

# **Prompting for few-shot learning**

AntNLP 2022 Fall Seminar  
Ding Xuanwen

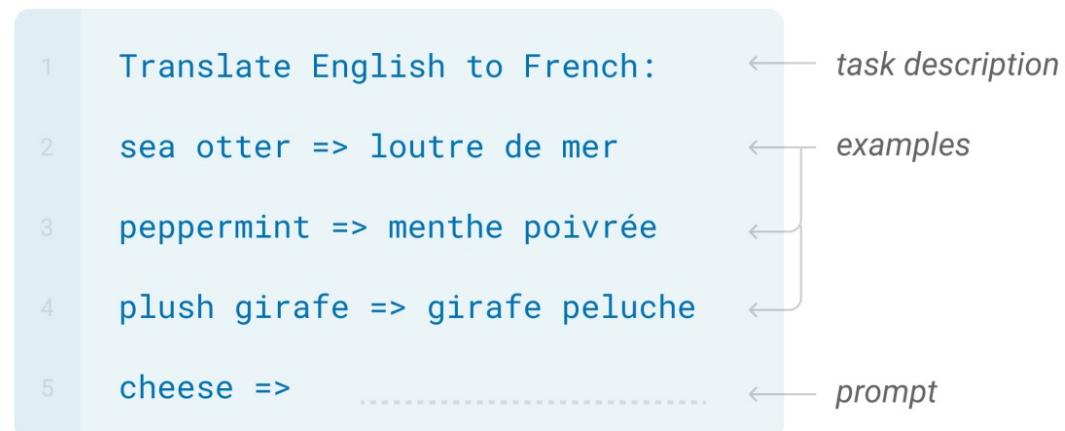
# 主要内容

- Background
- PET
- How Many Data Points is a Prompt Worth?
- LM-BFF
- Null-prompt
- True Few-Shot Learning with Language Models

# Few-shot for GPT3

## Few-shot

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



Almost state-of-the art on  
(some) SuperGLUE tasks with  
only 32 labeled examples!

## Few-shot for GPT3

### GPT

Natural language prompts  
gigantic model  
few in-context examples  
no-parameters updated

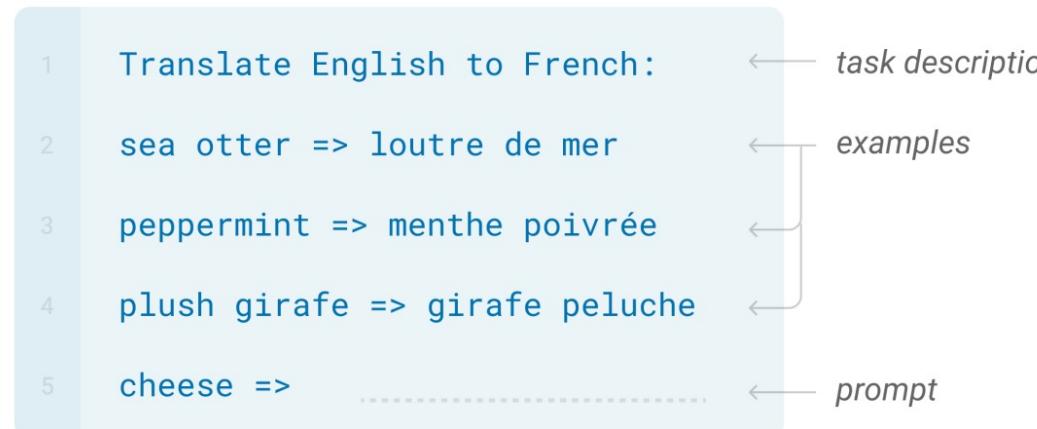
### BERT

110M Parameters  
1000x smaller than GPT3  
Generic [CLS] Token  
Fine-tuned to 2.5 to 400k examples for GLUE tasks

# Few-shot for GPT3

## Few-shot

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



👎 Needs huge models

👎 Limited number of labeled example

# PET : pattern-exploiting traning

## Exploiting Cloze Questions for Few Shot Text Classification and Natural Language Inference

Timo Schick<sup>1,2</sup> Hinrich Schütze<sup>1</sup>

<sup>1</sup> Center for Information and Language Processing, LMU Munich, Germany

<sup>2</sup> Sulzer GmbH, Munich, Germany

schickt@cis.lmu.de

### Abstract

Some NLP tasks can be solved in a fully unsupervised fashion by providing a pretrained language model with “task descriptions” in natural language (e.g., Radford et al., 2019). While this approach underperforms its supervised counterpart, we show in this work that the two ideas can be combined: We introduce Pattern-Exploiting Training (PET), a semi-supervised training procedure that reformulates input examples as cloze-style phrases to help language models understand a given task. These phrases are then used to assign soft labels to a large set of unlabeled examples. Finally, standard supervised training is performed on the resulting training set. For several tasks and languages, PET outperforms supervised training and strong semi-supervised approaches in low-resource settings by a large margin.<sup>1</sup>

### 1 Introduction

Learning from examples is the predominant approach for many NLP tasks: A model is trained

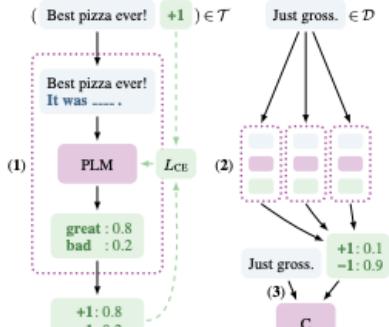


Figure 1: PET for sentiment classification. (1) A number of patterns encoding some form of task description are created to convert training examples to cloze questions; for each pattern, a pretrained language model is finetuned. (2) The ensemble of trained models annotates unlabeled data. (3) A classifier is trained on the resulting soft-labeled dataset.

👍 Needs normal-sized models

👎 Needs huge models

👍 Unlimited number of labeled examples

👎 Limited number of labeled example

# Pattern & verbalizer

<Notation>

- $M$  : Masked Language Model
- $V$  : Vocabulary
- $\_ \in V$  : Mask Token
- $\mathcal{L}, A$  : Labels, Target classification task
- $x = (s_1, \dots, s_k)$  : Input
- $P$  : Pattern
- $v : \mathcal{L} \rightarrow V$
- $(P, v)$  : Pattern-verbalizer pair(PVP)

<Example>

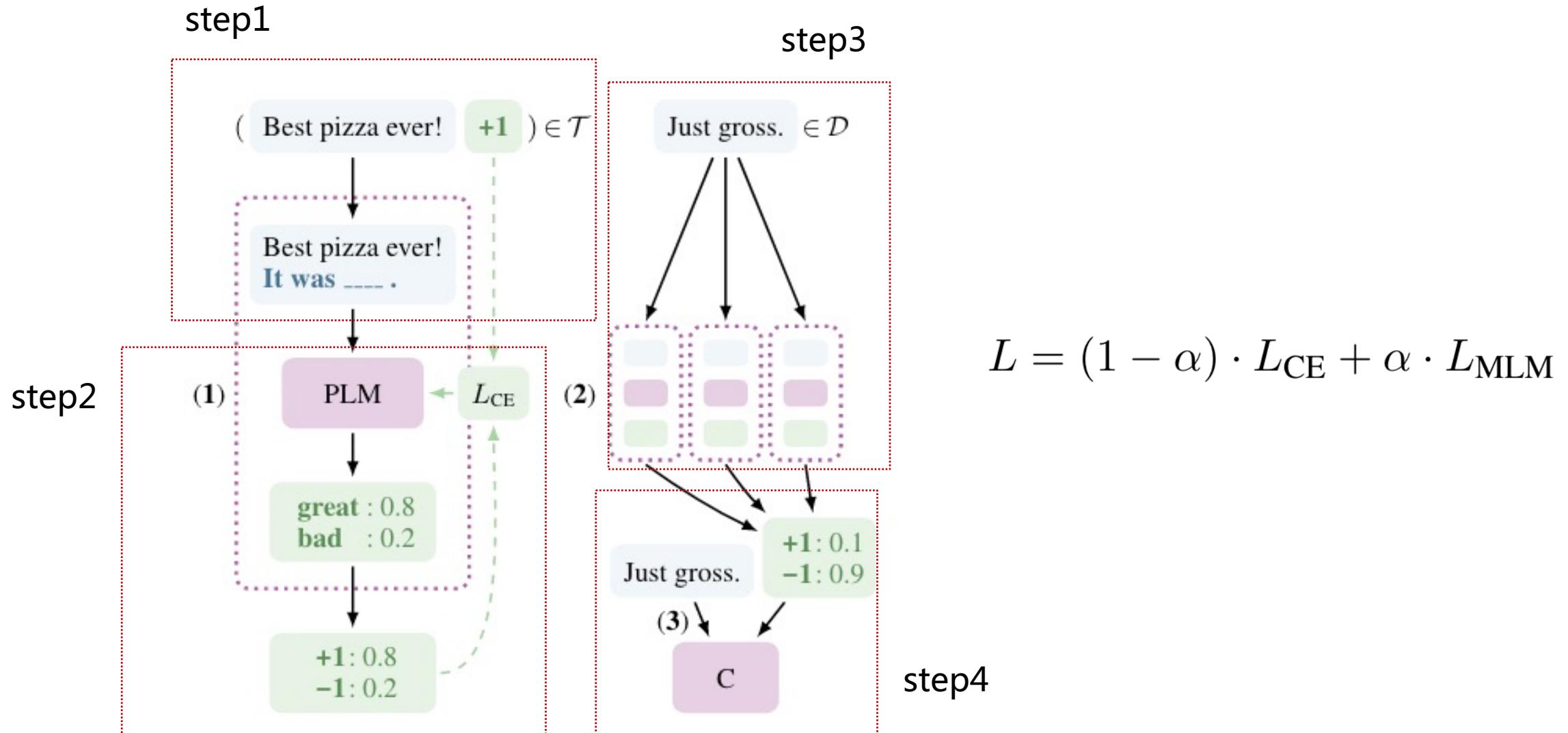
$$P(a, b) = a? \_, b.$$

$$v: y_0 \rightarrow Yes, y_1 \rightarrow No$$

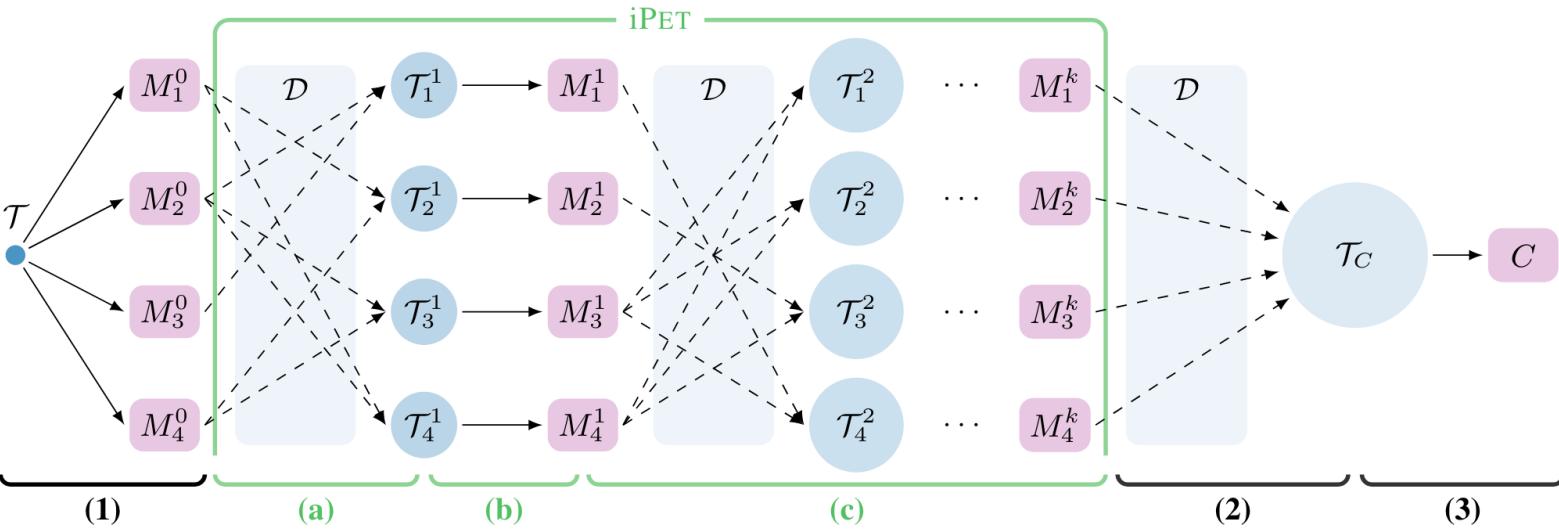
$x = (\text{Mia likes pies}, \text{Mia hates pie})$

