

Making Pre-trained Language Models Better Few-shot Learners

Advisor : Jia-Ling, Koh

Speaker : Hsiao-Ting Huang

Source : ACL'2021

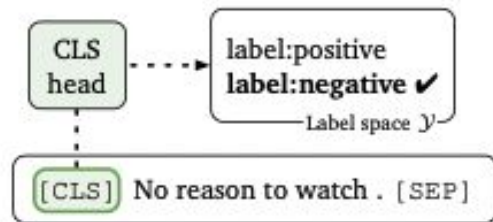
Date : 2023/05/02

Outline

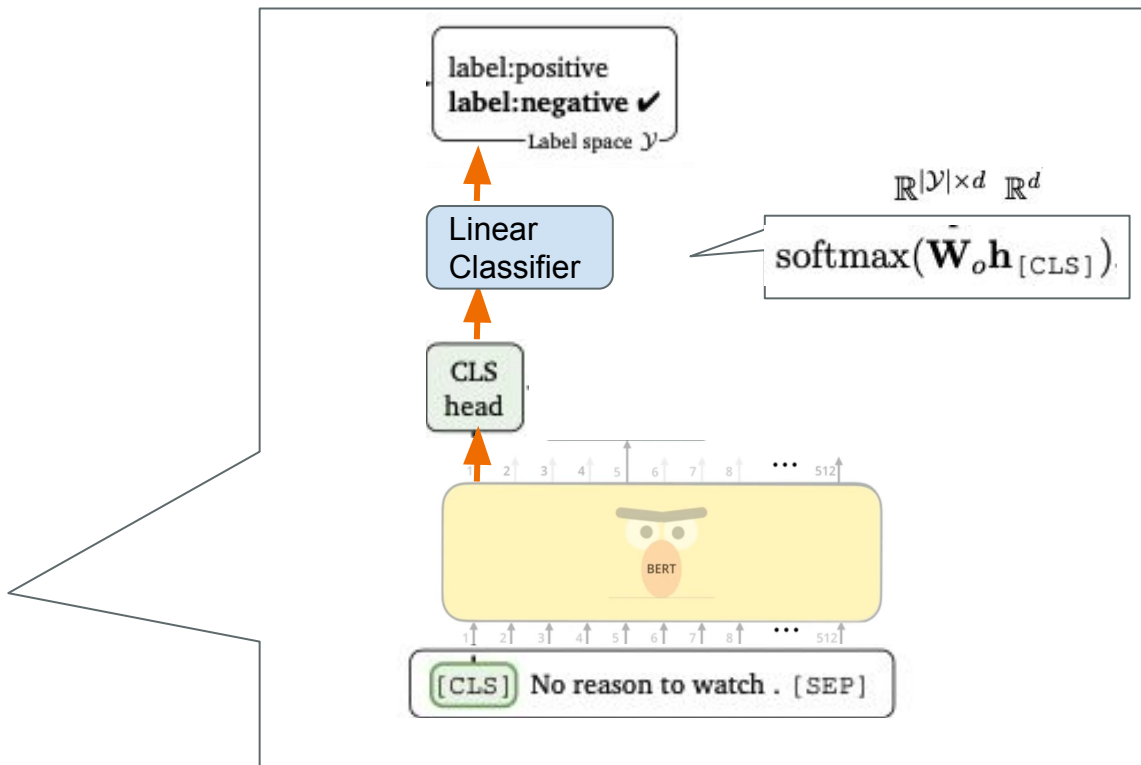
- Introduction
 - ❖ prompt-based fine-tuning
 - ❖ PET
 - ❖ in-context learning from GPT-3
- Method
 - ❖ Automatic Prompt Generation
 - ❖ Fine-tuning with Demonstrations
- Experiment
- Conclusion

Introduction :

