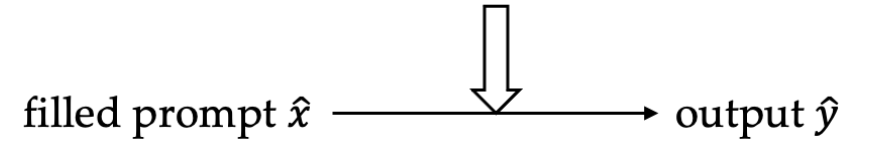
# Prompt-based Learning



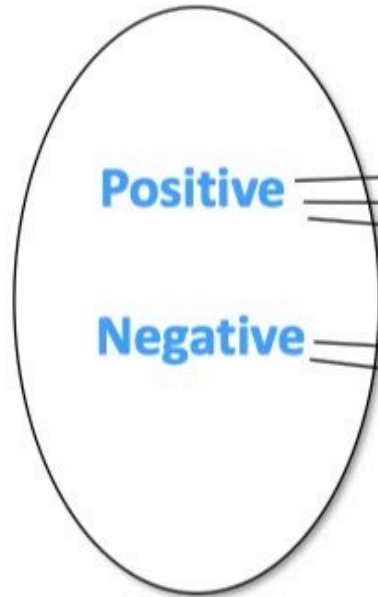
## Answer mapping

- Map the answer  $\hat{z}$  to the output  $\hat{y}$ .
- Multiple answers can result in the same output.

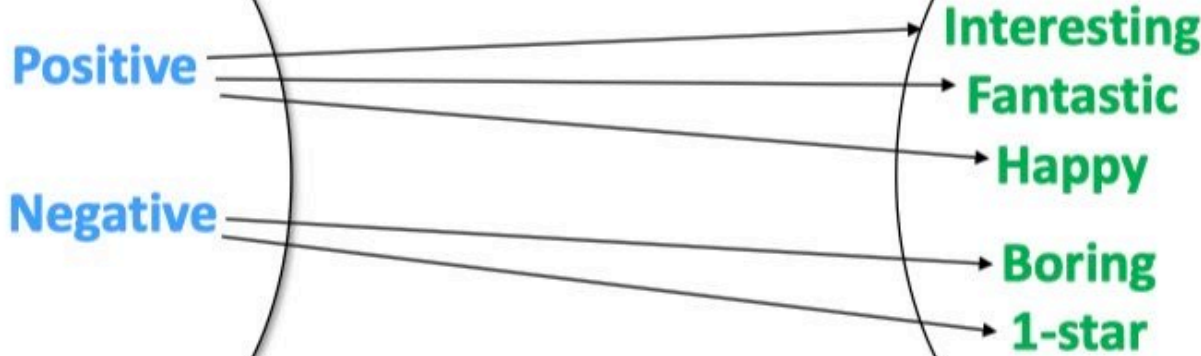
### Step III. Answer Mapping



### Label Space (Y)



### Answer Space (Z)

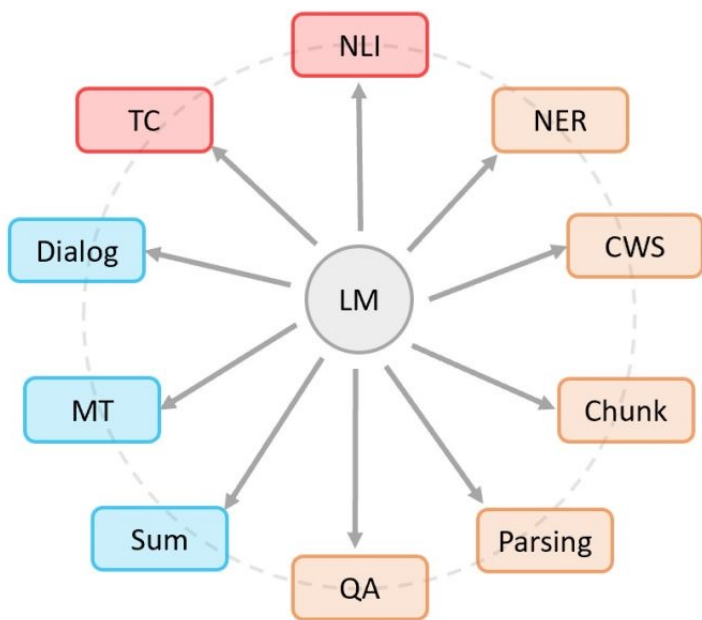




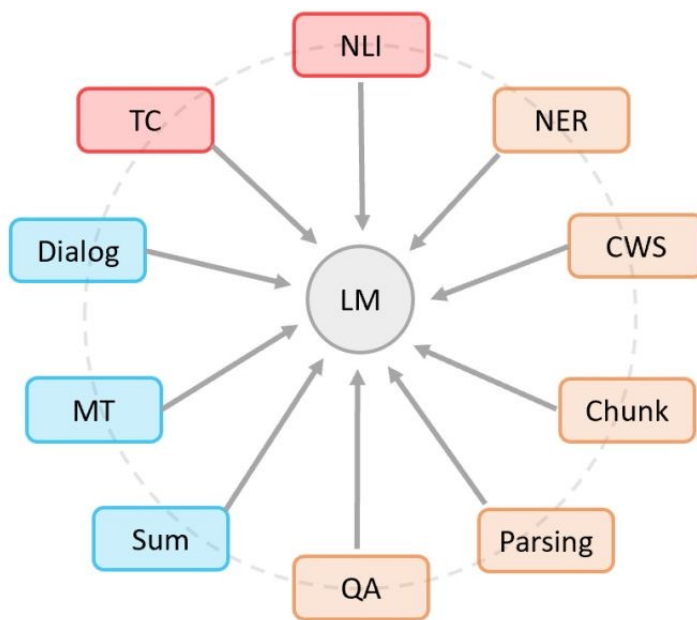
# Prompt-based Learning

16

- Pre-train  $\rightarrow$  Prompt  $\rightarrow$  Predict (**NLP Paradigm #4**)
- Narrow the gap between the pre-training task and the downstream task.