$P(x) = \text{Mia likes pie? \_, Mia hates pie.}$

# Overview



# PET & iPET



(1) 设定初始模型集

$$M^0 = \{M_1^0, \dots, M_n^0\}$$

(a) 创建一个soft-labeled dataset

$$s_{\mathbf{p}}(l \mid \mathbf{x}) = M(v(l) \mid P(\mathbf{x}))$$

score :

$$q_{\mathbf{p}}(l \mid \mathbf{x}) = \frac{e^{s_{\mathbf{p}}(l \mid \mathbf{x})}}{\sum_{l' \in \mathcal{L}} e^{s_{\mathbf{p}}(l' \mid \mathbf{x})}} \quad s_{\mathcal{M}}(l \mid \mathbf{x}) = \frac{1}{Z} \sum_{\mathbf{p} \in \mathcal{P}} w(\mathbf{p}) \cdot s_{\mathbf{p}}(l \mid \mathbf{x})$$

$$\text{模型选择 : } \mathcal{N} \subset \mathcal{M}^{j-1} \setminus \{M_i^{j-1}\}$$

$$\text{产生dataset: } \mathcal{T}_{\mathcal{N}} = \{(\mathbf{x}, \arg \max_{l \in \mathcal{L}} s_{\mathcal{N}}(l \mid \mathbf{x})) \mid \mathbf{x} \in \mathcal{D}\}$$

$$\text{新dataset大小 : } c_j(l) = d \cdot c_{j-1}(l)$$

## Experiments

**Dataset :** Yelp Reviews , AG' s News , Yahoo Questions , MNLI , x-stance

### Yelp Reviews

$$P_1(a) = \text{It was _____. } a \quad P_2(a) = \text{Just ____! } \| a$$

$$P_3(a) = a. \text{ All in all, it was _____.}$$

$$P_4(a) = a \| \text{In summary, the restaurant is _____.}$$

$$\begin{aligned} v(1) &= \text{terrible} & v(2) &= \text{bad} & v(3) &= \text{okay} \\ v(4) &= \text{good} & v(5) &= \text{great} \end{aligned}$$

### AG' s News

$$P_1(\mathbf{x}) = \text{____: } a \ b \quad P_2(\mathbf{x}) = a (\text{____}) b$$

$$P_3(\mathbf{x}) = \text{____ - } a \ b \quad P_4(\mathbf{x}) = a \ b (\text{____})$$

$$P_5(\mathbf{x}) = \text{____ News: } a \ b$$

$$P_6(\mathbf{x}) = [\text{Category: } \text{____}] a \ b$$

We use a verbalizer that maps 1–4 to “World”, “Sports”, “Business” and “Tech”, respectively.

# Experiments

Line	Examples	Method	Yelp	AG's	Yahoo	MNLI (m/mm)
1		unsupervised (avg)	$33.8 \pm 9.6$	$69.5 \pm 7.2$	$44.0 \pm 9.1$	$39.1 \pm 4.3 / 39.8 \pm 5.1$
2	$ \mathcal{T}  = 0$	unsupervised (max)	$40.8 \pm 0.0$	$79.4 \pm 0.0$	$56.4 \pm 0.0$	$43.8 \pm 0.0 / 45.0 \pm 0.0$
3		iPET	<b><math>56.7 \pm 0.2</math></b>	<b><math>87.5 \pm 0.1</math></b>	<b><math>70.7 \pm 0.1</math></b>	<b><math>53.6 \pm 0.1 / 54.2 \pm 0.1</math></b>
4		supervised	$21.1 \pm 1.6$	$25.0 \pm 0.1$	$10.1 \pm 0.1$	$34.2 \pm 2.1 / 34.1 \pm 2.0$
5	$ \mathcal{T}  = 10$	PET	$52.9 \pm 0.1$	$87.5 \pm 0.0$	$63.8 \pm 0.2$	$41.8 \pm 0.1 / 41.5 \pm 0.2$
6		iPET	<b><math>57.6 \pm 0.0</math></b>	<b><math>89.3 \pm 0.1</math></b>	<b><math>70.7 \pm 0.1</math></b>	<b><math>43.2 \pm 0.0 / 45.7 \pm 0.1</math></b>
7		supervised	$44.8 \pm 2.7$	$82.1 \pm 2.5$	$52.5 \pm 3.1$	$45.6 \pm 1.8 / 47.6 \pm 2.4$
8	$ \mathcal{T}  = 50$	PET	$60.0 \pm 0.1$	$86.3 \pm 0.0$	$66.2 \pm 0.1$	$63.9 \pm 0.0 / 64.2 \pm 0.0$
9		iPET	<b><math>60.7 \pm 0.1</math></b>	<b><math>88.4 \pm 0.1</math></b>	<b><math>69.7 \pm 0.0</math></b>	<b><math>67.4 \pm 0.3 / 68.3 \pm 0.3</math></b>
10		supervised	$53.0 \pm 3.1$	$86.0 \pm 0.7$	$62.9 \pm 0.9$	$47.9 \pm 2.8 / 51.2 \pm 2.6$
11	$ \mathcal{T}  = 100$	PET	$61.9 \pm 0.0$	$88.3 \pm 0.1$	$69.2 \pm 0.0$	$74.7 \pm 0.3 / 75.9 \pm 0.4$
12		iPET	<b><math>62.9 \pm 0.0</math></b>	<b><math>89.6 \pm 0.1</math></b>	<b><math>71.2 \pm 0.1</math></b>	<b><math>78.4 \pm 0.7 / 78.6 \pm 0.5</math></b>
13		supervised	$63.0 \pm 0.5$	<b><math>86.9 \pm 0.4</math></b>	$70.5 \pm 0.3$	$73.1 \pm 0.2 / 74.8 \pm 0.3$
14	$ \mathcal{T}  = 1000$	PET	<b><math>64.8 \pm 0.1</math></b>	<b><math>86.9 \pm 0.2</math></b>	<b><math>72.7 \pm 0.0</math></b>	<b><math>85.3 \pm 0.2 / 85.5 \pm 0.4</math></b>

Table 1: Average accuracy and standard deviation for RoBERTa (large) on Yelp, AG's News, Yahoo and MNLI (m:matched/mm:mismatched) for five training set sizes  $|\mathcal{T}|$ .

Ex.	Method	Yelp	AG's	Yahoo	MNLI
$ \mathcal{T}  = 10$	UDA	27.3	72.6	36.7	34.7
	MixText	20.4	81.1	20.6	32.9
	PET	48.8	84.1	59.0	39.5
	iPET	<b>52.9</b>	<b>87.5</b>	<b>67.0</b>	<b>42.1</b>
$ \mathcal{T}  = 50$	UDA	46.6	83.0	60.2	40.8
	MixText	31.3	84.8	61.5	34.8
	PET	55.3	86.4	63.3	55.1
	iPET	<b>56.7</b>	<b>87.3</b>	<b>66.4</b>	<b>56.3</b>

Table 2: Comparison of PET with two state-of-the-art semi-supervised methods using RoBERTa (base)

# Experiments

Examples	Method	De	Fr	It
$ \mathcal{T}  = 1000$	supervised	43.3	49.5	41.0
	PET	<b>66.4</b>	<b>68.7</b>	<b>64.7</b>
$ \mathcal{T}  = 2000$	supervised	57.4	62.1	52.8
	PET	<b>69.5</b>	<b>71.7</b>	<b>67.3</b>
$ \mathcal{T}  = 4000$	supervised	63.2	66.7	58.7
	PET	<b>71.7</b>	<b>74.0</b>	<b>69.5</b>
$\mathcal{T}_{\text{De}}, \mathcal{T}_{\text{Fr}}$	supervised	76.6	76.0	71.0
	PET	<b>77.9</b>	<b>79.0</b>	<b>73.6</b>
$\mathcal{T}_{\text{De}} + \mathcal{T}_{\text{Fr}}$	sup. (*)	76.8	76.7	70.2
	supervised	77.6	79.1	75.9
	PET	<b>78.8</b>	<b>80.6</b>	<b>77.2</b>

Table 3: Results on x-stance intra-target for XLM-R (base) trained on subsets of  $\mathcal{T}_{\text{De}}$  and  $\mathcal{T}_{\text{Fr}}$  and for joint training on all data ( $\mathcal{T}_{\text{De}} + \mathcal{T}_{\text{Fr}}$ ). (\*): Best results for mBERT reported in [Vamvas and Sennrich \(2020\)](#).

# Experiments

## Auxiliary Language Modeling

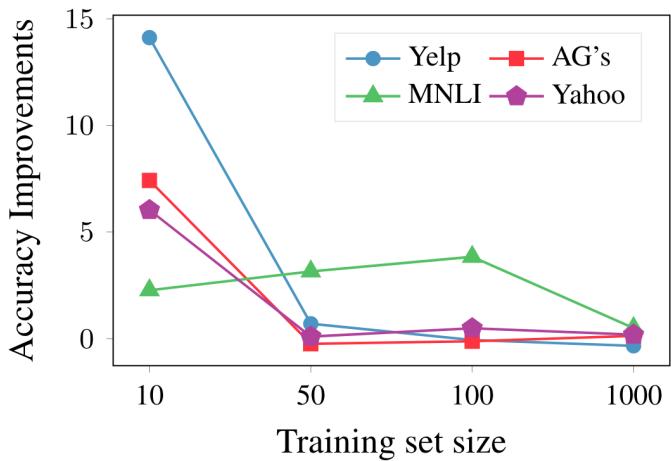


Figure 3: Accuracy improvements for PET due to adding  $L_{MLM}$  during training

## Iterative PET

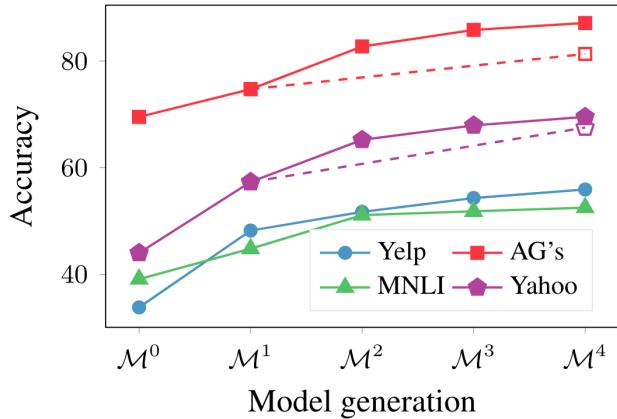


Figure 4: Average accuracy for each generation of models with iPET in a zero-shot setting. Accuracy on AG's News and Yahoo when skipping generation 2 and 3 is indicated through dashed lines.

