



**ICML**  
International Conference  
On Machine Learning

# Black-Box Tuning for Language-Model-as-a-Service

**Tianxiang Sun**

School of Computer Science, Fudan University

<https://txsun1997.github.io/>

# Acknowledgment



Tianxiang Sun<sup>1</sup>



Yunfan Shao<sup>1</sup>



Hong Qian<sup>2</sup>



Xuanjing Huang<sup>1</sup>



Xipeng Qiu<sup>1,3</sup>



Yang Yu<sup>4</sup>



Zhengfu He<sup>1</sup>

<sup>1</sup> Fudan University

<sup>2</sup> East China Normal University

<sup>3</sup> Peng Cheng Laboratory

<sup>4</sup> Nanjing University

# Acknowledgment



Tianxiang Sun<sup>1</sup>



Yunfan Shao<sup>1</sup>



Hong Qian<sup>2</sup>



Xuanjing Huang<sup>1</sup>



Xipeng Qiu<sup>1,3</sup>

**NLP/Deep Learning**



Yang Yu<sup>4</sup>



Zhengfu He<sup>1</sup>

<sup>1</sup> Fudan University

<sup>2</sup> East China Normal University

<sup>3</sup> Peng Cheng Laboratory

<sup>4</sup> Nanjing University

# Acknowledgment



Tianxiang Sun<sup>1</sup>



Yunfan Shao<sup>1</sup>



Hong Qian<sup>2</sup>



Xuanjing Huang<sup>1</sup>



Xipeng Qiu<sup>1,3</sup>

**RL/Black-Box  
Optimization**



Yang Yu<sup>4</sup>



Zhengfu He<sup>1</sup>

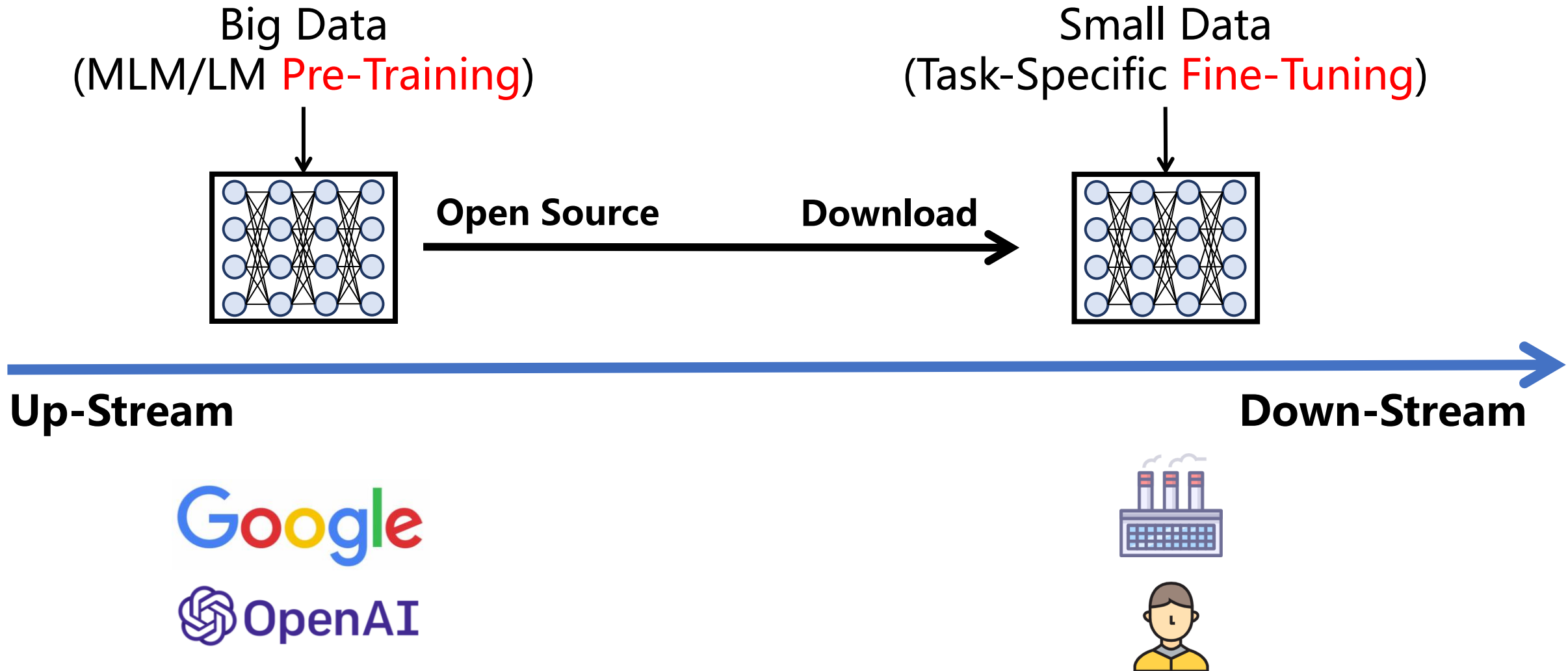
<sup>1</sup> Fudan University

<sup>2</sup> East China Normal University

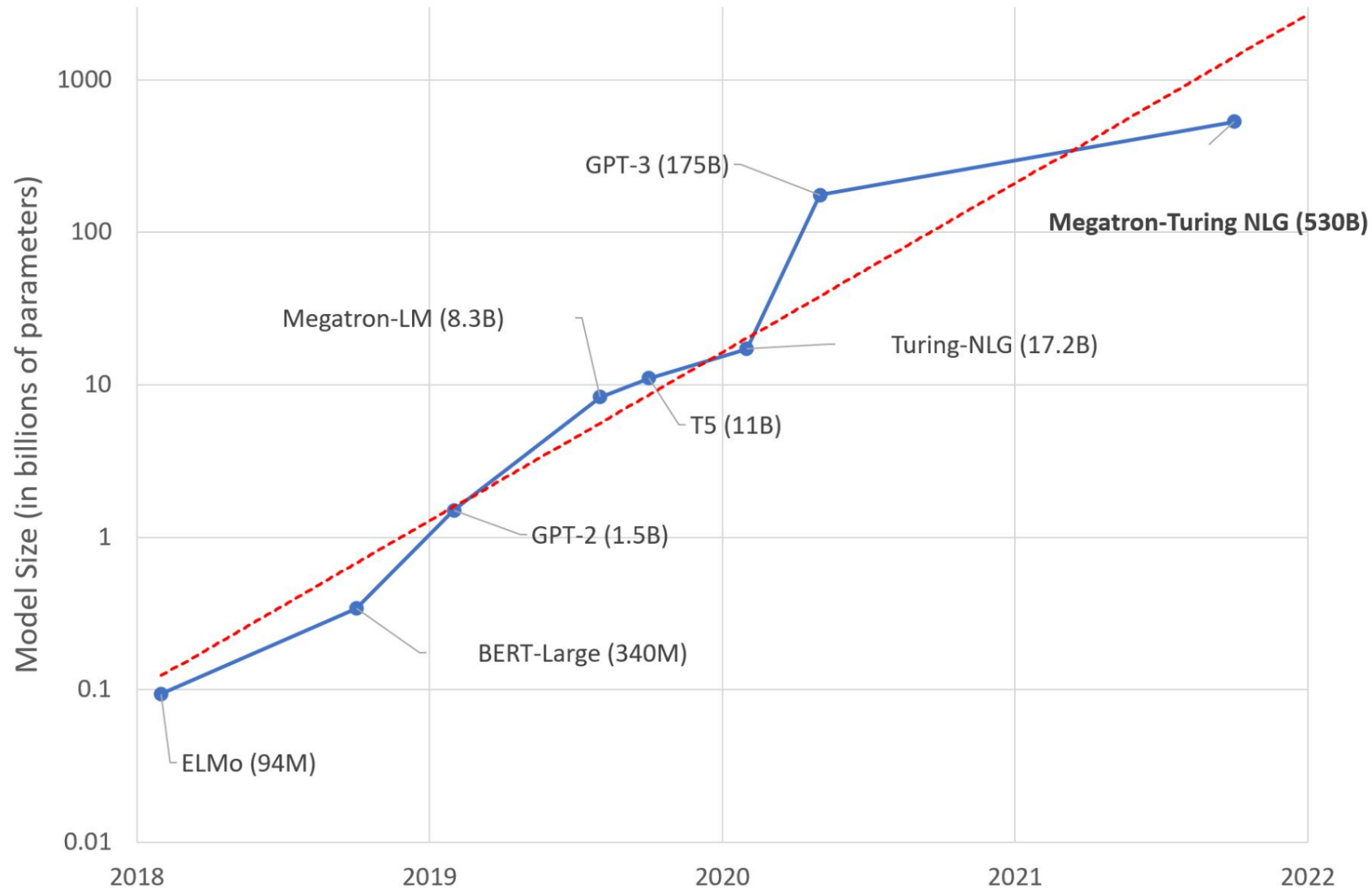
<sup>3</sup> Peng Cheng Laboratory

<sup>4</sup> Nanjing University

# Pre-train, then fine-tune



# When language models become larger...



# When language models become larger...

In the era of **large language models (LLMs)**...

- **Servers** often do not open-source the weights of LLMs due to commercial reasons
- **Users** usually do not have enough resources to run LLMs



# When language models become larger...

Why did OpenAI choose to release an API instead of open-sourcing the models?