Finetuning



Prompting

**Input:**  $x =$  I love this movie.

**Template:** [x]  
Overall, it was a [z] movie.

**Answer:**  
{fantastic:😊,  
boring:☹️}

**Prompting:**  $x' =$  I love this movie.  
Overall, it was a [z] movie.

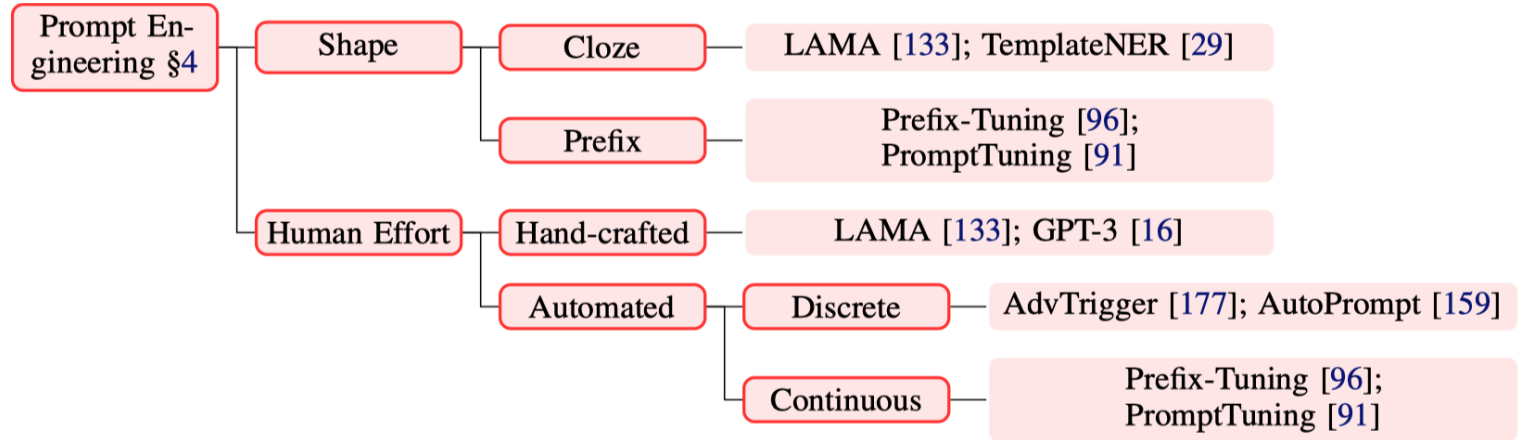
**Predicting:**  $x' =$  I love this movie.  
Overall, it was a **fantastic** movie.

**Mapping:** fantastic  $\Rightarrow$  😊

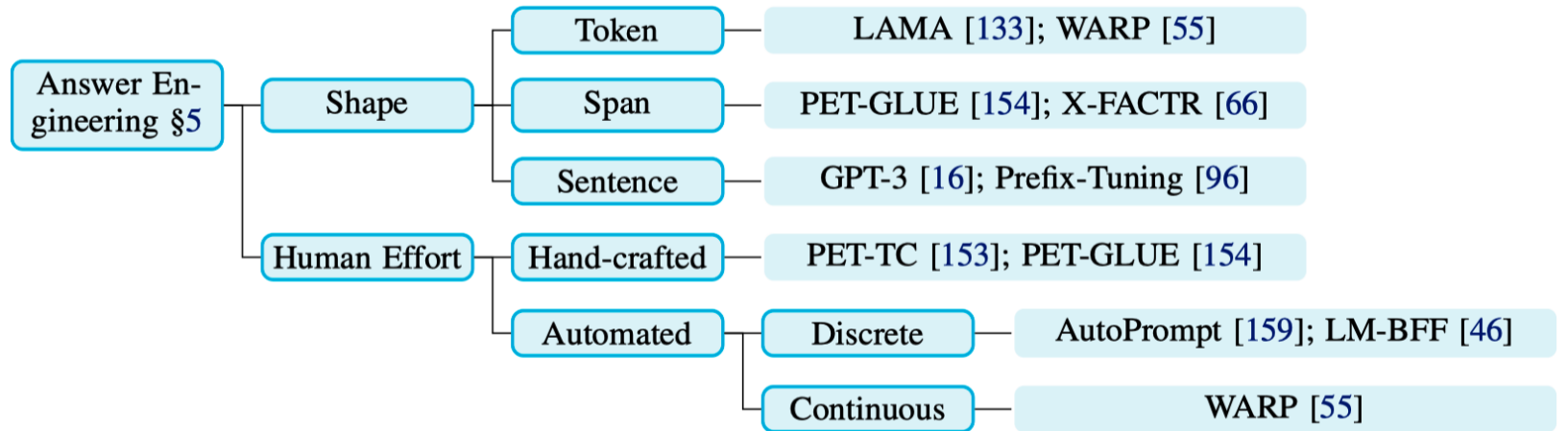
# Prompt-based Learning



- Human efforts
  - Prompt engineering



- Answer engineering



# Prompt-based Learning

- Human efforts
  - Prompt-based training strategies

zero-shot

Strategy	LM Params	Prompt Params		Example
		Additional	Tuned	
Promptless Fine-tuning	Tuned	-		ELMo [130], BERT [32], BART [94]
Tuning-free Prompting	Frozen	✗	✗	GPT-3 [16], AutoPrompt [159], LAMA [133]
Fixed-LM Prompt Tuning	Frozen	✓	Tuned	Prefix-Tuning [96], Prompt-Tuning [91]
Fixed-prompt LM Tuning	Tuned	✗	✗	PET-TC [153], PET-Gen [152], LM-BFF [46]
Prompt+LM Fine-tuning	Tuned	✓	Tuned	PADA [8], P-Tuning [103], PTR [56]

few-shot

# Prompt-based Learning

- **Making Pre-trained Language Models Better Few-shot Learners. (ACL 2021)**
  - "Finding the right prompts, however, is an art – requiring both domain expertise and an understanding of the language model’s inner workings."