## In-Domain Pretraining

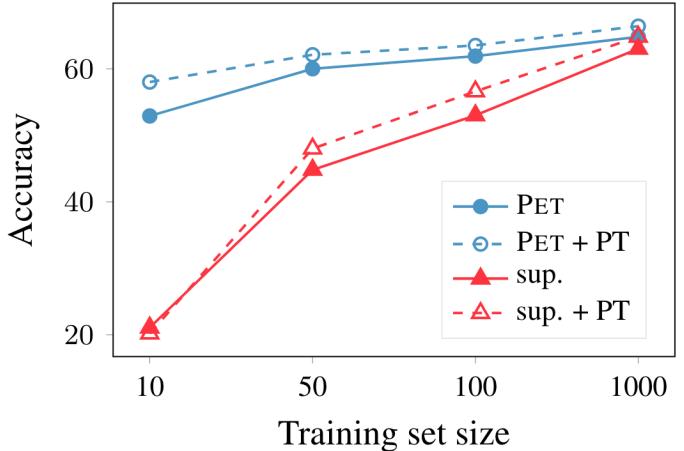


Figure 5: Accuracy of supervised learning (sup.) and PET both with and without pretraining (PT) on Yelp

# **How Many Data Points is a Prompt Worth?**

**Teven Le Scao**

Hugging Face

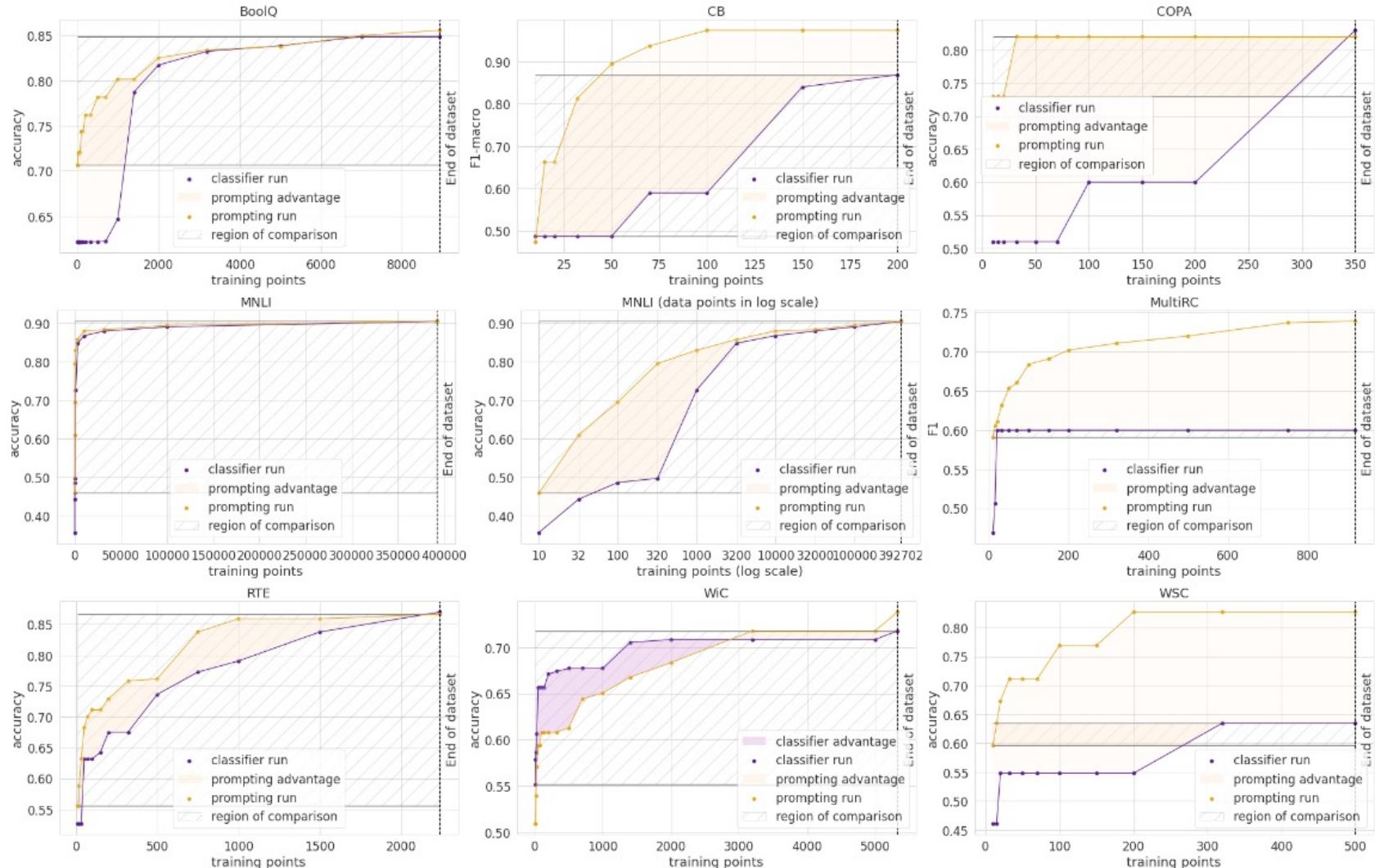
`teven@huggingface.co`

**Alexander M. Rush**

Hugging Face

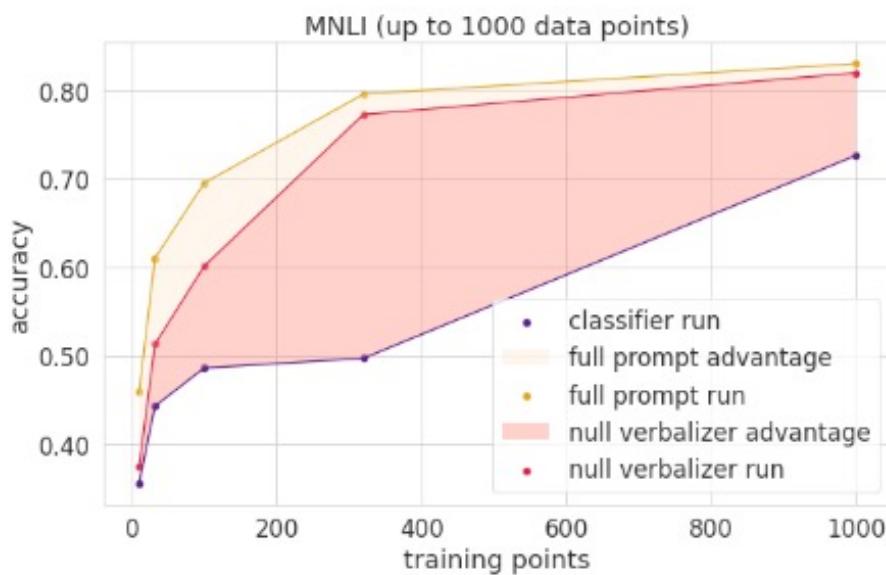
`sasha@huggingface.co`

# Heads vs Prompts

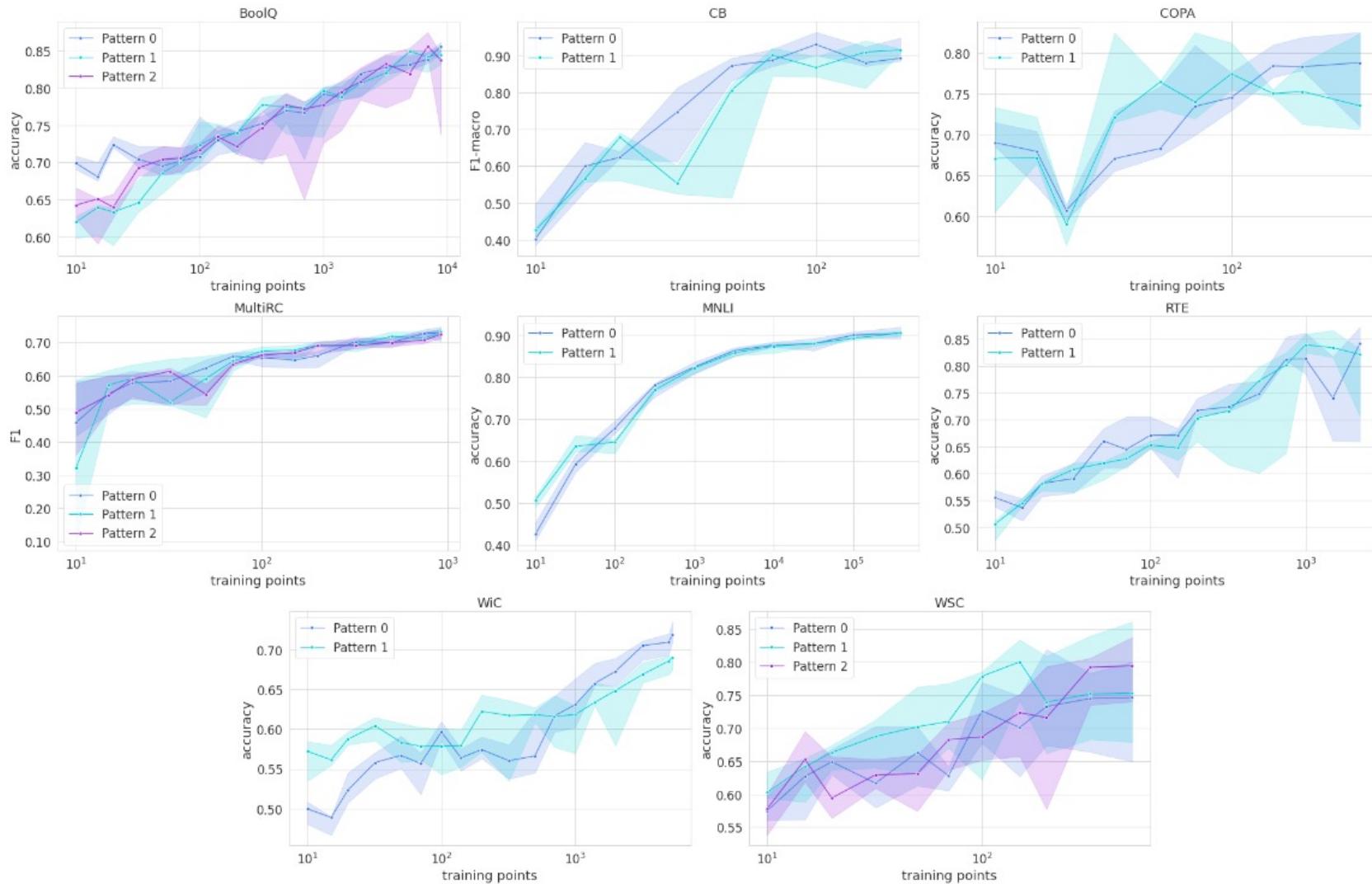


# Pattern vs Verbalizer

	Average Advantage (# Training Points)								
	MNLI	BoolQ	CB	COPA	MultiRC*	RTE	WiC	WSC	
$P \text{ vs } H$	$3506 \pm 536$	$752 \pm 46$	$90 \pm 2$	$288 \pm 242$	$384 \pm 378$	$282 \pm 34$	$-424 \pm 74$	$281 \pm 137$	
$P \text{ vs } N$	$150 \pm 252$	$299 \pm 81$	$78 \pm 2$	-	$74 \pm 56$	$404 \pm 68$	$-354 \pm 166$	-	
$N \text{ vs } H$	$3355 \pm 612$	$453 \pm 90$	$12 \pm 1$	-	$309 \pm 320$	$-122 \pm 62$	$-70 \pm 160$	-	



# Pattern vs Verbalizer



**LM-BFF:** better few-shot fine-tuning of language models

# Making Pre-trained Language Models Better Few-shot Learners

**Tianyu Gao<sup>†\*</sup>**   **Adam Fisch<sup>‡\*</sup>**   **Danqi Chen<sup>†</sup>**

<sup>†</sup>Princeton University   <sup>‡</sup>Massachusetts Institute of Technology

{tianyug, danqic}@cs.princeton.edu  
fisch@csail.mit.edu

## (Small) Language Models are Few-shot Learners



GPT-3

Gigantic language model(up to 175B paramters)