There are three main reasons we did this. First, commercializing the technology helps us pay for our ongoing AI research, safety, and policy efforts.

Second, many of the models underlying the API are very large, taking a lot of expertise to develop and deploy and making them very expensive to run. This makes it hard for anyone except larger companies to benefit from the underlying technology. We're hopeful that the API will make powerful AI systems more accessible to smaller businesses and organizations.

Third, the API model allows us to more easily respond to misuse of the technology. Since it is hard to predict the downstream use cases of our models, it feels inherently safer to release them via an API and broaden access over time, rather than release an open source model where access cannot be adjusted if it turns out to have harmful applications.



# When language models become larger...

In the era of **large language models (LLMs)**...

- **Servers** often do not open-source the weights of LLMs due to commercial reasons
- **Users** usually do not have enough resources to run LLMs

The emergent ability of LLMs

- Manually craft text prompt to query LLMs
- In-context learning (GPT-3, [Brown et al., 2020](#))

# When language models become larger...

## Why in-context learning?

- Generalization: One general purpose model for all tasks
- Backpropagation is expensive
- Commercial use

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

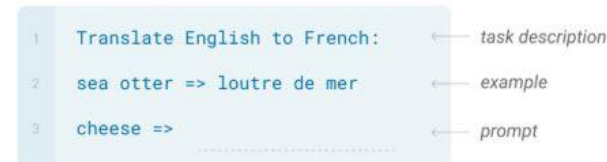
#### Zero-shot

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



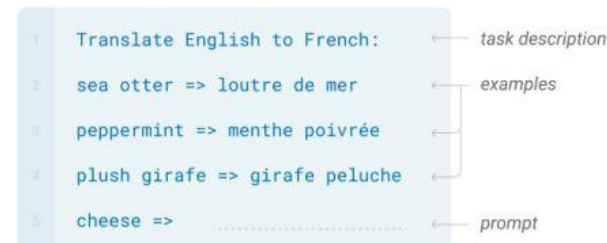
#### One-shot

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



#### Few-shot

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



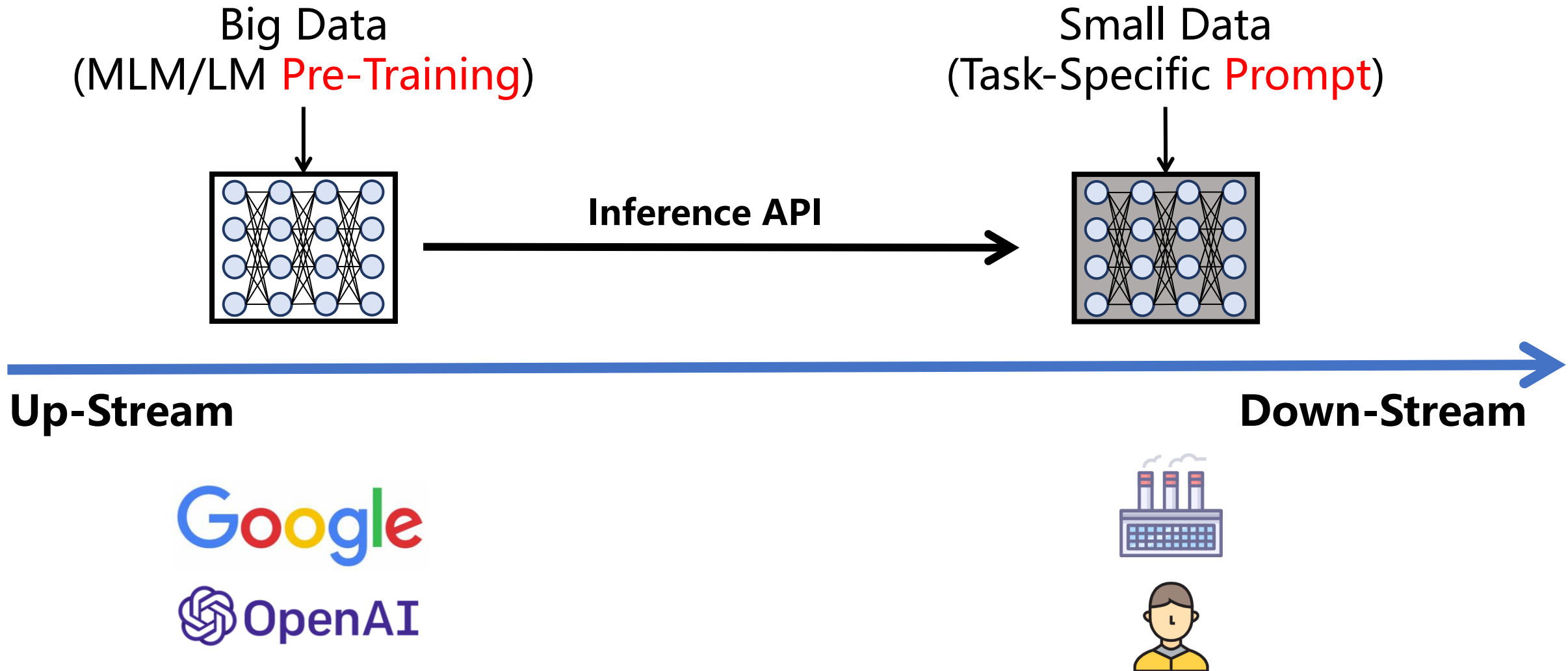
### Traditional fine-tuning (not used for GPT-3)