- standard fine-tuning

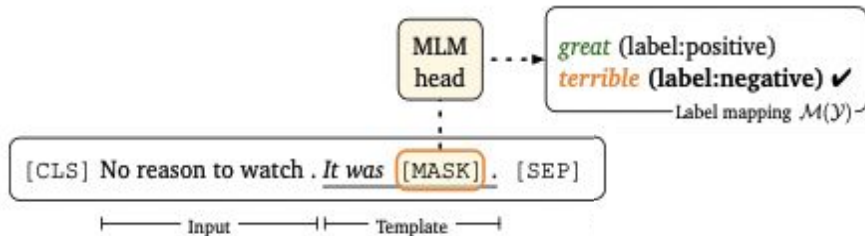


(b) Fine-tuning



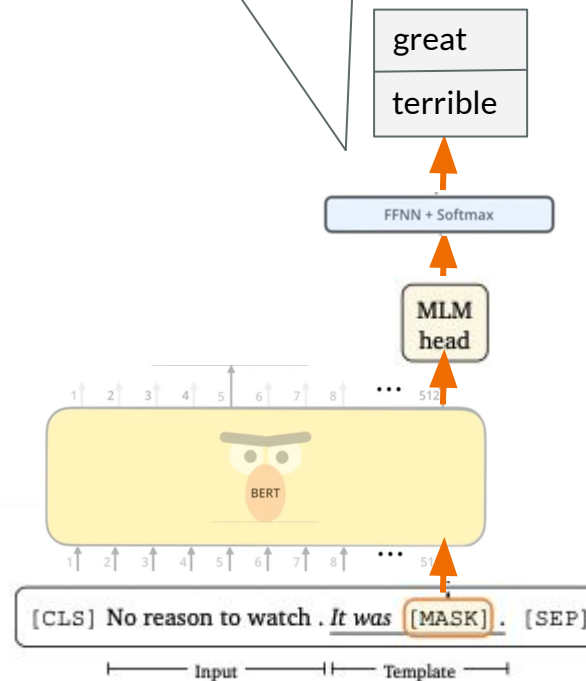
Introduction :