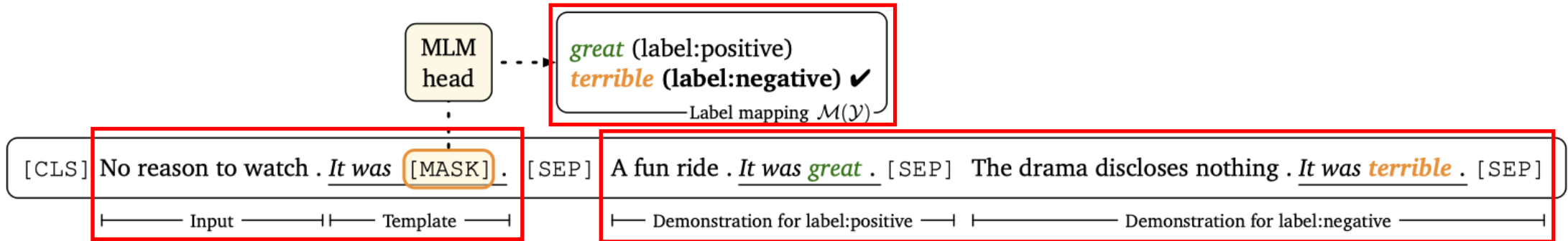
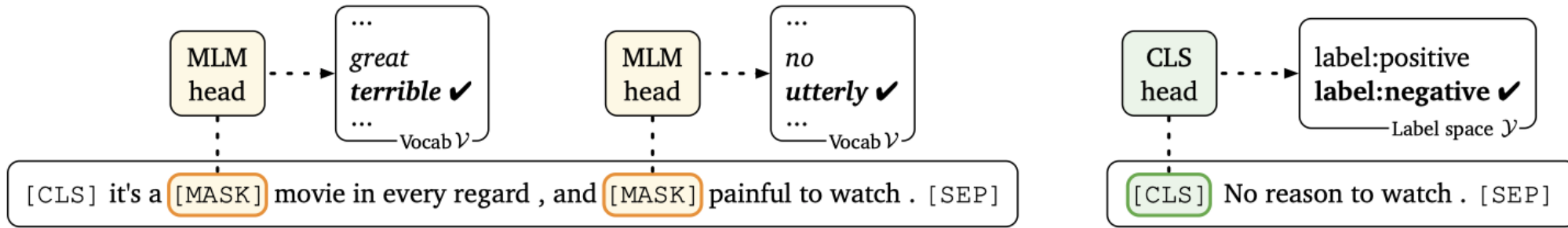
Template	Label words	Accuracy
SST-2 (positive/negative)		mean (std)
<S <sub>1</sub> > It was [MASK] .	great/terrible	<b>92.7 (0.9)</b>
<S <sub>1</sub> > It was [MASK] .	good/bad	92.5 (1.0)
<S <sub>1</sub> > It was [MASK] .	cat/dog	91.5 (1.4)
<S <sub>1</sub> > It was [MASK] .	dog/cat	86.2 (5.4)
<S <sub>1</sub> > It was [MASK] .	terrible/great	83.2 (6.9)
Fine-tuning	-	81.4 (3.8)
SNLI (entailment/neutral/contradiction)		mean (std)
<S <sub>1</sub> > ? [MASK] , <S <sub>2</sub> >	Yes/Maybe/No	<b>77.2 (3.7)</b>
<S <sub>1</sub> > . [MASK] , <S <sub>2</sub> >	Yes/Maybe/No	76.2 (3.3)
<S <sub>1</sub> > ? [MASK] <S <sub>2</sub> >	Yes/Maybe/No	74.9 (3.0)
<S <sub>1</sub> > <S <sub>2</sub> > [MASK]	Yes/Maybe/No	65.8 (2.4)
<S <sub>2</sub> > ? [MASK] , <S <sub>1</sub> >	Yes/Maybe/No	62.9 (4.1)
<S <sub>1</sub> > ? [MASK] , <S <sub>2</sub> >	Maybe/No/Yes	60.6 (4.8)
Fine-tuning	-	48.4 (4.8)



# Prompt-based Learning



## □ Making Pre-trained Language Models Better Few-shot Learners. (ACL 2021)





# Prompt-based Learning

21

- **Making Pre-trained Language Models Better Few-shot Learners. (ACL 2021)**
  - **Automatic selection of label words** (given a fixed template  $\mathcal{T}$ )
    - For each class  $c$ , select top  $k$  words that maximize the total probability of  $D_{train}^c$  using the initial PLM.

$$\text{Top-}k_{v \in \mathcal{V}} \left\{ \sum_{x_{in} \in \mathcal{D}_{train}^c} \log P_{\mathcal{L}}([\text{MASK}] = v \mid \mathcal{T}(x_{in})) \right\}$$

- Further find the top  $n$  words that maximize zero-shot accuracy on  $D_{train}$ .
- Finetune all top  $n$  assignments and select the best one on  $D_{dev}$ .