#### Fine-tuning

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



# Language-Model-as-a-Service (LMaaS)



# Language-Model-as-a-Service (LMaaS)

## GPT-3 Pricing

### Per-model prices

Ada Fastest

**\$0.0008** /1K tokens

Babbage

**\$0.0012** /1K tokens

Curie

**\$0.0060** /1K tokens

Davinci Most powerful

**\$0.0600** /1K tokens

# Language-Model-as-a-Service (LMaaS)

### Describe a layout.

Just describe any layout you want, and it'll try to render below!

Generate

### Products

Select product ▼

### Collections


- New
- Popular
- Upcoming
- Requested

### Categories


|                       |     |
|-----------------------|-----|
| All                   | 319 |
| A/B Testing           | 2   |
| Ad Generation         | 3   |
| AI Copywriting        | 37  |
| AI Writing Assistants | 1   |
| API Design            | 1   |
| Avatars               | 1   |
| Blog writing          | 2   |
| Book Writing          | 1   |

### New


Recently added GPT-3 apps




Customer Service  
ActiveChat.ai




Chatbots  
AI Buddy




Humor  
AI Guru




LegalTech  
aiLawDocs



Chatbots  
AskBrian



Developer Tools  
Azure OpenAI Service



Generative Art  
Botto







Image captioning  
ClipClap




Healthcare  
Curai




Code Generation  
DeepGenX




Summarization  
Delv AI




API Design  
Design an API with ...



Recruiting  
Drafted



Language Learning  
Duolingo



Research Assistants  
Elicit

# LMaaS in China

## 赛题说明

INTRODUCE

本次大赛的主题为基于悟道2.0大模型的创新应用开发，面向在校大学生、企业工程师、科研工作者等全球开发者全面开放。参赛选手需要依据对悟道能力的理解，结合社会关注的热点选择健康医疗、教育学习、社交生活、效率工具、环境自然或其他具有社会价值、产业价值的相关主题，提交一个潜在的智能应用方案并上线应用。

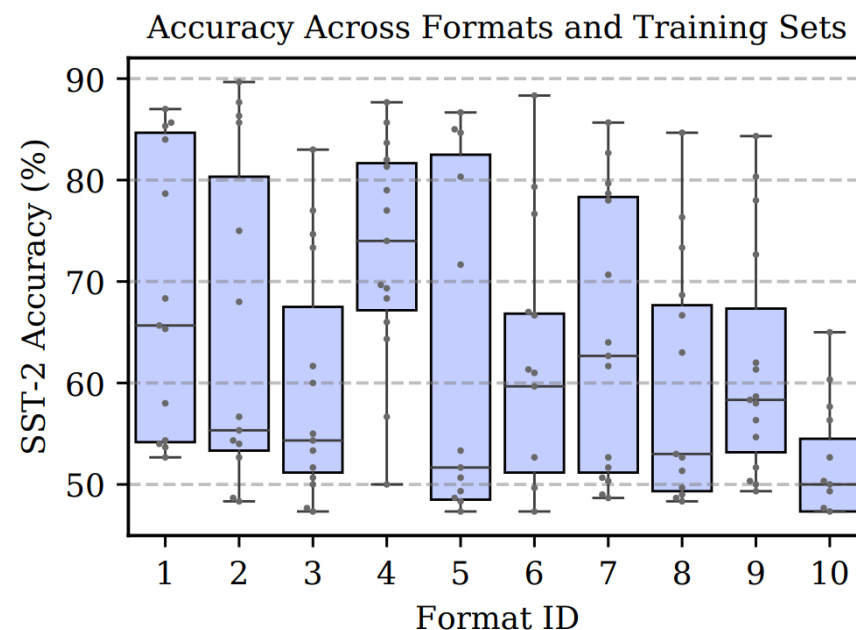
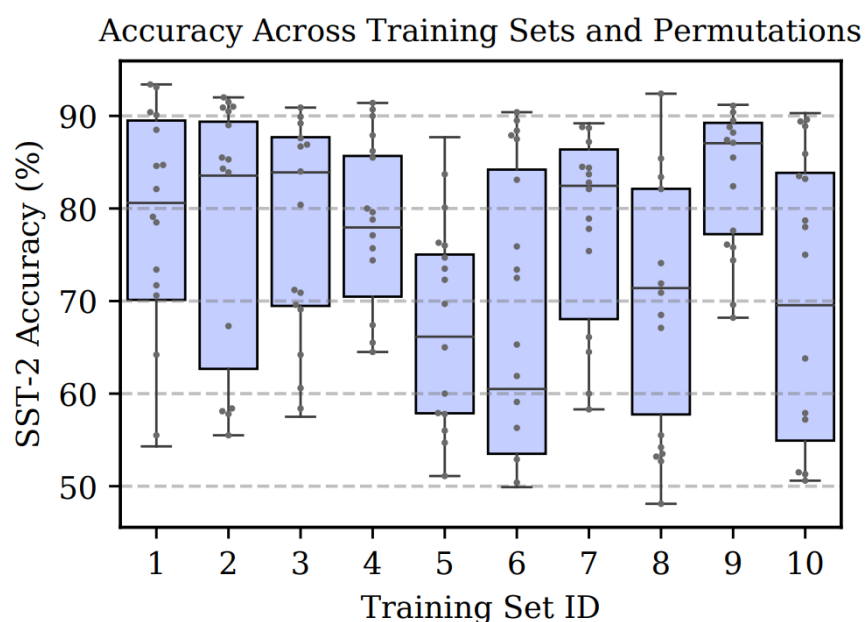
## API 文档说明