- prompt-based fine-tuning

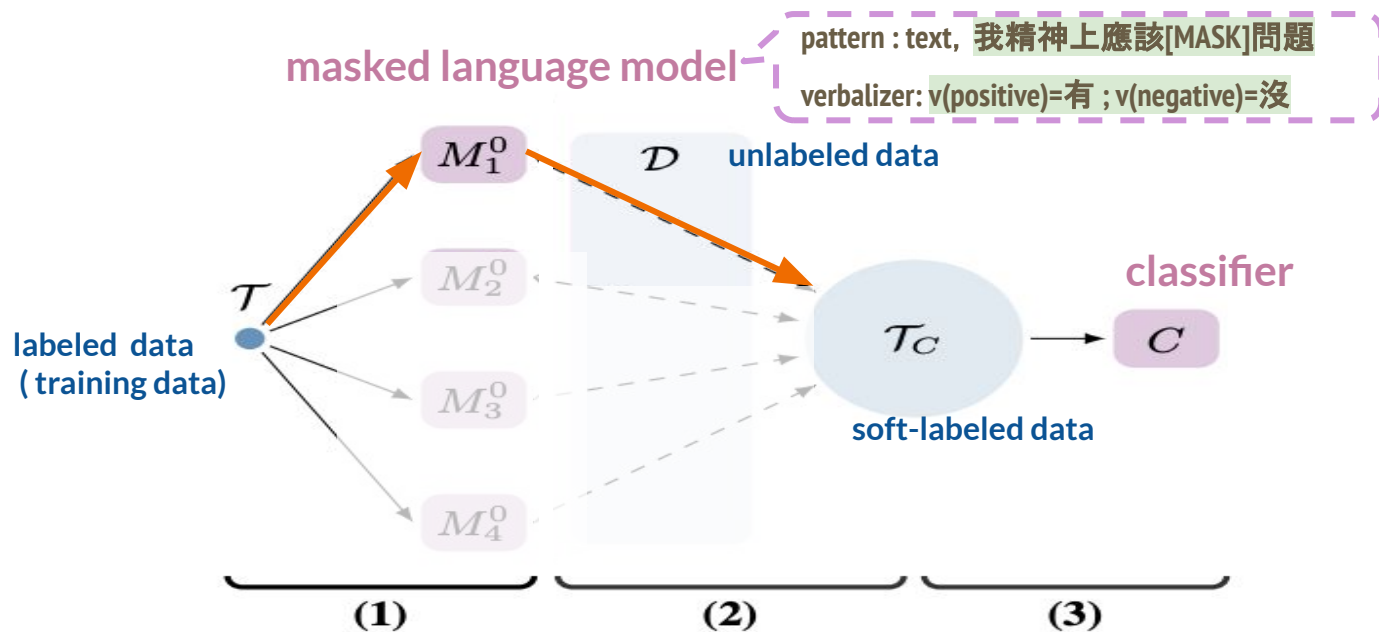


(c) Prompt-based fine-tuning

$$p(y | x_{in}) = p([\text{MASK}] = \mathcal{M}(y) | x_{\text{prompt}}) \\ = \frac{\exp(\mathbf{w}_{\mathcal{M}(y)} \cdot \mathbf{h}_{[\text{MASK}]})}{\sum_{y' \in \mathcal{Y}} \exp(\mathbf{w}_{\mathcal{M}(y')} \cdot \mathbf{h}_{[\text{MASK}]})}$$



Introduction : PET



- problem: Finding the right prompts, however, is an art

Introduction : Manual prompts

Template	Label words	Accuracy
SST-2 (positive/negative)		mean (std)
$\langle S_1 \rangle$ It was [MASK] .	great/terrible	92.7 (0.9)
$\langle S_1 \rangle$ It was [MASK] .	good/bad	92.5 (1.0)
$\langle S_1 \rangle$ It was [MASK] .	cat/dog	91.5 (1.4)
$\langle S_1 \rangle$ It was [MASK] .	dog/cat	86.2 (5.4)
$\langle S_1 \rangle$ It was [MASK] .	terrible/great	83.2 (6.9)
Fine-tuning	-	81.4 (3.8)

sentiment-classification

SNLI (entailment/neutral/contradiction)		mean (std)
$\langle S_1 \rangle$? [MASK] , $\langle S_2 \rangle$	Yes/Maybe/No	77.2 (3.7)
$\langle S_1 \rangle$. [MASK] , $\langle S_2 \rangle$	Yes/Maybe/No	76.2 (3.3)
$\langle S_1 \rangle$? [MASK] $\langle S_2 \rangle$	Yes/Maybe/No	74.9 (3.0)
$\langle S_1 \rangle$ $\langle S_2 \rangle$ [MASK]	Yes/Maybe/No	65.8 (2.4)
$\langle S_2 \rangle$? [MASK] , $\langle S_1 \rangle$	Yes/Maybe/No	62.9 (4.1)
$\langle S_1 \rangle$? [MASK] , $\langle S_2 \rangle$	Maybe/No/Yes	60.6 (4.8)
Fine-tuning	-	48.4 (4.8)

Natural Language Inference

Introduction : in-context learning from GPT-3

- Using in-context learning of GPT-3 for machine translation.

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

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

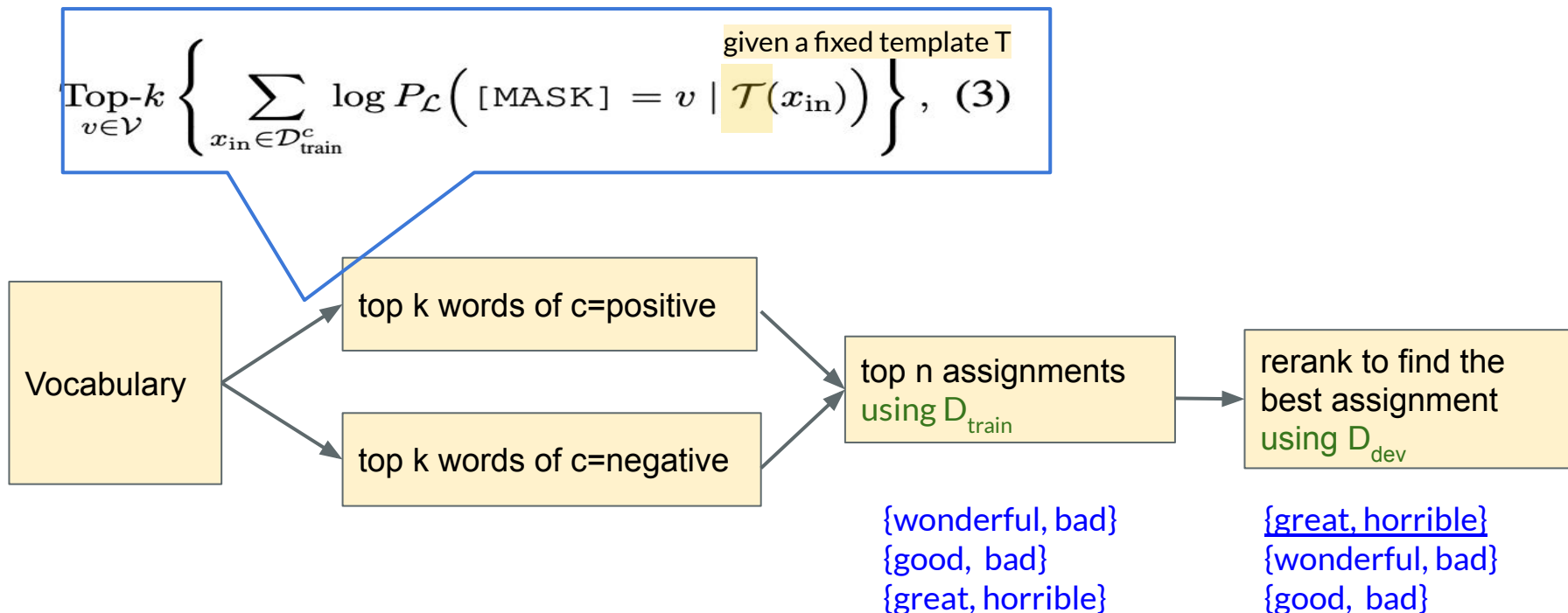
- problem:GPT-3 consists of 175B parameters