# Prompt-based Learning

## □ Making Pre-trained Language Models Better Few-shot Learners. (ACL 2021)

### □ Automatic generation of templates (given a fixed set of label words $\mathcal{M}(Y)$ )

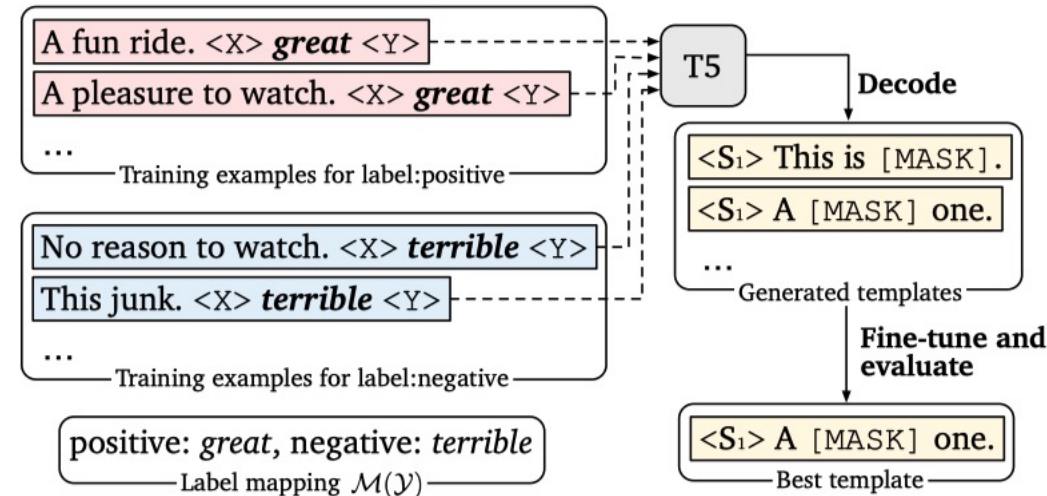
- Conduct simple conversions to each training sample  $(x_{in}, y) \in D_{train}$ .

$\langle S_1 \rangle \rightarrow \langle X \rangle \mathcal{M}(y) \langle Y \rangle \langle S_1 \rangle,$

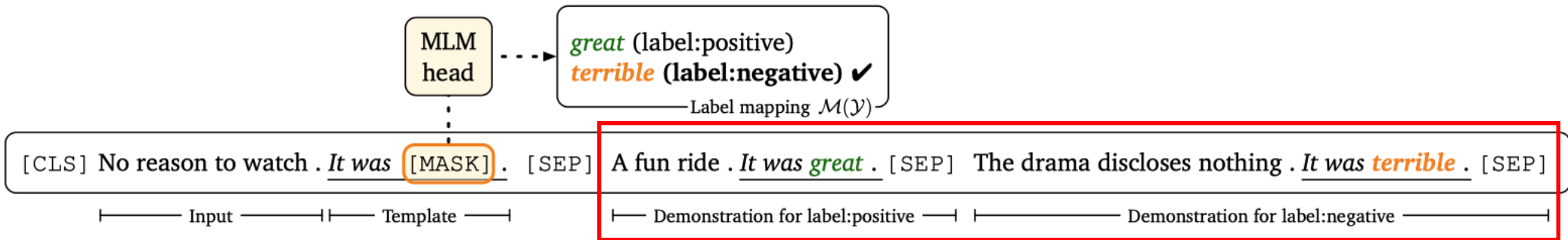
$\langle S_1 \rangle \rightarrow \langle S_1 \rangle \langle X \rangle \mathcal{M}(y) \langle Y \rangle,$

$\langle S_1 \rangle, \langle S_2 \rangle \rightarrow \langle S_1 \rangle \langle X \rangle \mathcal{M}(y) \langle Y \rangle \langle S_2 \rangle.$

- Use T5 to fill in missing spans.
- Use beam search to decode multiple templates.
- Finetune each generated templated on  $D_{train}$  and select the best one on  $D_{dev}$ .



- Making Pre-trained Language Models Better Few-shot Learners. (ACL 2021)
  - Finetune with demonstrations
    - Use the pre-trained SBERT to generate the embeddings of each training sample  $x_{in}^{(c)}$ .
    - Randomly sample one example  $(x_{in}^{(c)}, y^{(c)})$  from the top 50% samples that most similar to  $x_{in}$ .
    - Convert to filled prompt  $\tilde{\mathcal{J}}(x_{in}^{(c)}, y^{(c)})$  and concatenate with  $x_{in}$ .





## □ Making Pre-trained Language Models Better Few-shot Learners. (ACL 2021)

Task	Auto template	Auto label words
SST-2	(positive/negative)	
	$\langle S_1 \rangle$ A [MASK] one .	irresistible/pathetic
	$\langle S_1 \rangle$ A [MASK] piece .	wonderful/bad
	$\langle S_1 \rangle$ All in all [MASK] .	delicious/bad
SST-5	(very positive/positive/neutral/negative/very negative)	
	$\langle S_1 \rangle$ The movie is [MASK] .	wonderful/remarkable/hilarious/better/awful
	$\langle S_1 \rangle$ The music is [MASK] .	wonderful/perfect/hilarious/better/awful
	$\langle S_1 \rangle$ But it is [MASK] .	unforgettable/extraordinary/good/better/terrible
MR	(positive/negative)	
	It was [MASK] ! $\langle S_1 \rangle$	epic/terrible
	$\langle S_1 \rangle$ It's [MASK] .	epic/awful
	$\langle S_1 \rangle$ A [MASK] piece of work .	exquisite/horrible
CR	(positive/negative)	
	$\langle S_1 \rangle$ It's [MASK] !	fantastic/horrible
	$\langle S_1 \rangle$ The quality is [MASK] .	neat/pointless
	$\langle S_1 \rangle$ That is [MASK] .	magnificent/unacceptable