Freezing model parameters

Using manually-designed prompts and demonstrations



PET

Small language model

Fine-tuning model parameters

Using manually-designed prompts



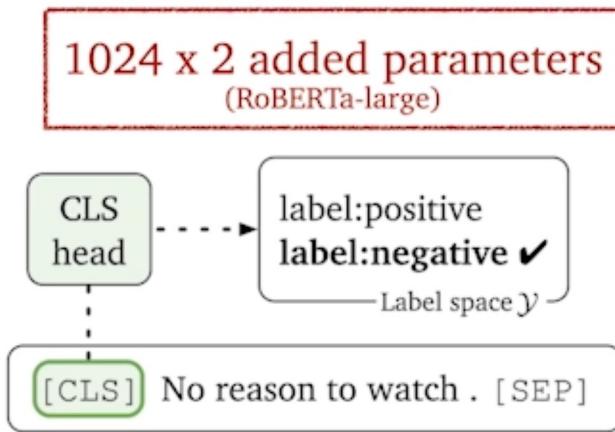
LM-BFF

Small language model

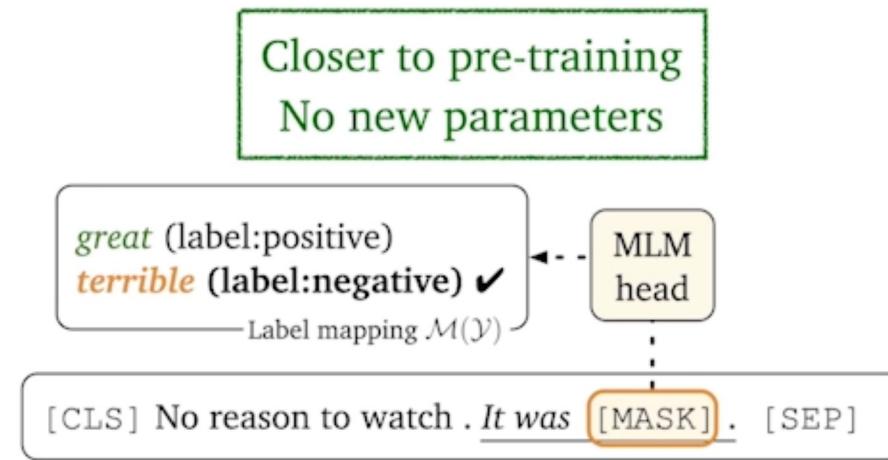
Fine-tuning model parameters

Using **automatically-searched** prompts and **demonstration**

# Prompt-based Fine-tuning



Fine-tuning



Prompt-based fine-tuning

## Manual Label Words : the Good and the Bad

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

## Manual Templates: the Good and the Bad

SNLI (entailment/neutral/contradiction)	mean (std)
$\langle S_1 \rangle ? [MASK] , \langle S_2 \rangle$	Yes/Maybe/No <b>77.2 (3.7)</b>
$\langle S_1 \rangle . [MASK] , \langle S_2 \rangle$	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)

# Automatic Prompt Generation

## Recent Work

1. domain expertise
2. an understanding of the language model's inner workings
3. manual prompts are likely to be suboptimal

## Our goals

1. reduce the human involvement required to design prompts
2. find more optimal settings than those that we manually choose

## Automatic selection of label words

1. For each class  $c$ , select the top  $k$  vocabulary words based on their **conditional likelihood**

Given the manual template:<S> It was [MASK].

$$\text{Top-}k \left\{ \sum_{v \in \mathcal{V}} \log P_{\mathcal{L}} \left( [\text{MASK}] = v \mid \mathcal{T}(x_{\text{in}}) \right) \right\}$$

label : positive  
good  
great  
perfect  
...

label : negative  
awful  
bad  
terrible  
...

2. find the top  $n$  assignments over the pruned space that maximize zero-shot accuracy on  $\mathcal{D}_{train}$

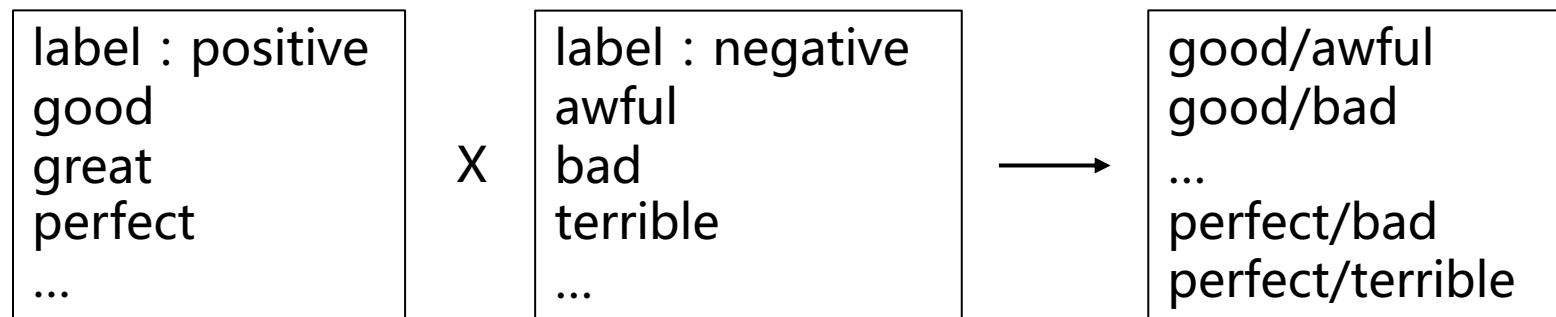
3. fine-tune all top  $n$  assignments, and re-rank to find the best one using  $\mathcal{D}_{dev}$

## Automatic selection of label words

1. For each class  $c$ , select the top  $k$  vocabulary words based on their conditional likelihood

2. find the top  $n$  assignments over the pruned space that maximize **zero-shot** accuracy on  $\mathcal{D}_{train}$

Given the manual template:<S> It was [MASK].

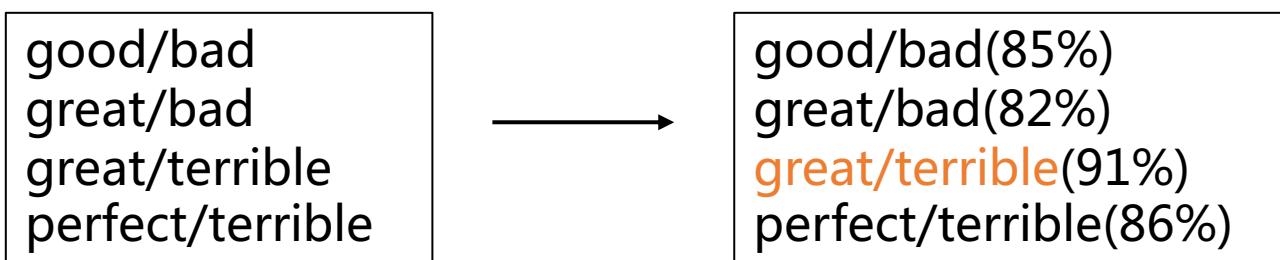


3. fine-tune all top  $n$  assignments, and re-rank to find the best one using  $\mathcal{D}_{dev}$

## Automatic selection of label words

1. For each class  $c$ , select the top  $k$  vocabulary words based on their conditional likelihood
2. find the top  $n$  assignments over the pruned space that maximize zero-shot accuracy on  $\mathcal{D}_{train}$
3. **fine-tune** all top  $n$  assignments, and re-rank to find the best one using  $\mathcal{D}_{dev}$

Given the manual template:<S> It was [MASK].

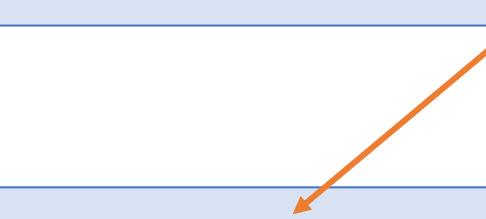


## Automatic generation of templates

T5

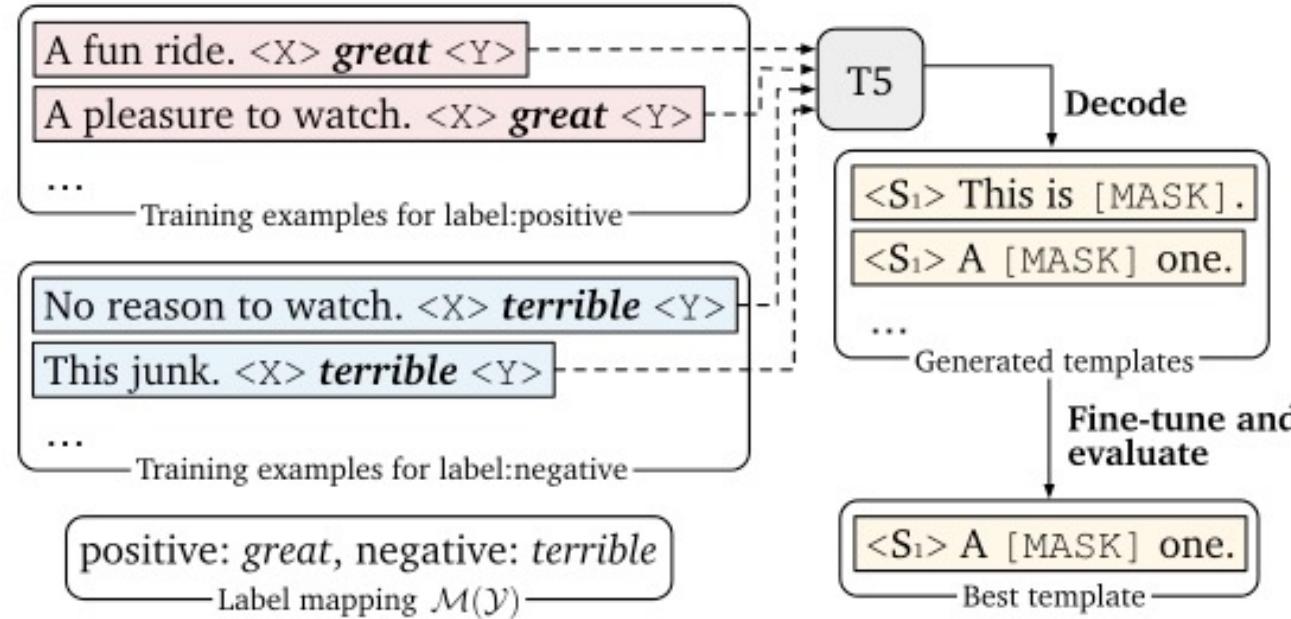
**Input :** Thank you <X> me to your party <Y> week

**Output :** <X> for inviting <Y> last <Z>



T5 takes a **fill-in-the blank** pretraining objective.

# Automatic generation of templates



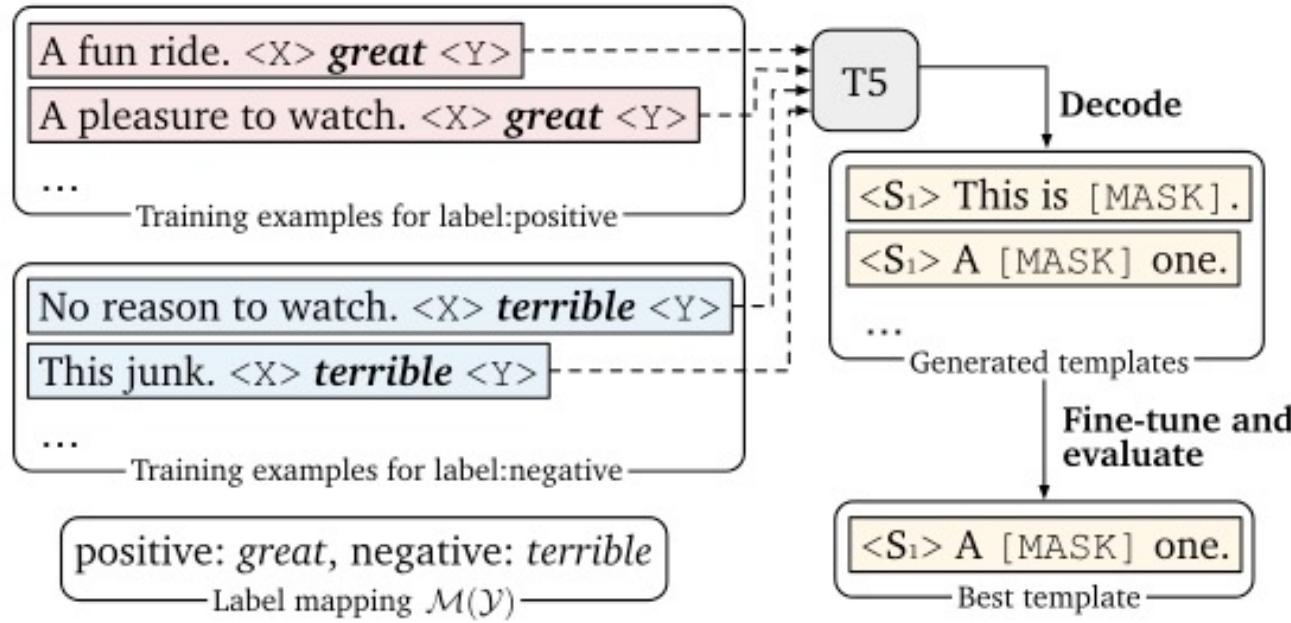
$$\mathcal{T}_g(x_{\text{in}}, y) :$$

$$\begin{aligned} <S_1> &\longrightarrow <X> \mathcal{M}(y) <Y> <S_1>, \\ <S_1> &\longrightarrow <S_1> <X> \mathcal{M}(y) <Y>, \\ <S_1>, <S_2> &\longrightarrow <S_1> <X> \mathcal{M}(y) <Y> <S_2>. \end{aligned}$$

condition the model on all training examples  
at each position in auto-regressive decoding