本次比赛共提供 9 个应用的 API 接口，每支队伍每天访问次数有限。

- ✓ CogView API 文档，100 次 / 天。
- ✓ 宋词 API 文档，100 次 / 天。
- ✓ 藏头诗 API 文档，100 次 / 天。
- ✓ 问答 API 文档，100 次 / 天。
- ✓ 获取图像的特征向量，1000 次 / 天。
- ✓ 写诗 API 文档，100 次 / 天。
- ✓ 获取文本的特征向量，1000 次 / 天。
- ✓ 新闻 API 文档，100 次 / 天。
- ✓ 快速写诗 API 文档，100 次 / 天。

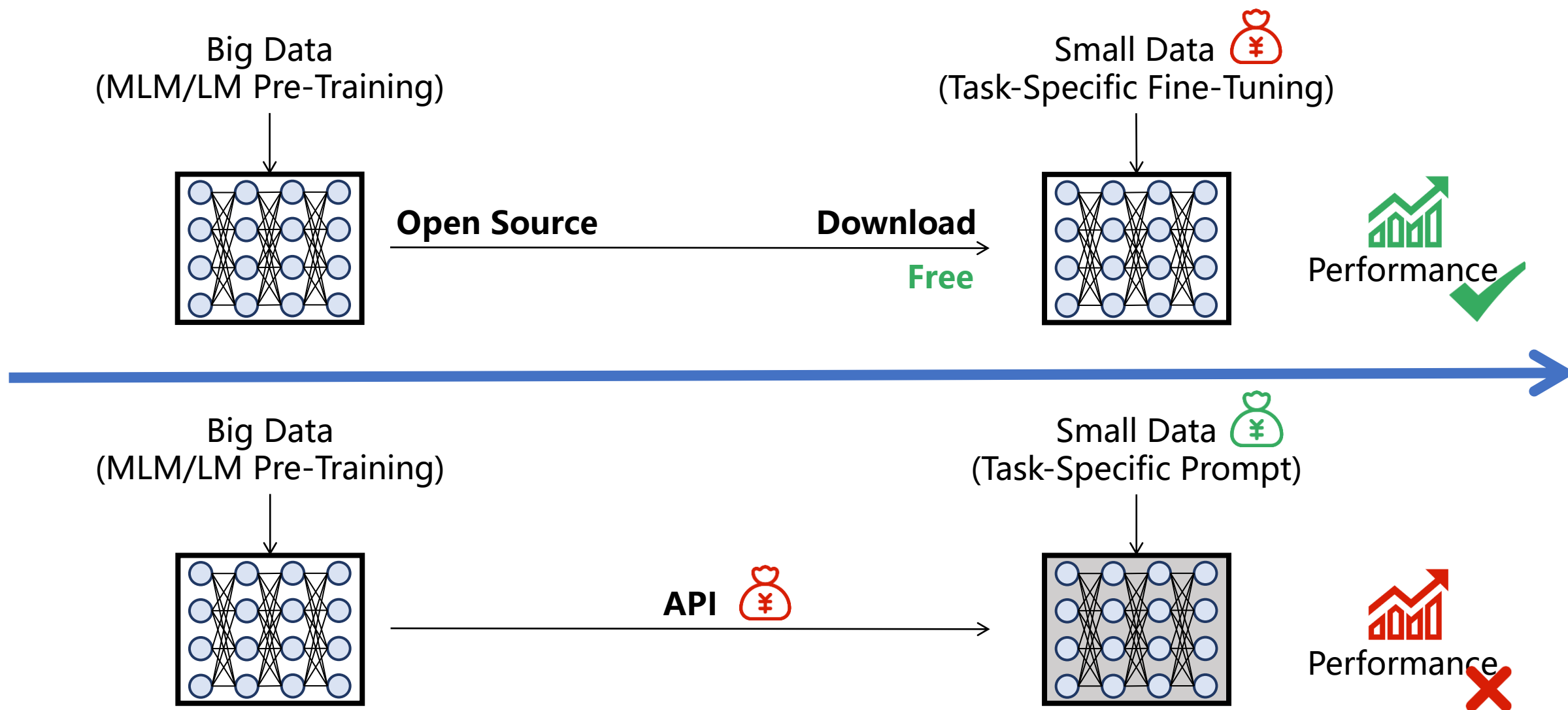
# However...

The performance of **manual prompt** and **in-context learning** highly depend on the choice of prompt and demonstrations, and lags far behind model tuning.





# Grounding LLMs From the Cloud



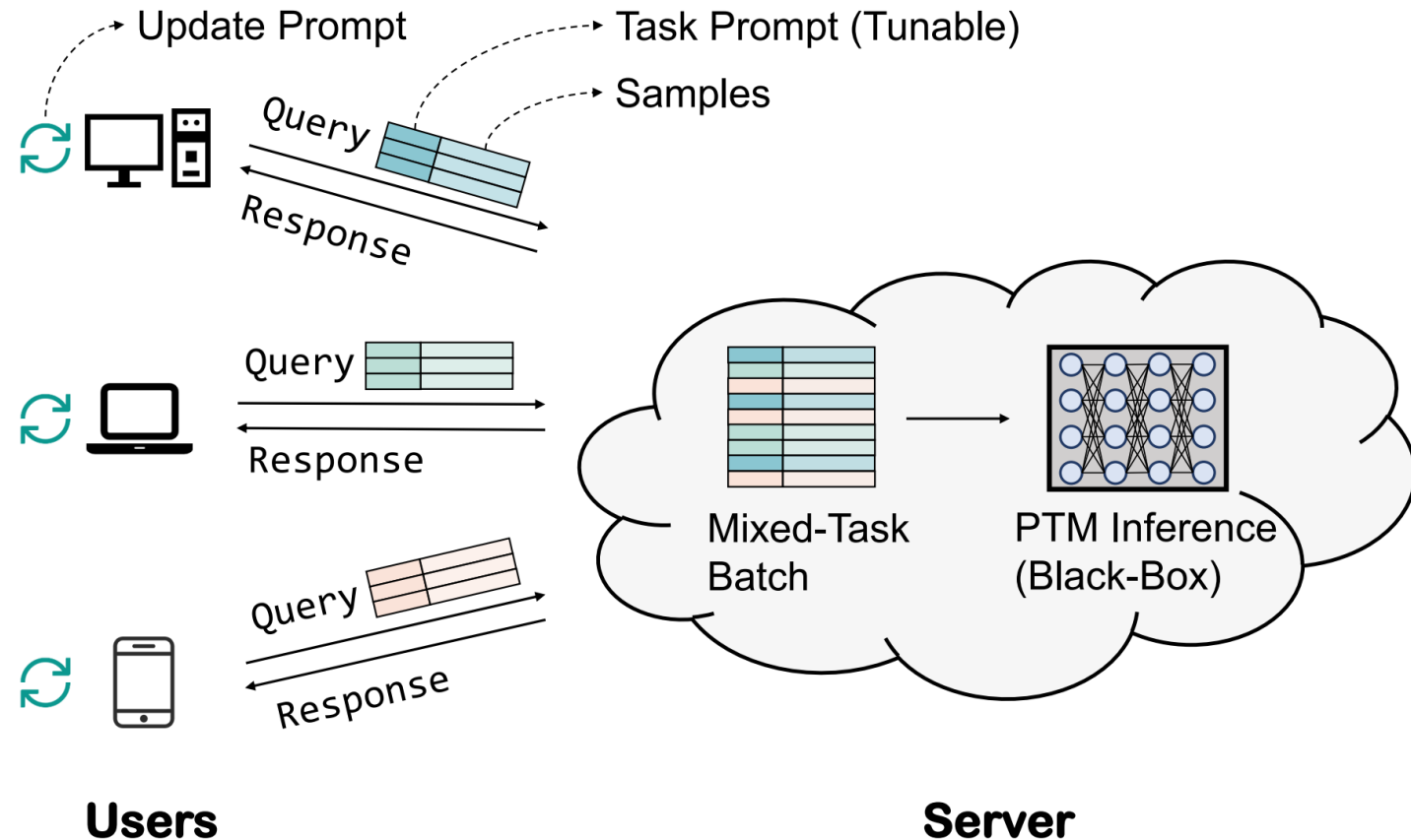
**Performance is the Key for grounding. (Who are the users?)**