## □ Making Pre-trained Language Models Better Few-shot Learners. (ACL 2021)

	SST-2 (acc)	SST-5 (acc)	MR (acc)	CR (acc)	MPQA (acc)	Subj (acc)	TREC (acc)	CoLA (Matt.)
Majority <sup>†</sup>	50.9	23.1	50.0	50.0	50.0	50.0	18.8	0.0
Prompt-based zero-shot <sup>‡</sup>	83.6	35.0	80.8	79.5	67.6	51.4	32.0	2.0
“GPT-3” in-context learning	84.8 (1.3)	30.6 (0.9)	80.5 (1.7)	87.4 (0.8)	63.8 (2.1)	53.6 (1.0)	26.2 (2.4)	-1.5 (2.4)
Fine-tuning	81.4 (3.8)	43.9 (2.0)	76.9 (5.9)	75.8 (3.2)	72.0 (3.8)	90.8 (1.8)	88.8 (2.1)	<b>33.9</b> (14.3)
Prompt-based FT (man)	92.7 (0.9)	47.4 (2.5)	87.0 (1.2)	90.3 (1.0)	84.7 (2.2)	91.2 (1.1)	84.8 (5.1)	9.3 (7.3)
+ demonstrations	92.6 (0.5)	<b>50.6</b> (1.4)	86.6 (2.2)	90.2 (1.2)	<b>87.0</b> (1.1)	<b>92.3</b> (0.8)	87.5 (3.2)	18.7 (8.8)
Prompt-based FT (auto)	92.3 (1.0)	49.2 (1.6)	85.5 (2.8)	89.0 (1.4)	85.8 (1.9)	91.2 (1.1)	88.2 (2.0)	14.0 (14.1)
+ demonstrations	<b>93.0</b> (0.6)	49.5 (1.7)	<b>87.7</b> (1.4)	<b>91.0</b> (0.9)	86.5 (2.6)	91.4 (1.8)	<b>89.4</b> (1.7)	21.8 (15.9)
Fine-tuning (full) <sup>†</sup>	95.0	58.7	90.8	89.4	87.8	97.0	97.4	62.6
	MNLI (acc)	MNLI-mm (acc)	SNLI (acc)	QNLI (acc)	RTE (acc)	MRPC (F1)	QQP (F1)	STS-B (Pear.)
Majority <sup>†</sup>	32.7	33.0	33.8	49.5	52.7	81.2	0.0	-
Prompt-based zero-shot <sup>‡</sup>	50.8	51.7	49.5	50.8	51.3	61.9	49.7	-3.2
“GPT-3” in-context learning	52.0 (0.7)	53.4 (0.6)	47.1 (0.6)	53.8 (0.4)	60.4 (1.4)	45.7 (6.0)	36.1 (5.2)	14.3 (2.8)
Fine-tuning	45.8 (6.4)	47.8 (6.8)	48.4 (4.8)	60.2 (6.5)	54.4 (3.9)	76.6 (2.5)	60.7 (4.3)	53.5 (8.5)
Prompt-based FT (man)	68.3 (2.3)	70.5 (1.9)	77.2 (3.7)	64.5 (4.2)	69.1 (3.6)	74.5 (5.3)	65.5 (5.3)	71.0 (7.0)
+ demonstrations	<b>70.7</b> (1.3)	<b>72.0</b> (1.2)	<b>79.7</b> (1.5)	<b>69.2</b> (1.9)	68.7 (2.3)	77.8 (2.0)	<b>69.8</b> (1.8)	73.5 (5.1)
Prompt-based FT (auto)	68.3 (2.5)	70.1 (2.6)	77.1 (2.1)	68.3 (7.4)	<b>73.9</b> (2.2)	76.2 (2.3)	67.0 (3.0)	75.0 (3.3)
+ demonstrations	70.0 (3.6)	<b>72.0</b> (3.1)	77.5 (3.5)	68.5 (5.4)	71.1 (5.3)	<b>78.1</b> (3.4)	67.7 (5.8)	<b>76.4</b> (6.2)
Fine-tuning (full) <sup>†</sup>	89.8	89.5	92.6	93.3	80.9	91.4	81.7	91.9

- Effectiveness
  - Finetuning
  - Domain-adaptive finetuning
- Efficiency
  - Prompt-based learning
  - **Parameter-efficient finetuning**
- Conclusion



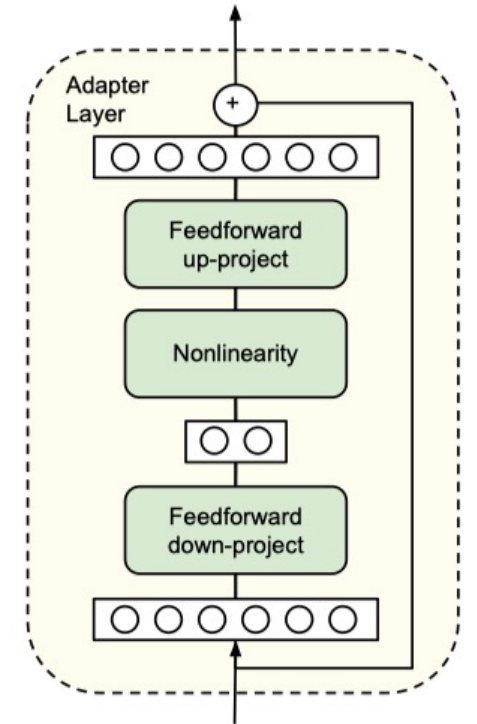
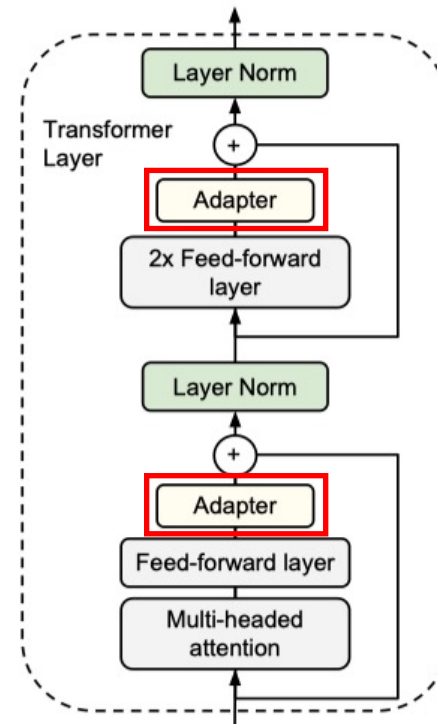
# Parameter-efficient Finetuning

27

- ❑ **Motivation:** Finetuning the entire PLM is parameter inefficient.
- ❑ **Solution:** Fix the PLM and only finetune a few additional parameters.
  - ❑ Adapter
  - ❑ Prefix-tuning & Prompt-tuning
  - ❑ Low-rank adaptation