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

# Automatic generation of templates

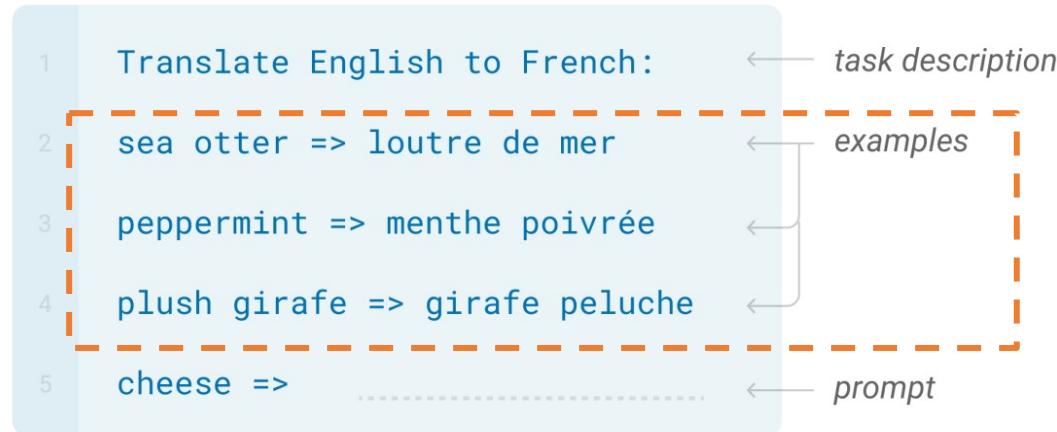


take **beam search** with a large width (=100) to get high-quality templates

# Fine-tuning with Demonstrations

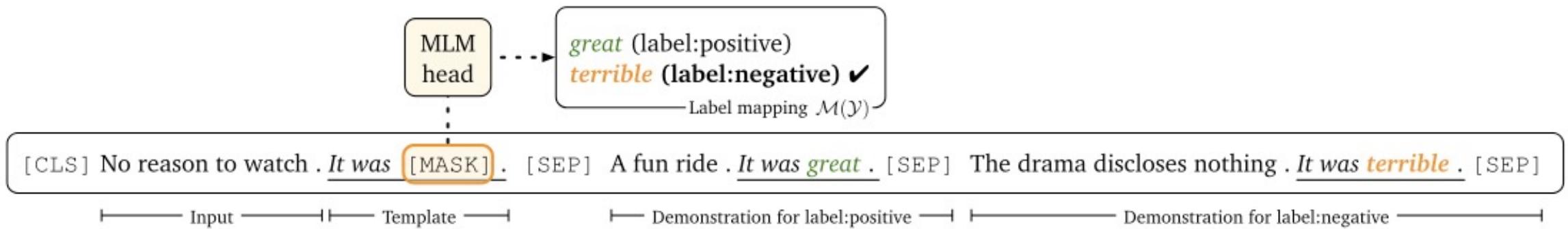
## Few-shot

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



GPT-3 style in-context learning:  
Randomly sample examples and fill them in the context

# Fine-tuning with Demonstrations

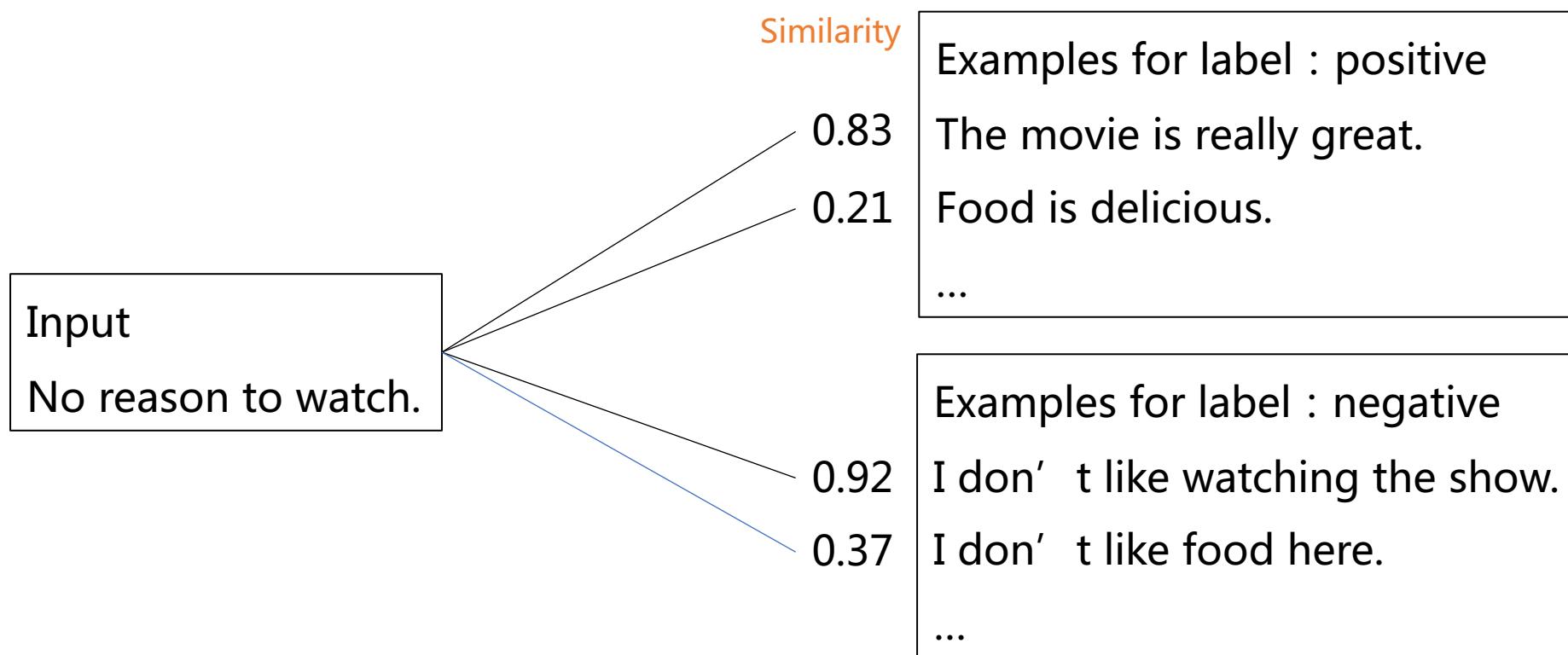


(c) Prompt-based fine-tuning with demonstrations (our approach)

# Selectively Sampling Demonstrations

select demonstrations that are **semantically close** to the input :

measure the cosine similarities between the input and all the training examples  
by using a **pre-trained sentence encoder**(Sentence-BERT)



## Experiments Setup

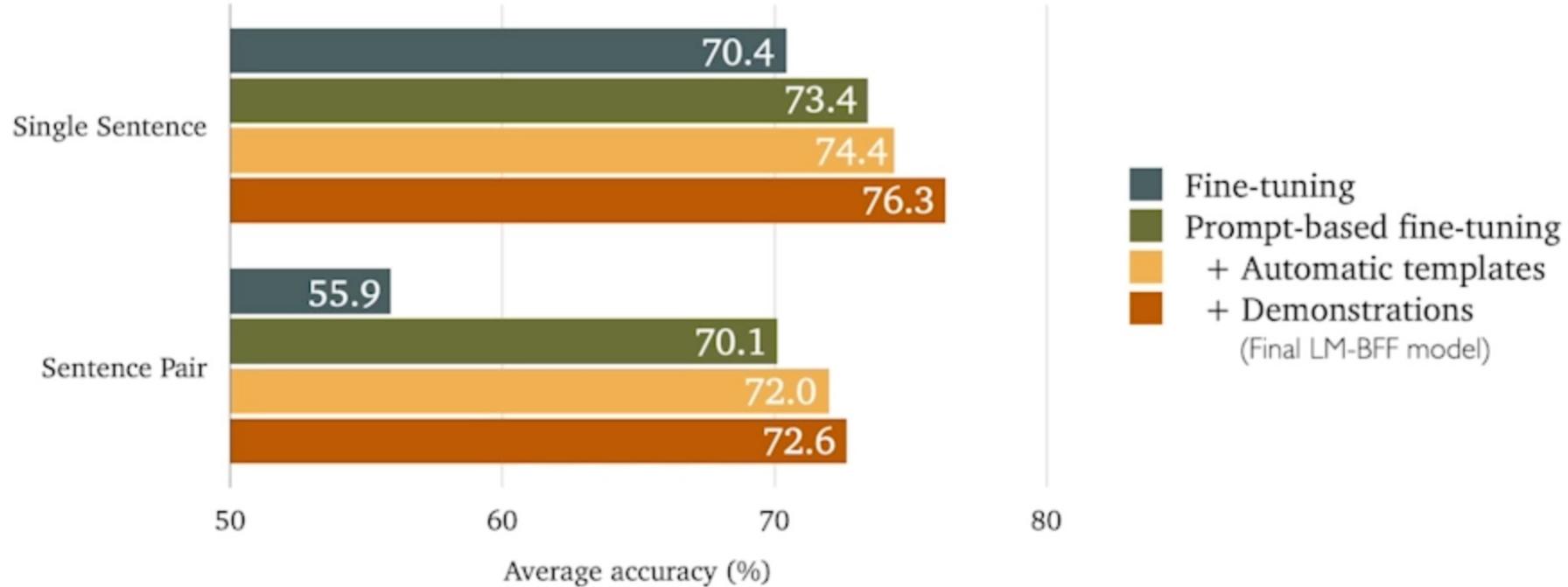
16 examples per class for both training and development sets

- original dev. set : **not true** few-shot learning
- no dev. set : “**shooting in the dark**” for hyper-parameters and early stopping