# To make LLMs benefit more people...

Can we optimize the prompt with the API feedback? (without expensive backpropagation)

Objective:

$$\mathbf{p}^* = \arg \min_{\mathbf{p} \in \mathcal{P}} \mathcal{L}(f(\mathbf{p}; \tilde{X}), \tilde{Y})$$



# A challenge of high dimensionality

Considering optimization of the continuous prompt, the dimensionality can be **tens of thousands** (say we are going to optimize 50 prompt tokens, each with 1k dimensions, there are 50k parameters to be optimized.)

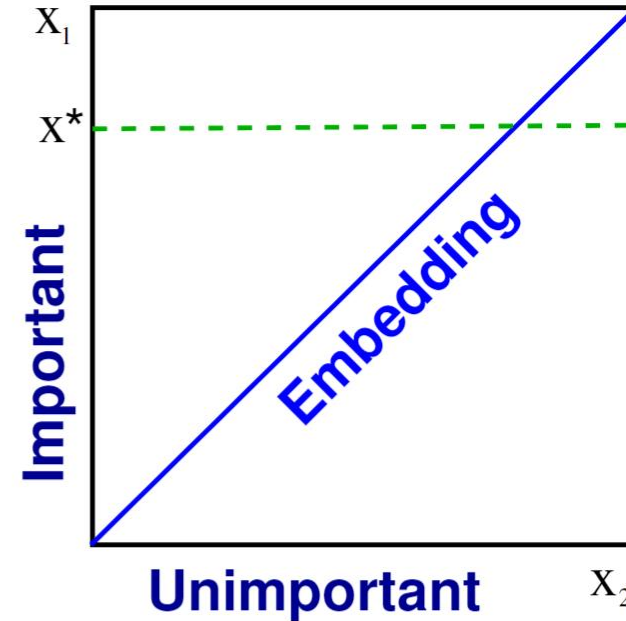
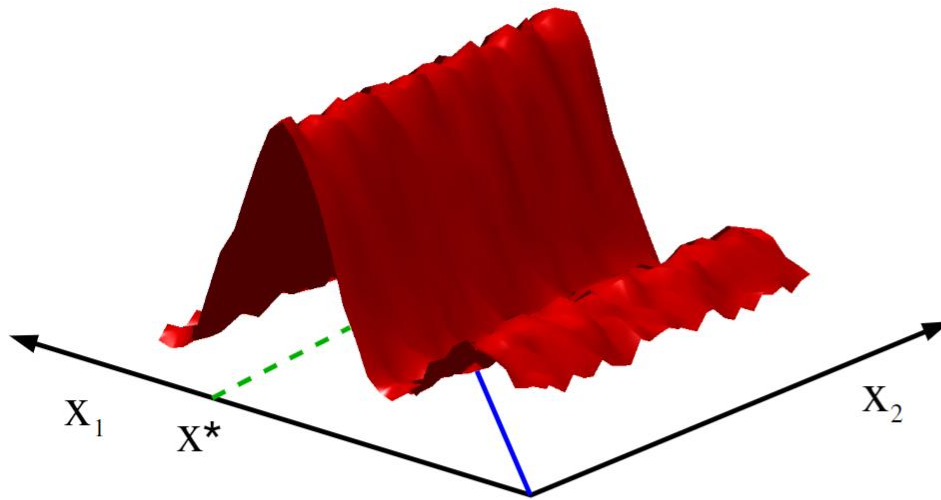
**Derivative-free optimization (DFO)** can struggle with high-dimensional problems, **except for** the case when the problem has a low **intrinsic dimensionality**.

Note: Intrinsic dimensionality is the minimal number of parameters needed to represent the problem