# Parameter-efficient Finetuning

- ❑ **Parameter-Efficient Transfer Learning for NLP (ICML 2019)**
  - ❑ Add some small adapter modules between layers of the PLM.
  - ❑ A new set of adapters are added and finetuned for every new task.
  - ❑ Adapter modules
    - ❑ **small number of parameters**
      - ❑ origin dimension  $d$ , projection dimension  $m$
      - ❑ total number of parameters =  $2md + d + m$
    - ❑ **near-identity initialization**
      - ❑ skip-connection
      - ❑ near-zero initialization for projection layers

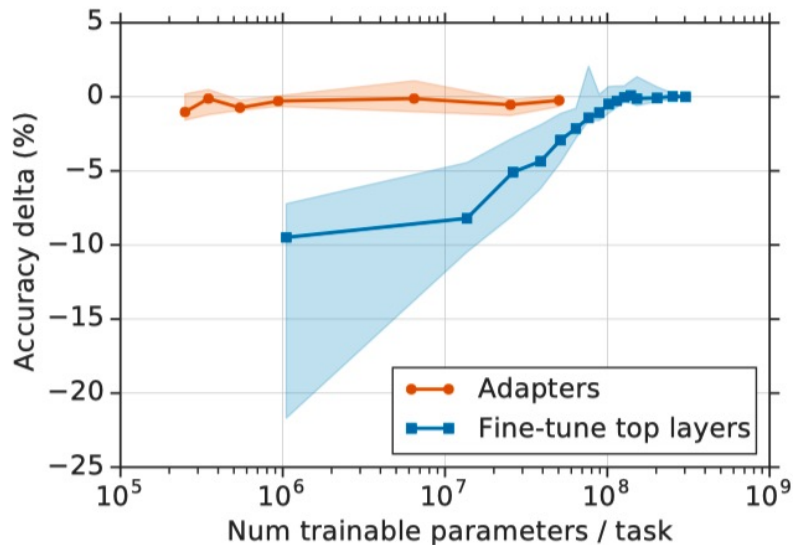


# Parameter-efficient Finetuning

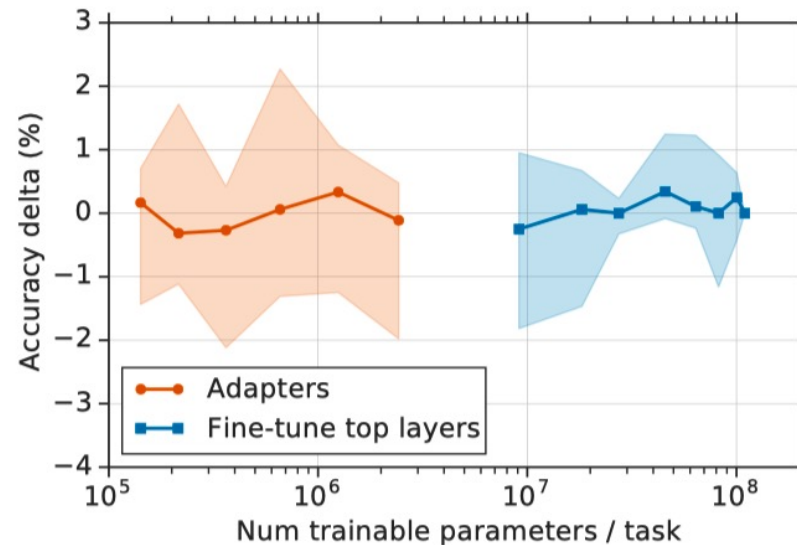
## Parameter-Efficient Transfer Learning for NLP (ICML 2019)

	Total num params	Trained params / task	CoLA	SST	MRPC	STS-B	QQP	MNLI <sub>m</sub>	MNLI <sub>mm</sub>	QNLI	RTE	Total
BERT <sub>LARGE</sub>	9.0×	100%	60.5	94.9	89.3	87.6	72.1	86.7	85.9	91.1	70.1	80.4
Adapters (8-256)	1.3×	3.6%	59.5	94.0	89.5	86.9	71.8	84.9	85.1	90.7	71.5	80.0
Adapters (64)	1.2×	2.1%	56.9	94.2	89.6	87.3	71.8	85.3	84.6	91.4	68.8	79.6

GLUE (BERT<sub>LARGE</sub>)



Additional Tasks (BERT<sub>BASE</sub>)



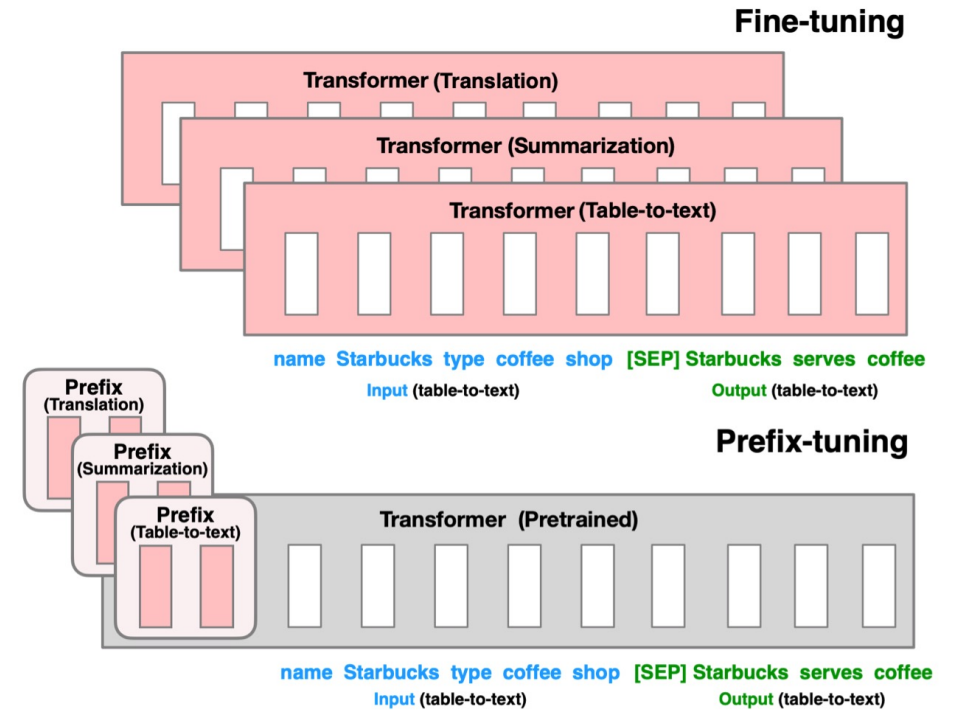
# Parameter-efficient Finetuning

## Prefix-Tuning

- Inspired by prompting: having a proper context can steer the LM without changing its parameters.
- Optimize a small **continuous** task-specific vector (called the **prefix**).