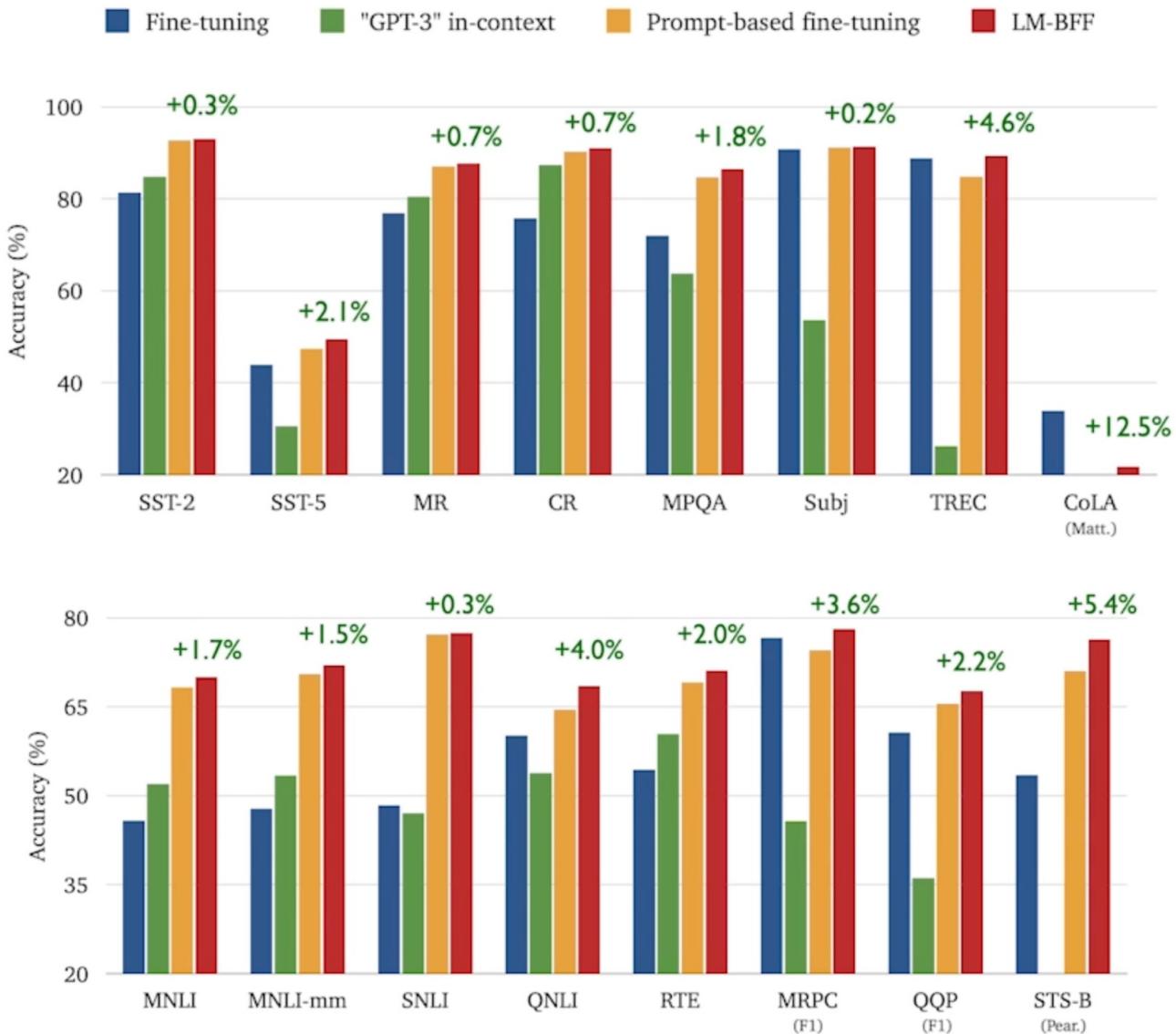
Sample 5 few-shot sets for each dataset and average the results

- performance has **a large variance** given different few-shot sets.

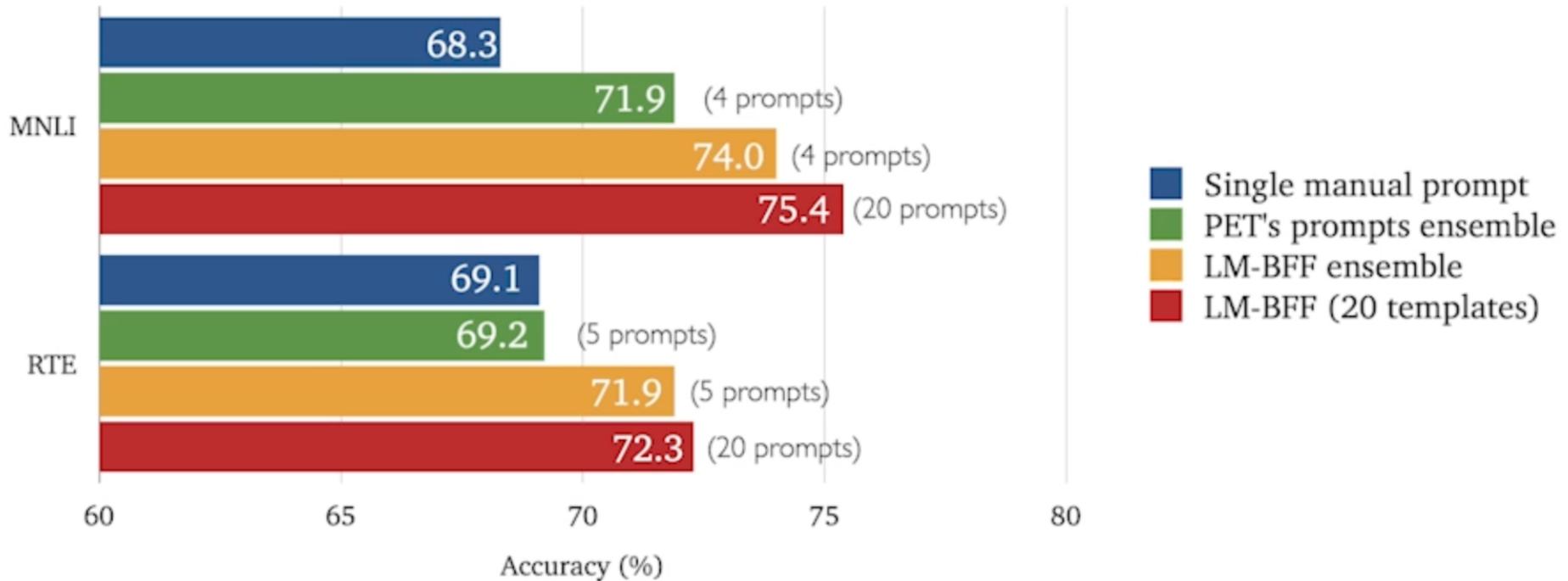
## Main Results: Singal prompts



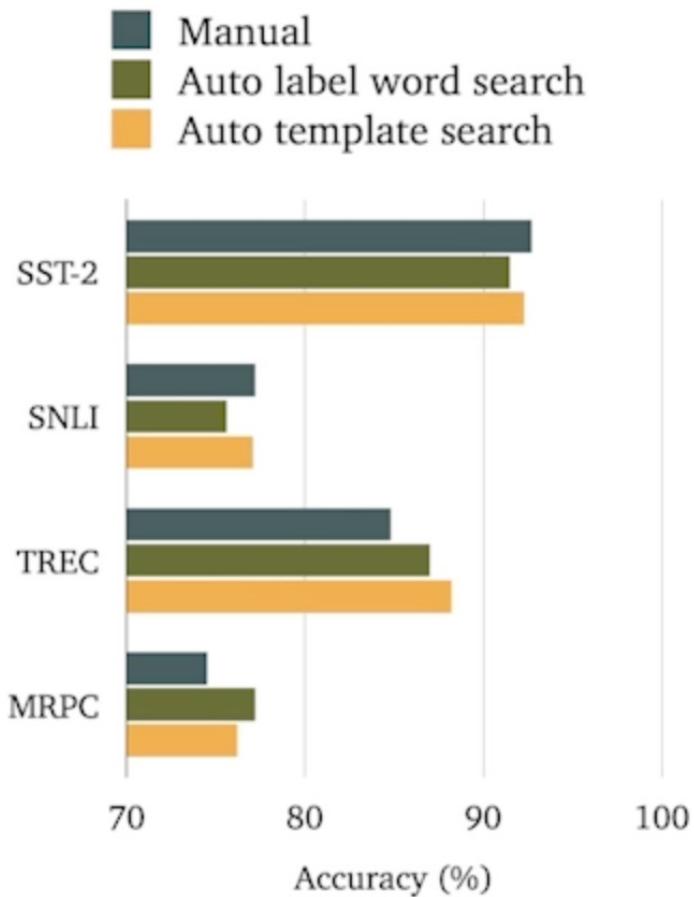
# Main Results: Singal prompts



## Main Results: Ensemble



# Ablation Study: Automatic Prompt Search

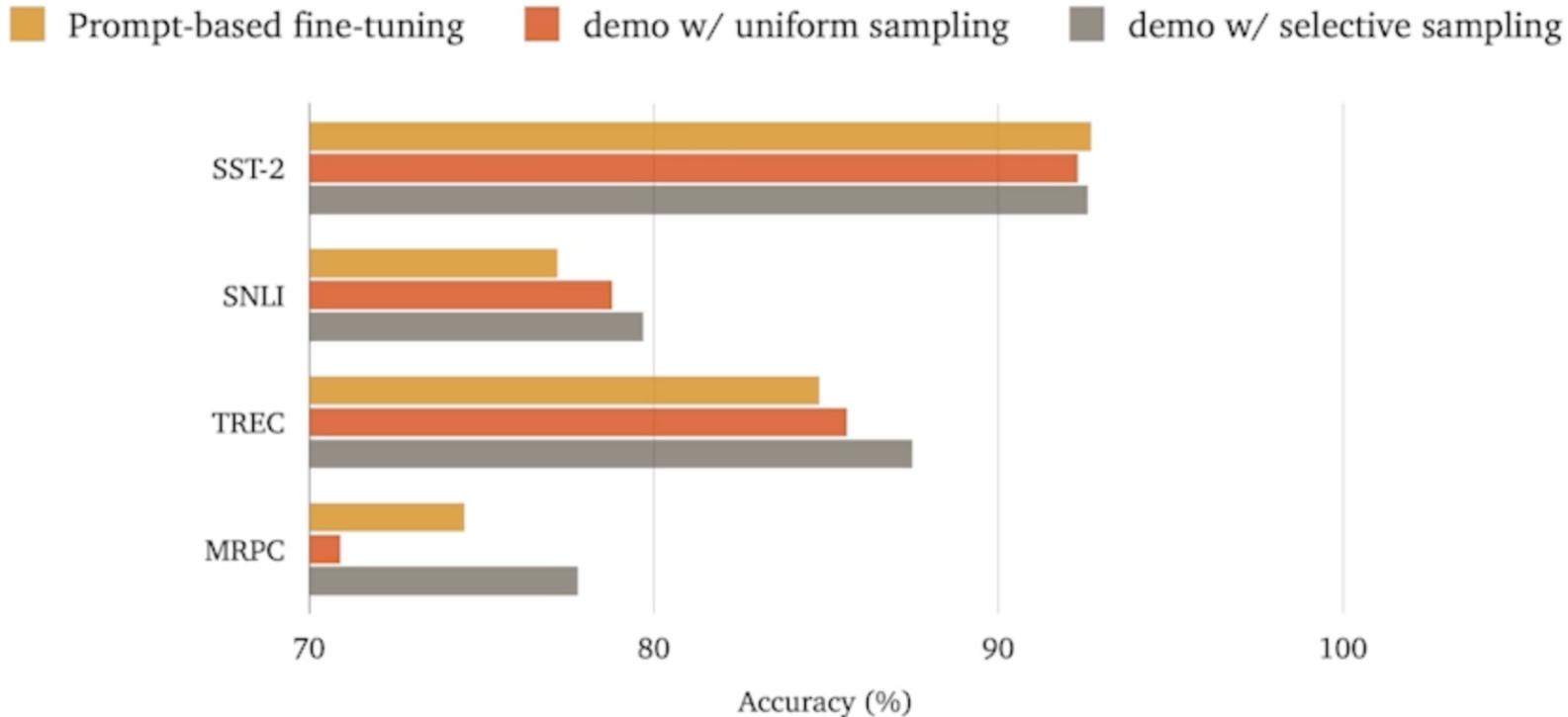


SST-2 (positive/negative)	
$\mathcal{M}(\mathcal{Y}) = \{\text{great, terrible}\}$	
#1. $\langle S_1 \rangle$ A [MASK] one .	
#2. $\langle S_1 \rangle$ A [MASK] piece .	
#3. $\langle S_1 \rangle$ All in all [MASK] .	