# A challenge of high dimensionality

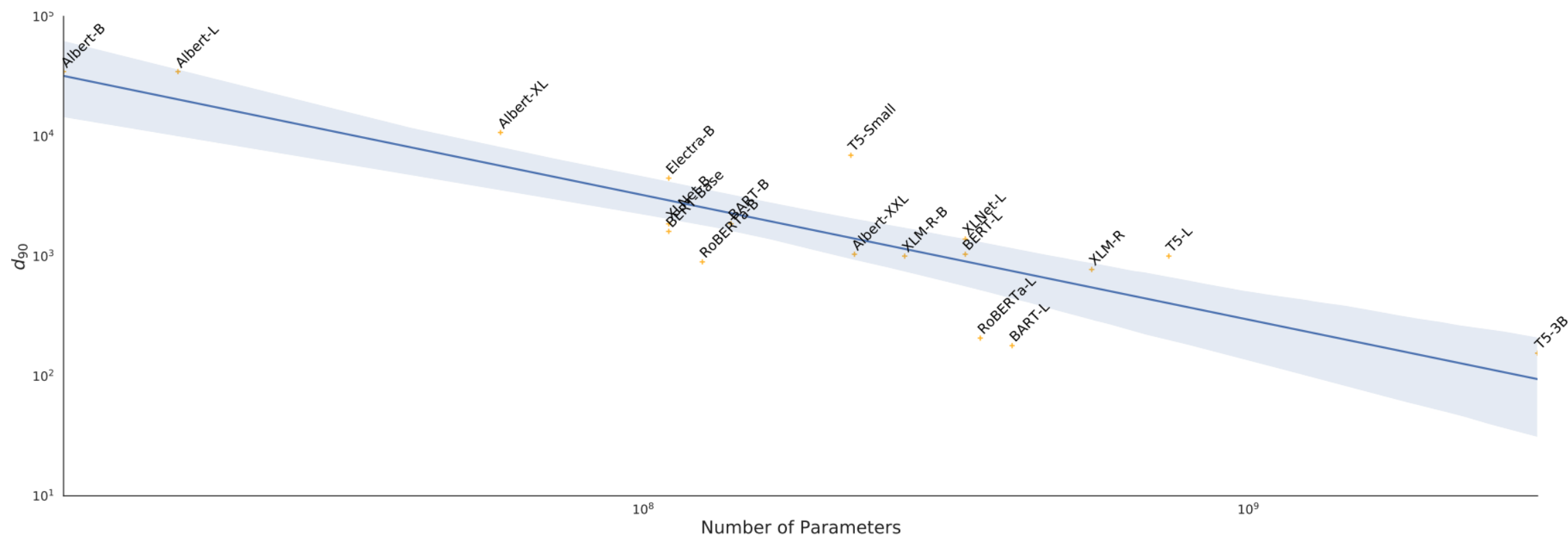
An example

- The objective to be optimized has two dimensions but only one matters
- In that case we can perform optimization with **random embedding**

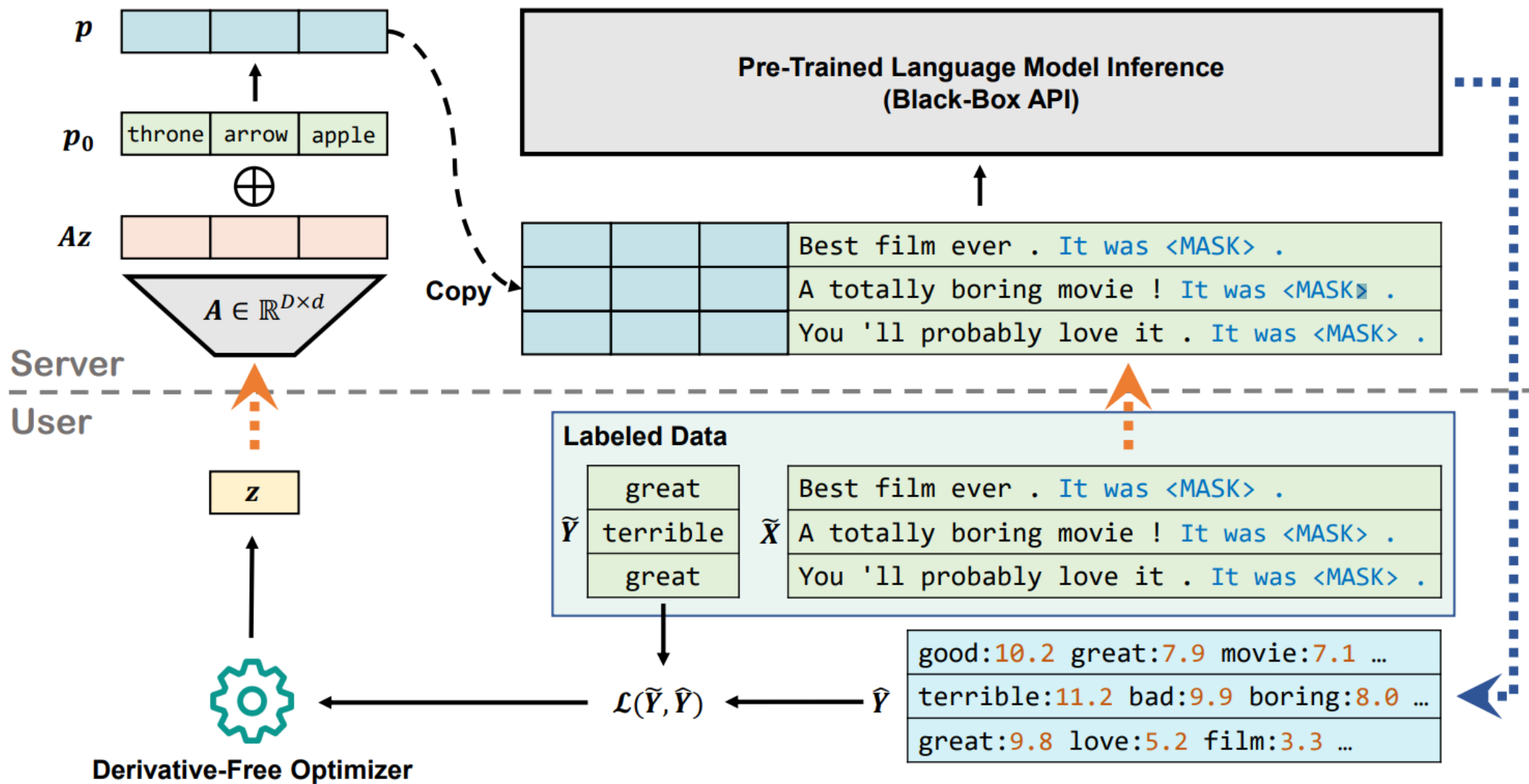


# Fortunately...

LLMs have a very low intrinsic dimensionality!



# Black-Box Tuning



# Black-Box Tuning

## The CMA-ES (Covariance Matrix Adaptation Evolution Strategy)

---

### The CMA-ES (Evolution Strategy with Covariance Matrix Adaptation)

---

Consider  $P^{(t)} = \mathcal{N}(\boldsymbol{\mu}^{(t)}, \sigma^{(t)^2} \mathbf{C}^{(t)})$  where  $\boldsymbol{\mu}^{(t)} \in \mathbb{R}^n$ ,  $\sigma^{(t)} \in \mathbb{R}_+$ ,  $\mathbf{C}^{(t)} \in \mathbb{R}^{n \times n}$

- $\boldsymbol{\mu}^{(t)} \rightarrow \boldsymbol{\mu}^{(t+1)}$ : Maximum likelihood update, i.e.  $P(\mathbf{x}_{\text{selected}}^{(t)} | \boldsymbol{\mu}^{(t+1)}) \rightarrow \max$
- $\mathbf{C}^{(t)} \rightarrow \mathbf{C}^{(t+1)}$ : Maximum likelihood update, i.e.  $P(\frac{\mathbf{x}_{\text{selected}}^{(t)} - \boldsymbol{\mu}^{(t)}}{\sigma^{(t)}} | \mathbf{C}^{(t+1)}) \rightarrow \max$ , under consideration of prior  $\mathbf{C}^{(t)}$  (otherwise  $\mathbf{C}^{(t+1)}$  becomes singular).
- $\sigma^{(t)} \rightarrow \sigma^{(t+1)}$ : Update to achieve conjugate perpendicularity, i.e. conceptually  $(\boldsymbol{\mu}^{(t+2)} - \boldsymbol{\mu}^{(t+1)})^T \mathbf{C}^{(t)-1} (\boldsymbol{\mu}^{(t+1)} - \boldsymbol{\mu}^{(t)}) / \sigma^{(t+1)^2} \rightarrow 0$