$$h_i = \begin{cases} P_\theta[i, :], & \text{if } i \in P_{\text{idx}}, \\ \text{LM}_\phi(z_i, h_{<i}), & \text{otherwise.} \end{cases}$$

$$P_\theta[i, :] = \text{MLP}_\theta(P'_\theta[i, :])$$

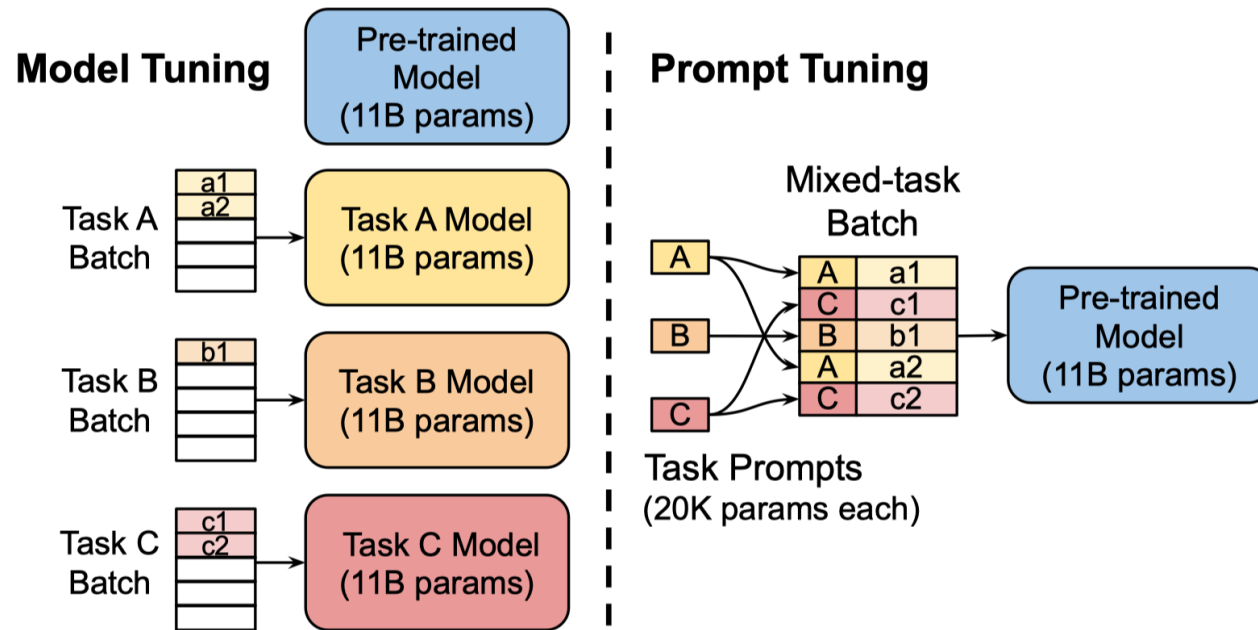




# Parameter-efficient Finetuning

## □ Prompt-Tuning

- A simplification of prefix-tuning.
- Only allow an additional  $k$  tunable tokens per downstream task to be prepended to the input text.





# Parameter-efficient Finetuning

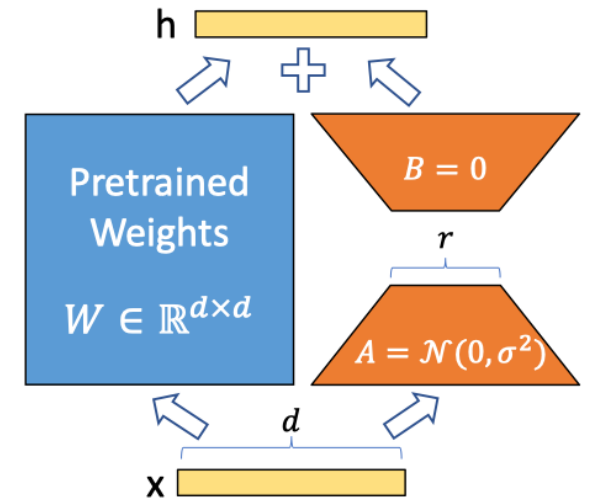
32

- LoRA: Low-Rank Adaptation of Language Models (2021)
  - Hypothesize: The change in weights during model adaptation has a low "intrinsic rank".
  - Inject trainable rank decomposition matrices into each layer of the PLM.

$$W_0 + \Delta W = W_0 + BA$$

$$W_0 \in \mathbb{R}^{d \times k} \quad A \in \mathbb{R}^{r \times k} \quad B \in \mathbb{R}^{d \times r}$$

- Only adapt the attention weights (i.e.,  $W_q, W_k, W_v, W_o$ ).
- No additional inference latency.
- Cannot put samples of different tasks into the same batch.