Outline

- Introduction
 - ❖ prompt-based fine-tuning
 - ❖ in-context learning from GPT-3
 - ❖ PET
- **Method**
 - ❖ Automatic Prompt Generation
 - ❖ Fine-tuning with Demonstrations
- Experiment
- Conclusion

Automatic Prompt Generation :

- Automatic selection of label words



Automatic Prompt Generation :

- Automatic generation of templates

T5 input:

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

A fun ride. $\langle X \rangle$ **great** $\langle Y \rangle$
 A pleasure to watch. $\langle X \rangle$ **great** $\langle Y \rangle$
 ...
 Training examples for label: **positive**

No reason to watch. $\langle X \rangle$ **terrible** $\langle Y \rangle$
 This junk. $\langle X \rangle$ **terrible** $\langle Y \rangle$
 ...
 Training examples for label: **negative**

positive: **great**, negative: **terrible**
 Label mapping $\mathcal{M}(y)$

fix label word $\mathcal{M}(y)$

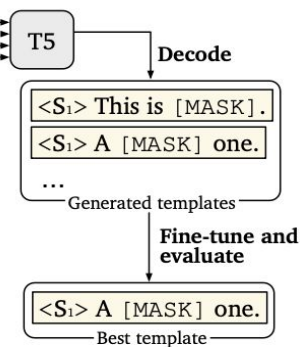
goal of T5 output:

$$\sum_{(x_{in}, y) \in \mathcal{D}_{train}} \log \bar{P}_{T5}(\mathcal{T} | \mathcal{T}_g(x_{in}, y))$$

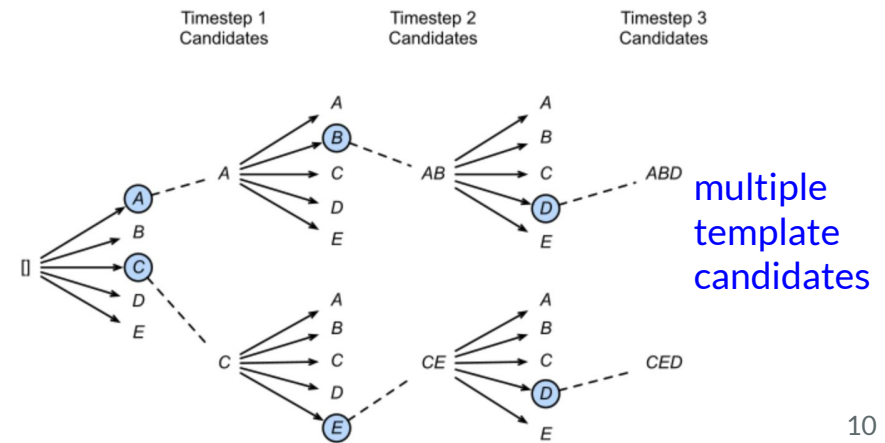


$$\sum_{j=1}^{|\mathcal{T}|} \sum_{(x_{in}, y) \in \mathcal{D}_{train}} \log P_{T5}(t_j | t_1, \dots, t_{j-1}, \mathcal{T}_g(x_{in}, y)), \quad (4)$$