SNLI (entailment/neutral/contradiction)	
$\mathcal{M}(\mathcal{Y}) = \{\text{Yes, Maybe, No}\}$	
#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 $\langle S_2 \rangle$	

# Ablation Study: Sampling Demonstrations



## Summary

- LM-BFF
  - Automatically-searched prompt-based fine-tuning
  - A tailored way for incorporating demonstrations
- A rigorous few-shot evaluation on 15 tasks
  - LM-BFF dramatically outperforms standard fine-tuning
- Limitations
  - Large variance
  - Automatic search still relies on either manual templates or label words
  - Favors tasks that can be posed as fill-in-the-blanks

# Cutting Down on Prompts and Parameters: Simple Few-Shot Learning with Language Models

**Robert L. Logan IV<sup>1</sup> Ivana Balažević<sup>\*2</sup> Eric Wallace<sup>3</sup>**

**Fabio Petroni<sup>4</sup> Sameer Singh<sup>1</sup> Sebastian Riedel<sup>4,5</sup>**

<sup>1</sup>UC Irvine    <sup>2</sup>University of Edinburgh

<sup>3</sup>UC Berkeley    <sup>4</sup>Facebook AI Research    <sup>5</sup>University College London

[{rlogan,sameer}@uci.edu](mailto:{rlogan,sameer}@uci.edu)  [ivana.balazevic@ed.ac.uk](mailto:ivana.balazevic@ed.ac.uk)

[ericwallace@berkeley.edu](mailto:ericwallace@berkeley.edu)  [{fabio petroni,sriedel}@fb.com](mailto:{fabio petroni,sriedel}@fb.com)