# Experiments

## Datasets and processing details

| Category        | Dataset   | $ \mathcal{V} $ | $ \text{Train} $ | $ \text{Test} $ | Type       | Template  | Label words  |
|-----------------|-----------|-----------------|------------------|-----------------|------------|---|--|
| single-sentence | SST-2     | 2               | 67k              | 0.9k            | sentiment  | $\langle S \rangle$ . It was [MASK].                  | great, bad   |
|                 | Yelp P.   | 2               | 560k             | 38k             | sentiment  | $\langle S \rangle$ . It was [MASK].                  | great, bad   |
|                 | AG's News | 4               | 120k             | 7.6k            | topic      | [MASK] News: $\langle S \rangle$                      | World, Sports, Business, Tech  |
|                 | DBPedia   | 14              | 560k             | 70k             | topic      | [Category: [MASK]] $\langle S \rangle$                | Company, Education, Artist, Athlete, Office, Transportation, Building, Natural, Village, Animal, Plant, Album, Film, Written |
| sentence-pair   | MRPC      | 2               | 3.7k             | 0.4k            | paraphrase | $\langle S_1 \rangle$ ? [MASK], $\langle S_2 \rangle$ | Yes, No  |
|                 | RTE       | 2               | 2.5k             | 0.3k            | NLI        | $\langle S_1 \rangle$ ? [MASK], $\langle S_2 \rangle$ | Yes, No  |
|                 | SNLI      | 3               | 549k             | 9.8k            | NLI        | $\langle S_1 \rangle$ ? [MASK], $\langle S_2 \rangle$ | Yes, Maybe, No   |

# Experiments

16-shot (per class) learning with RoBERTa-large (350M)

| Method                        | SST-2<br>acc     | Yelp P.<br>acc   | AG's News<br>acc  | DBPedia<br>acc   | MRPC<br>F1       | SNLI<br>acc      | RTE<br>acc       | Avg.         |
|-------------------------------|------------------|------------------|-------------------|------------------|------------------|------------------|------------------|--------------|
| <i>Gradient-Based Methods</i> |                  |                  |                   |                  |                  |                  |                  |              |
| Prompt Tuning                 | 68.23 $\pm$ 3.78 | 61.02 $\pm$ 6.65 | 84.81 $\pm$ 0.66  | 87.75 $\pm$ 1.48 | 51.61 $\pm$ 8.67 | 36.13 $\pm$ 1.51 | 54.69 $\pm$ 3.79 | 63.46        |
| + Pre-trained prompt          | /                | /                | /                 | /                | 77.48 $\pm$ 4.85 | 64.55 $\pm$ 2.43 | 77.13 $\pm$ 0.83 | 74.42        |
| P-Tuning v2                   | 64.33 $\pm$ 3.05 | 92.63 $\pm$ 1.39 | 83.46 $\pm$ 1.01  | 97.05 $\pm$ 0.41 | 68.14 $\pm$ 3.89 | 36.89 $\pm$ 0.79 | 50.78 $\pm$ 2.28 | 70.47        |
| Model Tuning                  | 85.39 $\pm$ 2.84 | 91.82 $\pm$ 0.79 | 86.36 $\pm$ 1.85  | 97.98 $\pm$ 0.14 | 77.35 $\pm$ 5.70 | 54.64 $\pm$ 5.29 | 58.60 $\pm$ 6.21 | 78.88        |
| <i>Gradient-Free Methods</i>  |                  |                  |                   |                  |                  |                  |                  |              |
| Manual Prompt                 | 79.82            | 89.65            | 76.96             | 41.33            | 67.40            | 31.11            | 51.62            | 62.56        |
| In-Context Learning           | 79.79 $\pm$ 3.06 | 85.38 $\pm$ 3.92 | 62.21 $\pm$ 13.46 | 34.83 $\pm$ 7.59 | 45.81 $\pm$ 6.67 | 47.11 $\pm$ 0.63 | 60.36 $\pm$ 1.56 | 59.36        |
| Feature-MLP                   | 64.80 $\pm$ 1.78 | 79.20 $\pm$ 2.26 | 70.77 $\pm$ 0.67  | 87.78 $\pm$ 0.61 | 68.40 $\pm$ 0.86 | 42.01 $\pm$ 0.33 | 53.43 $\pm$ 1.57 | 66.63        |
| Feature-BiLSTM                | 65.95 $\pm$ 0.99 | 74.68 $\pm$ 0.10 | 77.28 $\pm$ 2.83  | 90.37 $\pm$ 3.10 | 71.55 $\pm$ 7.10 | 46.02 $\pm$ 0.38 | 52.17 $\pm$ 0.25 | 68.29        |
| <b>Black-Box Tuning</b>       | 89.56 $\pm$ 0.25 | 91.50 $\pm$ 0.16 | 81.51 $\pm$ 0.79  | 87.80 $\pm$ 1.53 | 61.56 $\pm$ 4.34 | 46.58 $\pm$ 1.33 | 52.59 $\pm$ 2.21 | 73.01        |
| + Pre-trained prompt          | /                | /                | /                 | /                | 75.51 $\pm$ 5.54 | 83.83 $\pm$ 0.21 | 77.62 $\pm$ 1.30 | <b>83.90</b> |

# Experiments

## Detailed comparison on SST-2 and AG News

|                                      | Deployment-<br>Efficient | As-A-<br>Service | Test<br>Accuracy | Training<br>Time               | Memory<br>User | Footprint<br>Server | Upload<br>per query | Download<br>per query |
|--------------------------------------|--------------------------|------------------|------------------|--------------------------------|----------------|---------------------|---------------------|-----------------------|
| SST-2 (max sequence length: 47)      |                          |                  |                  |                                |                |                     |                     |                       |
| Prompt Tuning                        | ✓                        | ×                | 72.6             | 15.9 mins                      | -              | 5.3 GB              | -                   | -                     |
| Model Tuning                         | ×                        | ×                | 87.8             | 9.8 mins                       | -              | 7.3 GB              | -                   | -                     |
| Feature-MLP                          | ✓                        | ✓                | 63.8             | 7.0 mins                       | 20 MB          | 2.8 GB              | 4 KB                | 128 KB                |
| Feature-BiLSTM                       | ✓                        | ✓                | 66.2             | 9.3 mins                       | 410 MB         | 2.8 GB              | 4 KB                | 6016 KB               |
| Black-Box Tuning                     | ✓                        | ✓                | 89.4             | 10.1 (6.1 <sup>*</sup> ) mins  | 30 MB          | 3.0 GB              | 6 KB                | 0.25 KB               |
| AG's News (max sequence length: 107) |                          |                  |                  |                                |                |                     |                     |                       |
| Prompt Tuning                        | ✓                        | ×                | 84.0             | 30.2 mins                      | -              | 7.7 GB              | -                   | -                     |
| Model Tuning                         | ×                        | ×                | 88.4             | 13.1 mins                      | -              | 7.3 GB              | -                   | -                     |
| Feature-MLP                          | ✓                        | ✓                | 71.0             | 13.5 mins                      | 20 MB          | 3.6 GB              | 20 KB               | 256 KB                |
| Feature-BiLSTM                       | ✓                        | ✓                | 73.1             | 19.7 mins                      | 500 MB         | 3.6 GB              | 20 KB               | 27392 KB              |
| Black-Box Tuning                     | ✓                        | ✓                | 82.6             | 21.0 (17.7 <sup>*</sup> ) mins | 30 MB          | 4.6 GB              | 22 KB               | 1 KB                  |