$(t_1, \dots, t_{|\mathcal{T}|})$ are the template tokens.



beam search (beam width = 2)



Fine-tuning with Demonstrations :

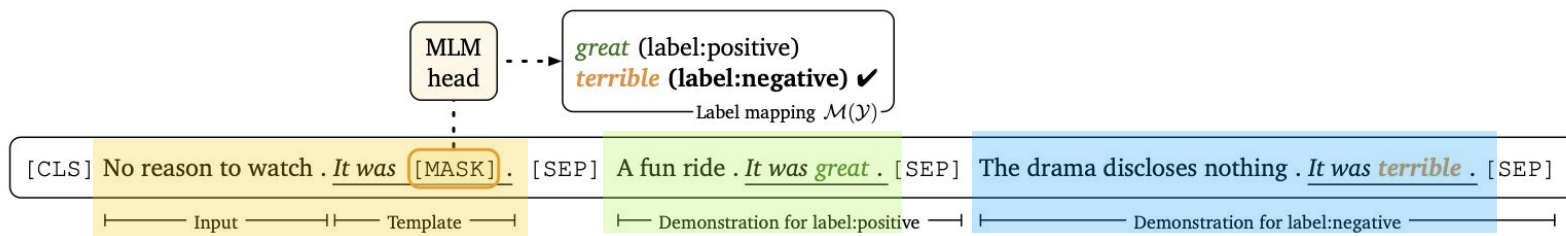
- Sampling similar demonstrations



- Training examples as demonstrations

$$\mathcal{T}(x_{in}) \oplus \tilde{\mathcal{T}}(x_{in}^{(1)}, y^{(1)}) \oplus \dots \oplus \tilde{\mathcal{T}}(x_{in}^{(|\mathcal{Y}|)}, y^{(|\mathcal{Y}|)}).$$

sample from the top $r = 50\%$ instances for each class



(c) Prompt-based fine-tuning with demonstrations (our approach)

Outline

- Introduction
 - ❖ prompt-based fine-tuning
 - ❖ in-context learning from GPT-3
 - ❖ PET
- Method
 - ❖ Automatic Prompt Generation
 - ❖ Fine-tuning with Demonstrations
- Experiment
- Conclusion

Datasets-SST-2

Category	Dataset	$ \mathcal{Y} $	L	#Train	#Test	Type	Labels (classification tasks)
	SST-2	2	19	6,920	872	sentiment	positive, negative

sentence (string)	label (class label)
"hide new secretions from the parental units "	0 (negative)
"contains no wit , only labored gags "	0 (negative)
"that loves its characters and communicates something rather beautiful about human nature "	1 (positive)

- manual prompt:

$\langle S_1 \rangle$ It was [MASK] .	positive: great, negative: terrible
---------------------------------------	-------------------------------------

- auto prompt:

$\langle S_1 \rangle$ A [MASK] one .	irresistible/pathetic
$\langle S_1 \rangle$ A [MASK] piece .	wonderful/bad
$\langle S_1 \rangle$ All in all [MASK] .	delicious/bad

Datasets-TREC

Category	Dataset	$ \mathcal{Y} $	L	#Train	#Test	Type	Labels (classification tasks)
	TREC	6	10	5,452	500	question cls.	abbr., entity, description, human, loc., num.

text (string)	coarse_label label)
"How did serfdom develop in and then leave Russia ?"	2 (DESC)
"What films featured the character Popeye Doyle ?"	1 (ENTY)
"How many points make up a perfect fivepin bowling score ?"	5 (NUM)
"Who was the inventor of silly putty ?"	3 (HUM)
"What is the highest waterfall in the United States ?"	4 (LOC)
"What does the abbreviation AIDS stand for ?"	0 (ABBR)

- manual prompt:

[MASK] : $\langle S_1 \rangle$

abbreviation: Expression, entity: Entity, description: Description
human: Human, location: Location, numeric: Number

- auto prompt:

Q: [MASK] : $\langle S_1 \rangle$
 $\langle S_1 \rangle$ Why [MASK] ?
 $\langle S_1 \rangle$ Answer: [MASK] .