# Motivation

Data  
Efficiency

Search  
Efficiency

Memory  
Efficiency

**Improve**

**Few-Shot Language Model Performance**

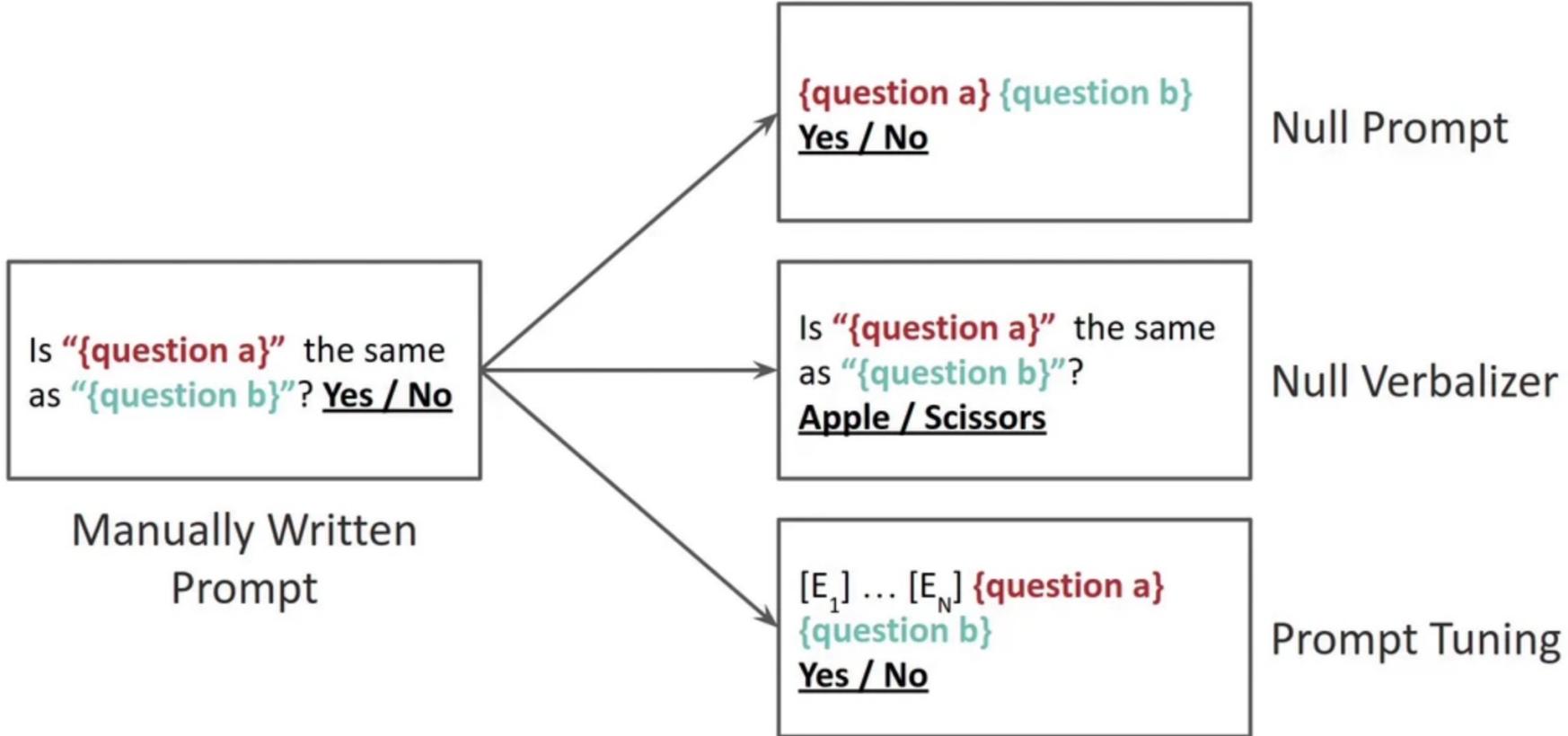
**by**

**Streamlining Prompt Engineering**

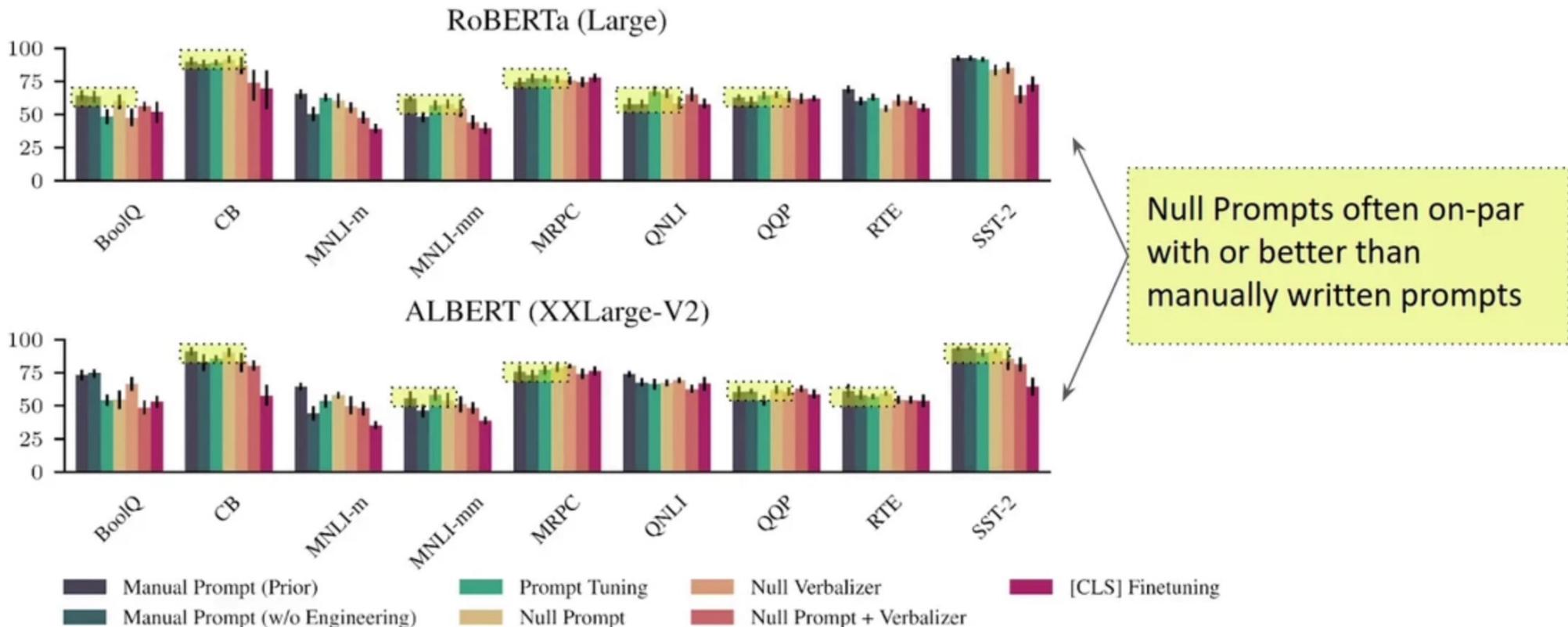
**and**

**Reducing Parameter Counts**

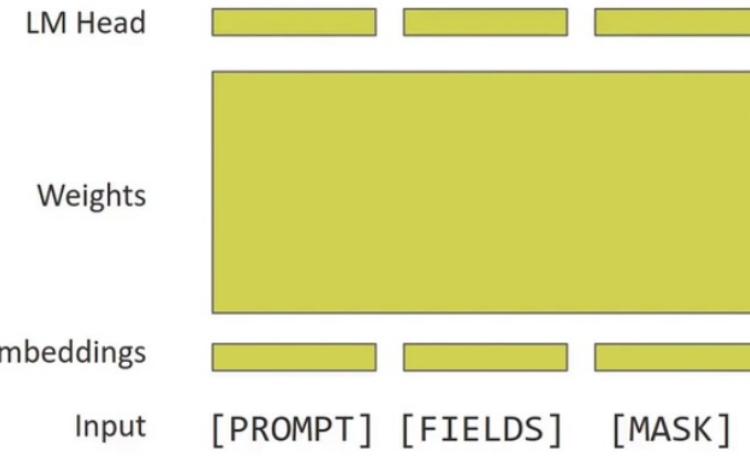
# Simplifying Prompt Engineering



# Simplifying Prompt Engineering



# Improving Parameter Efficiency

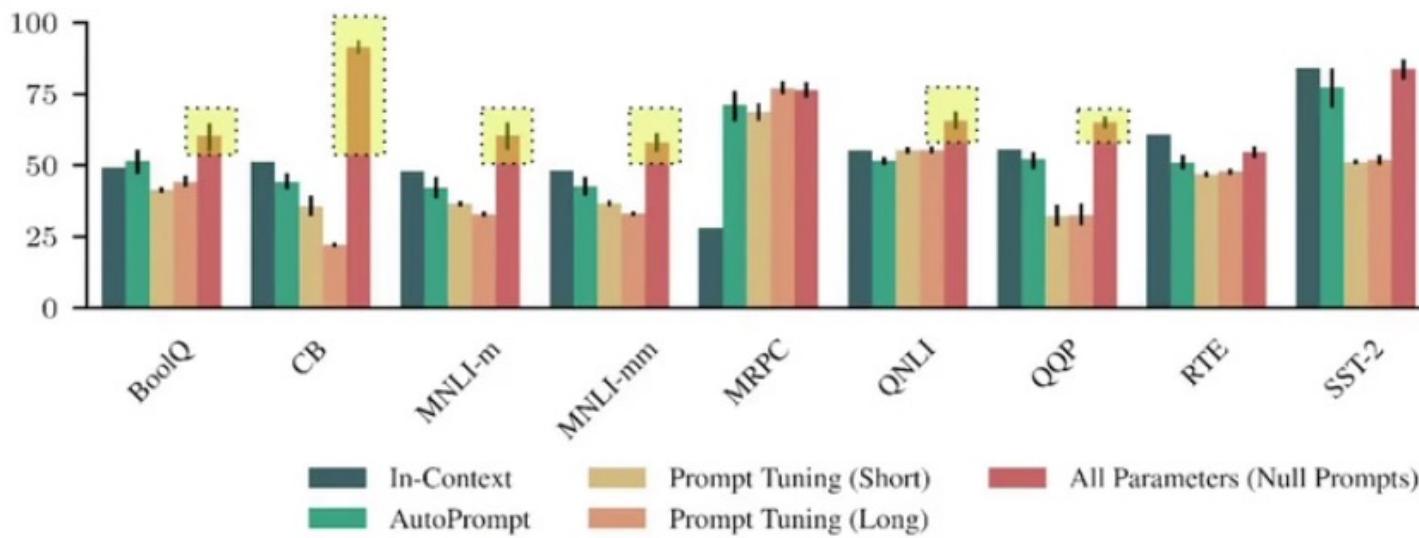


Prompt-Based Finetuning



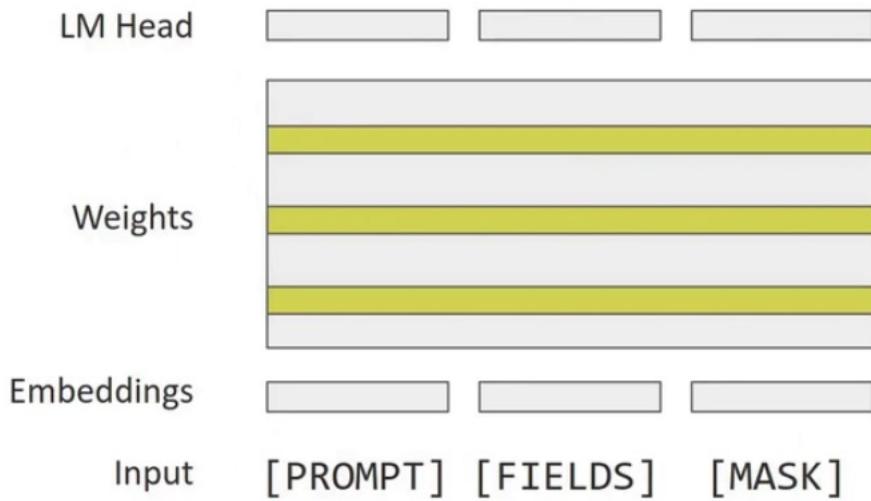
Prompt-Only Tuning

# Improving Parameter Efficiency



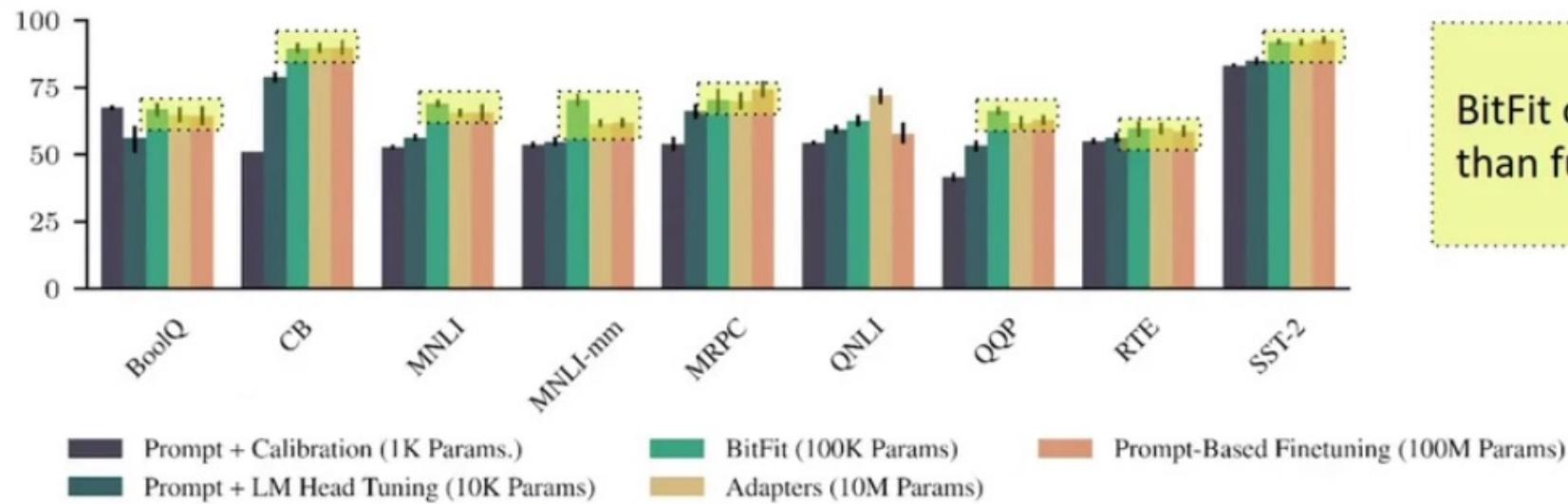
Prompt-Only Tuning  
Substantially Decreases  
Performance

# Improving Parameter Efficiency



BitFit  
&  
Adapters

# Improving Parameter Efficiency



BitFit on-par with or better  
than full finetuning

---

# True Few-Shot Learning with Language Models

---

Ethan Perez<sup>1</sup>, Douwe Kiela<sup>2</sup>, Kyunghyun Cho<sup>13</sup>

<sup>1</sup>New York University, <sup>2</sup>Facebook AI Research,

<sup>3</sup>CIFAR Fellow in Learning in Machines & Brains

perez@nyu.edu

Learning Scenario	Many Train Distributions	Many Train Examples	Many Val Examples
Data-Rich Supervised	✗	✓	✓
Multi-Dist. Few-Shot	✓	✗	✗
Tuned Few-Shot	✗	✗	✓
True Few-Shot	✗	✗	✗

## model selection criteria

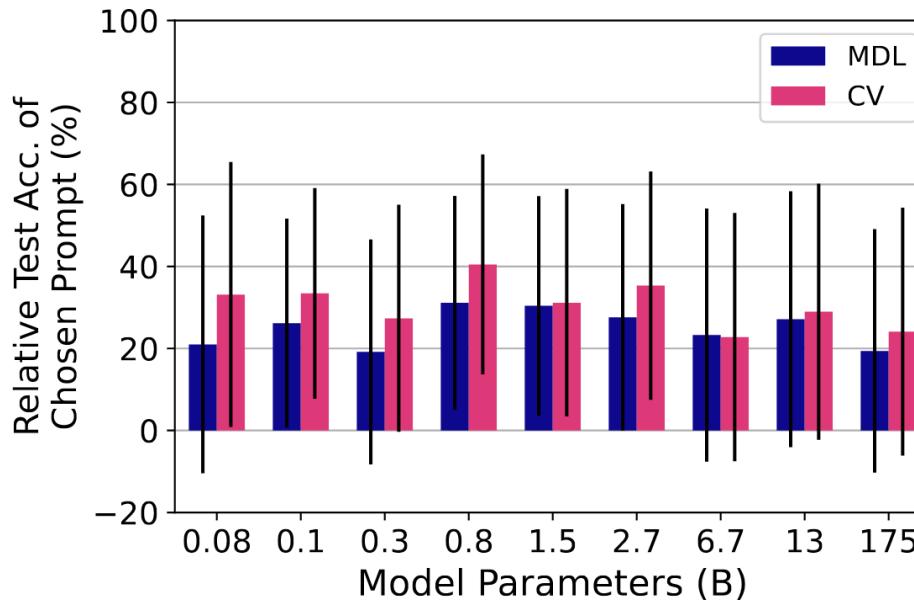
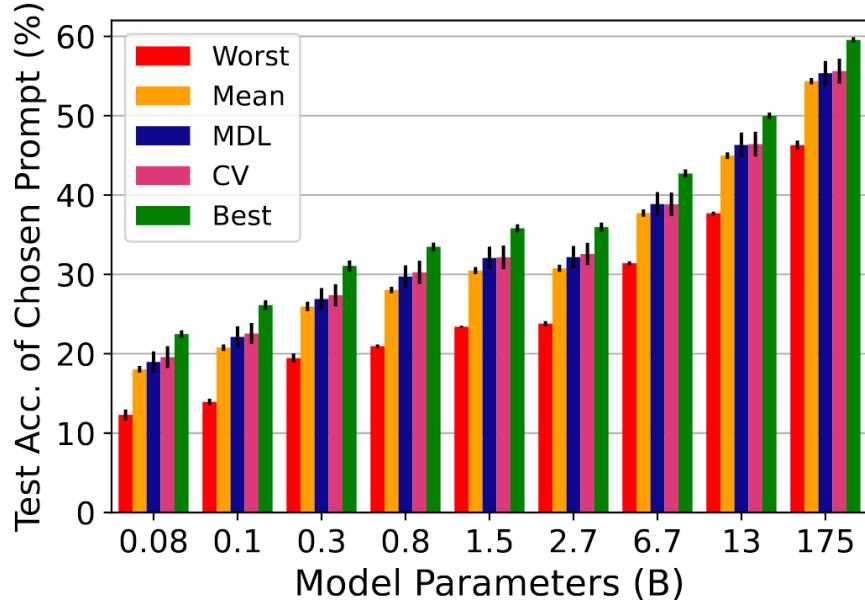
Cross-Validation (CV)

$$\text{CV}(\mathcal{A}, R, F) = \mathbb{E}_{k \sim \text{Unif}(1, K)} \left[ \mathcal{L}\left(\mathcal{A}(F(D_{\text{train}})_{\neg k}, R); F(D_{\text{train}})_k\right) \right]$$

Minimum description length ( MDL )

$$\text{MDL}(\mathcal{A}, R, F) = \mathbb{E}_{k \sim \text{Unif}(1, K)} \left[ \mathcal{L}\left(\mathcal{A}(F(D_{\text{train}})_{1:k-1}, R); F(D_{\text{train}})_k\right) \right]$$

# How well does prompt selection do in true few-shot learning?



# True Few-Shot Hyperparameter Selection

	<b>BoolQ</b> Acc	<b>CB</b> Acc/F1	<b>COPA</b> Acc	<b>RTE</b> Acc	<b>WiC</b> Acc	<b>WSC</b> Acc	<b>MultiRC</b> EM/F1	<b>ReCoRD</b> EM/F1	<b>Avg</b>
<b>Worst</b>	75.0 <sub>4.8</sub>	79.5 <sub>2.3</sub> /67.3 <sub>7.8</sub>	76.8 <sub>2.2</sub>	63.2 <sub>4.0</sub>	49.0 <sub>1.3</sub>	77.2 <sub>1.8</sub>	38.5 <sub>7.4</sub> /80.0 <sub>2.9</sub>	76.2 <sub>1.8</sub> /86.5 <sub>1.2</sub>	69.4 <sub>1.5</sub>
<b>Mean</b>	79.0 <sub>1.5</sub>	85.9 <sub>2.3</sub> /74.5 <sub>11.0</sub>	81.1 <sub>2.9</sub>	70.8 <sub>2.5</sub>	51.5 <sub>1.8</sub>	82.5 <sub>2.7</sub>	44.2 <sub>6.6</sub> /82.3 <sub>2.7</sub>	78.3 <sub>1.3</sub> /87.8 <sub>0.8</sub>	73.9 <sub>1.2</sub>
<b>MDL</b>	76.5 <sub>5.8</sub>	85.7 <sub>5.6</sub> /74.8 <sub>13.4</sub>	82.0 <sub>2.9</sub>	70.4 <sub>8.5</sub>	52.2 <sub>3.0</sub>	82.0 <sub>3.1</sub>	39.7 <sub>8.1</sub> /80.6 <sub>3.2</sub>	78.9 <sub>0.7</sub> /88.2 <sub>0.4</sub>	73.4 <sub>2.8</sub>
<b>CV</b>	78.9 <sub>2.4</sub>	83.9 <sub>5.3</sub> /69.2 <sub>10.3</sub>	80.5 <sub>3.3</sub>	68.7 <sub>7.0</sub>	51.1 <sub>1.6</sub>	83.1 <sub>2.6</sub>	41.9 <sub>7.2</sub> /81.4 <sub>3.1</sub>	78.7 <sub>1.6</sub> /88.1 <sub>1.0</sub>	73.0 <sub>2.1</sub>
<b>Best</b>	80.9 <sub>1.0</sub>	89.8 <sub>3.1</sub> /79.8 <sub>13.4</sub>	84.8 <sub>4.5</sub>	76.7 <sub>1.8</sub>	54.1 <sub>2.3</sub>	86.6 <sub>1.8</sub>	46.8 <sub>6.9</sub> /83.4 <sub>2.9</sub>	80.4 <sub>1.1</sub> /89.2 <sub>0.7</sub>	77.2 <sub>0.9</sub>
<b>ADAPET [42]</b>	80.3	89.3 / 86.8	89.0	76.5	54.4	81.7	39.2 / 80.1	85.4 / 92.1	77.3
<b>iPET [9]</b>	80.6	92.9 / 92.4	95.0	74.0	52.2	80.1	33.0 / 74.0	86.0 / 86.5	76.8
<b>PET [9]</b>	79.4	85.1 / 59.4	95.0	69.8	52.4	80.1	37.9 / 77.3	86.0 / 86.5	74.1
<b>GPT-3 [2]</b>	77.5	82.1 / 57.2	92.0	72.9	55.3	75.0	32.5 / 74.8	89.0 / 90.1	73.2