# Parameter-efficient Finetuning



33

## □ Low-Rank Adaptation of Language Models (2021)

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB <sub>base</sub> (FT)*	125.0M	<b>87.6</b>	94.8	90.2	<b>63.6</b>	92.8	<b>91.9</b>	78.7	91.2	86.4
RoB <sub>base</sub> (BitFit)*	0.1M	84.7	93.7	<b>92.7</b>	62.0	91.8	84.0	81.5	90.8	85.2
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.3M	87.1 $\pm$ 0	94.2 $\pm$ 1	88.5 $\pm$ 1.1	60.8 $\pm$ 4	93.1 $\pm$ 1	90.2 $\pm$ 0	71.5 $\pm$ 2.7	89.7 $\pm$ 3	84.4
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.9M	87.3 $\pm$ 1	94.7 $\pm$ 3	88.4 $\pm$ 1	62.6 $\pm$ 9	93.0 $\pm$ 2	90.6 $\pm$ 0	75.9 $\pm$ 2.2	90.3 $\pm$ 1	85.4
RoB <sub>base</sub> (LoRA)	0.3M	87.5 $\pm$ 3	<b>95.1<math>\pm</math>2</b>	89.7 $\pm$ 7	63.4 $\pm$ 1.2	<b>93.3<math>\pm</math>3</b>	90.8 $\pm$ 1	<b>86.6<math>\pm</math>7</b>	<b>91.5<math>\pm</math>2</b>	<b>87.2</b>
RoB <sub>large</sub> (FT)*	355.0M	90.2	<b>96.4</b>	<b>90.9</b>	68.0	94.7	<b>92.2</b>	86.6	92.4	88.9
RoB <sub>large</sub> (LoRA)	0.8M	<b>90.6<math>\pm</math>2</b>	96.2 $\pm$ 5	<b>90.9<math>\pm</math>1.2</b>	<b>68.2<math>\pm</math>1.9</b>	<b>94.9<math>\pm</math>3</b>	91.6 $\pm$ 1	<b>87.4<math>\pm</math>2.5</b>	<b>92.6<math>\pm</math>2</b>	<b>89.0</b>
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	3.0M	90.2 $\pm$ 3	96.1 $\pm$ 3	90.2 $\pm$ 7	<b>68.3<math>\pm</math>1.0</b>	<b>94.8<math>\pm</math>2</b>	<b>91.9<math>\pm</math>1</b>	83.8 $\pm$ 2.9	92.1 $\pm$ 7	88.4
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	0.8M	<b>90.5<math>\pm</math>3</b>	<b>96.6<math>\pm</math>2</b>	89.7 $\pm$ 1.2	67.8 $\pm$ 2.5	<b>94.8<math>\pm</math>3</b>	91.7 $\pm$ 2	80.1 $\pm$ 2.9	91.9 $\pm$ 4	87.9
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	6.0M	89.9 $\pm$ 5	96.2 $\pm$ 3	88.7 $\pm$ 2.9	66.5 $\pm$ 4.4	94.7 $\pm$ 2	92.1 $\pm$ 1	83.4 $\pm$ 1.1	91.0 $\pm$ 1.7	87.8
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	0.8M	90.3 $\pm$ 3	96.3 $\pm$ 5	87.7 $\pm$ 1.7	66.3 $\pm$ 2.0	94.7 $\pm$ 2	91.5 $\pm$ 1	72.9 $\pm$ 2.9	91.5 $\pm$ 5	86.4
RoB <sub>large</sub> (LoRA)†	0.8M	<b>90.6<math>\pm</math>2</b>	96.2 $\pm$ 5	<b>90.2<math>\pm</math>1.0</b>	68.2 $\pm$ 1.9	<b>94.8<math>\pm</math>3</b>	91.6 $\pm$ 2	<b>85.2<math>\pm</math>1.1</b>	<b>92.3<math>\pm</math>5</b>	<b>88.6</b>
DeB <sub>XXL</sub> (FT)*	1500.0M	91.8	<b>97.2</b>	92.0	72.0	<b>96.0</b>	92.7	93.9	92.9	91.1
DeB <sub>XXL</sub> (LoRA)	4.7M	<b>91.9<math>\pm</math>2</b>	96.9 $\pm$ 2	<b>92.6<math>\pm</math>6</b>	<b>72.4<math>\pm</math>1.1</b>	<b>96.0<math>\pm</math>1</b>	<b>92.9<math>\pm</math>1</b>	<b>94.9<math>\pm</math>4</b>	<b>93.0<math>\pm</math>2</b>	<b>91.3</b>

- Effectiveness
  - Finetuning
  - Domain-adaptive finetuning
- Efficiency
  - Prompt-based learning
  - Parameter-efficient finetuning
- **Conclusion**

- ❑ **Finetuning** requires carefully selected learning rate and output layers.
- ❑ **Domain-adaptive finetuning** is usually helpful when applying the PLM to a new downstream domain.
- ❑ **Prompt-based learning** is a hot new paradigm with little theoretical analysis. How to design better templates and answers is still an open question.
- ❑ **Parameter-efficient finetuning** is also a hot research topic due to its high efficiency and applicability in few-shot settings.

- ❑ How to Fine-Tune BERT for Text Classification? Sun et al. 2019. *arXiv preprint arXiv: 1905.05583*.
- ❑ Improving BERT Fine-Tuning via Self-Ensemble and Self-Distillation. Xu et al. 2020. *arXiv preprint arXiv: 2002.10345*.
- ❑ Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. Gururangan et al. 2020. In *ACL*.
- ❑ Language Models are Few-Shot Learners. Brown et al. 2020. *arXiv preprint arXiv: 2005.14165*.
- ❑ Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. Liu et al. 2021. *arXiv preprint arXiv: 2107.13586*.
- ❑ Making Pre-trained Language Models Better Few-shot Learners. Gao et al. 2021. In *ACL*.
- ❑ Parameter-Efficient Transfer Learning for NLP. Houlsby et al. 2019. In *ICML*.
- ❑ Prefix-Tuning: Optimizing Continuous Prompts for Generation. Li et al. 2021. In *ACL*.
- ❑ The Power of Scale for Parameter-Efficient Prompt Tuning. Lester et al. 2021. In *EMNLP*.
- ❑ LoRA: Low-Rank Adaptation of Large Language Models. Hu et al. 2021. *arXiv preprint arXiv: 2106.09685*.



Thanks  
Q&A