Application/Advisor/Discussion/Culture/Assignment/Minute
Production/AE/Context/Artist/Assignment/Minute
Personality/Advisor/Conclusion/Hum/Assignment/Minute

Datasets-MNLI

Category	Dataset	$ \mathcal{Y} $	L	#Train	#Test	Type	Labels (classification tasks)
	TREC	6	10	5,452	500	question cls.	abbr., entity, description, human, loc., num.

sentence1	sentence2	label
"Fun for adults and children."	"Fun for only children."	2 (contradiction)
"Issues in Data Synthesis."	"Problems in data synthesis."	0 (entailment)

- manual prompt:

$\langle S_1 \rangle$? [MASK] , $\langle S_2 \rangle$	entailment: Yes, natural: Maybe, contradiction: No
--	--

- auto prompt:

$\langle S_1 \rangle$. [MASK] , you are right , $\langle S_2 \rangle$	Fine/Plus/Otherwise
$\langle S_1 \rangle$. [MASK] you're right $\langle S_2 \rangle$	There/Plus/Otherwise
$\langle S_1 \rangle$. [MASK] ! $\langle S_2 \rangle$	Meaning/Plus/Otherwise

Experiment

	SST-2 (acc)	TREC (acc)	MNLI (acc)	SNLI (acc)	RTE (acc)	MRPC (F1)
Majority [†]	50.9	18.8	32.7	33.8	52.7	81.2
Prompt-based zero-shot [‡]	83.6	32.0	50.8	49.5	51.3	61.9
“GPT-3” in-context learning	84.8 (1.3)	26.2 (2.4)	52.0 (0.7)	47.1 (0.6)	60.4 (1.4)	45.7 (6.0)
Fine-tuning	81.4 (3.8)	88.8 (2.1)	45.8 (6.4)	48.4 (4.8)	54.4 (3.9)	76.6 (2.5)
Prompt-based FT (man) + demonstrations	92.7 (0.9)	84.8 (5.1)	68.3 (2.3)	77.2 (3.7)	69.1 (3.6)	74.5 (5.3)
Prompt-based FT (auto) + demonstrations	92.6 (0.5)	87.5 (3.2)	70.7 (1.3)	79.7 (1.5)	68.7 (2.3)	77.8 (2.0)
Prompt-based FT (auto) + demonstrations	92.3 (1.0)	88.2 (2.0)	68.3 (2.5)	77.1 (2.1)	73.9 (2.2)	76.2 (2.3)
Prompt-based FT (auto) + demonstrations	93.0 (0.6)	89.4 (1.7)	70.0 (3.6)	77.5 (3.5)	71.1 (5.3)	78.1 (3.4)
Fine-tuning (full) [†]	95.0	97.4	89.8	92.6	80.9	91.4

Experiment - ensemble model

	Prompt-based Fine-tuning	MNLI	RTE
manual prompt	Our single manual \mathcal{P}	68.3 (2.3)	69.1 (3.6)
	\mathcal{P}_{PET}	71.9 (1.5)	69.2 (4.0)
auto prompt	$\mathcal{P}_{\text{ours}}, \mathcal{P}_{\text{ours}} = \mathcal{P}_{\text{PET}} $ + demonstrations	70.4 (3.1)	73.0 (3.2)
	$\mathcal{P}_{\text{ours}}, \mathcal{P}_{\text{ours}} = 20$ + demonstrations	74.0 (1.9)	71.9 (4.6)
	$\mathcal{P}_{\text{ours}}, \mathcal{P}_{\text{ours}} = 20$ + demonstrations	72.7 (2.5)	73.1 (3.3)
		75.4 (1.6)	72.3 (4.5)

Table 4: Ensemble models using manual prompts from PET (Schick and Schütze, 2021a,b) and our automatic templates. PET uses 4 prompts for MNLI and 5 for RTE. We also use an equal number of templates in $|\mathcal{P}_{\text{ours}}| = |\mathcal{P}_{\text{PET}}|$ for a fair comparison.

Experiment - manual prompts vs. auto prompt

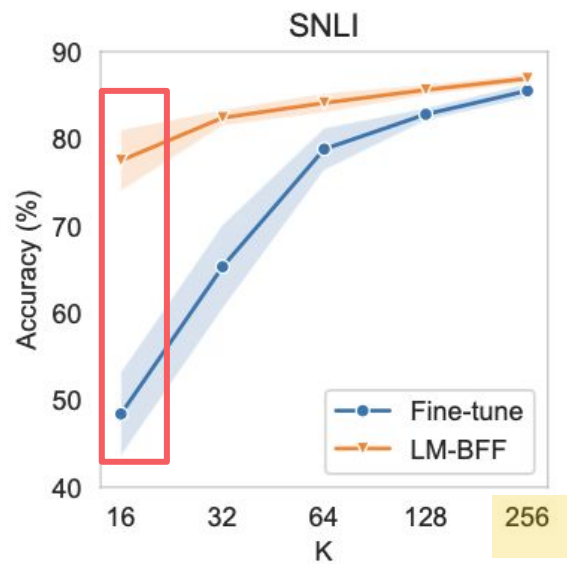
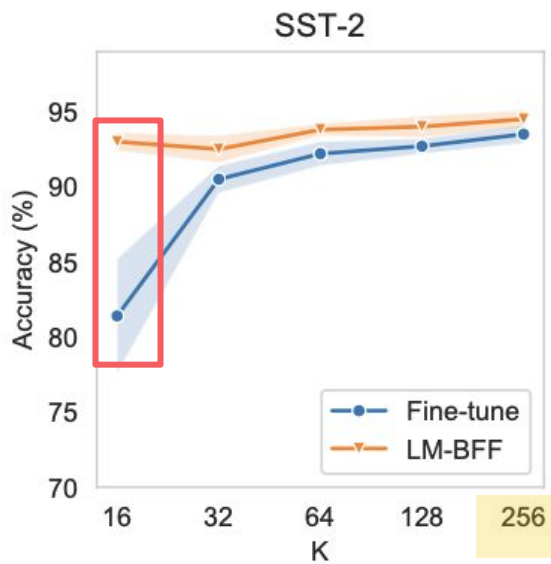
		SST-2	SNLI	TREC	MRPC
manual prompt	Manual	92.7	77.2	84.8	74.5
	Auto T	92.3	77.1	88.2	76.2
auto prompt	Auto L	91.5	75.6	87.0	77.2
	Auto T + L	92.1	77.0	89.2	74.0

SST-2	(positive/negative)
Auto T	$\mathcal{M}(\mathcal{Y}) = \{\text{great, terrible}\}$ #1. $\langle S_1 \rangle$ A [MASK] one . #2. $\langle S_1 \rangle$ A [MASK] piece . #3. $\langle S_1 \rangle$ All in all [MASK] .
Auto L	$\mathcal{T}(x_{in}) = \langle S_1 \rangle$ It was [MASK] . #1. irresistible/pathetic #2. wonderful/bad #3. delicious/bad
SNLI	(entailment/neutral/contradiction)
Auto T	$\mathcal{M}(\mathcal{Y}) = \{\text{Yes, Maybe, No}\}$ #1. $\langle S_1 \rangle$. [MASK] no , $\langle S_2 \rangle$ #2. $\langle S_1 \rangle$. [MASK] , in this case $\langle S_2 \rangle$ #3. $\langle S_1 \rangle$. [MASK] this time $\langle S_2 \rangle$
Auto L	$\mathcal{T}(x_{in}) = \langle S_1 \rangle$? [MASK] , $\langle S_2 \rangle$ #1. Alright/Watch/Except #2. Hi/Watch/Worse #3. Regardless/Fortunately/Unless

Experiment - Impact of demonstration sampling strategies

	SST-2	SNLI	TREC	MRPC
Prompt-based FT (man)	92.7	77.2	84.8	74.5
<u>random sample</u> for each class	92.3	78.8	85.6	70.9
sample from the <u>top r = 50%</u> for each class	92.7	79.5	83.4	76.6
	92.6	79.7	87.5	77.8

Experiment - fine-tuning vs our LM-BFF



Conclusion

- presented LM-BFF, a set of simple but effective techniques for fine-tuning language models using only a few examples
- (1) use prompt-based fine-tuning with automatically searched prompts
- (2) include selected task demonstrations (training examples) as part of the input context.