# Forward Is All You Need?

## Limitations of black-box tuning:

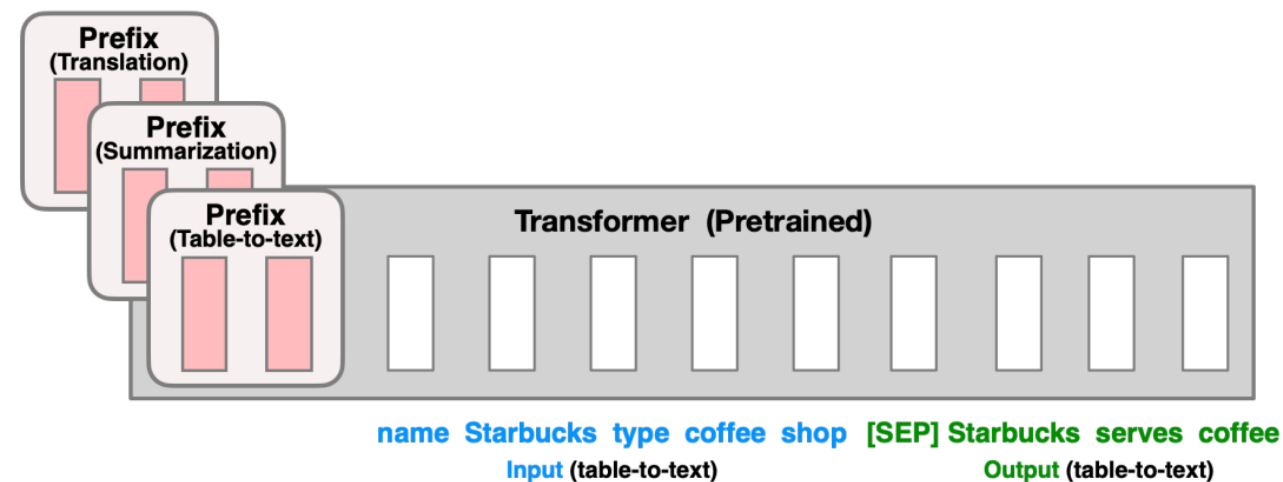
- Slow convergence on many-label classification (e.g., DBPedia)
- Requirement of prompt pre-training (gradient) on difficult tasks (e.g., SNLI)

## Current version of black-box tuning is just a lower bound:

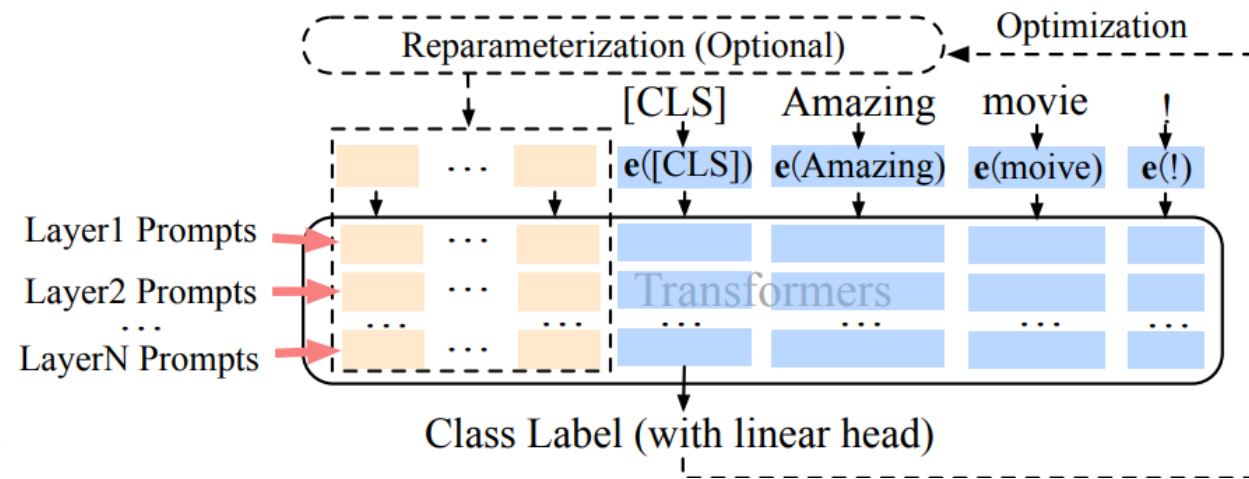
- Prompt/verbalizer engineering, prompt ensemble, prompt pre-training...
- Better derivative-free algorithms
- Pre-trained random embedding
- ...

# Can We Go Deeper?

## The ``Deep Prompt Tuning``



Prefix Tuning (Li and Liang, ACL 2021)



P-Tuning v2 (Liu et al., ACL 2022)

# Can We Go Deeper?

The challenge, again, is the **high dimensionality**

- Say we are going to optimize **50 prompt tokens** at each layer of RoBERTa-large, each with **1k dimensions**, there are  $50k \times 24 = 1.2M$  parameters to be optimized
- Besides, the prompt parameters at different layers are **heterogenous** and therefore we can not simply use the random embedding to solve it

# Take A Closer Look Into the Forward Pass

Thanks to the residual connections in modern LLMs, the forward computation can be decomposed as an **additive form**

An example of a 3-layer model:

$$\begin{aligned} f(\mathbf{x}_1) &= f_3(\mathbf{x}_3) + \mathbf{x}_3 \\ &= f_3(\mathbf{x}_3) + f_2(\mathbf{x}_2) + \mathbf{x}_2 \\ &= f_3(\mathbf{x}_3) + f_2(\mathbf{x}_2) + f_1(\mathbf{x}_1) + \mathbf{x}_1 \end{aligned}$$

Therefore, the optimization can be decomposed into multiple sub-problems!



# Take A Closer Look Into the Forward Pass

A general formulation of ``deep black-box tuning``:

$$f(\mathbf{x}_1, \mathbf{p}) = [\mathbf{A}_1 \mathbf{z}_1 + \mathbf{p}_1^0; \mathbf{x}_1] \\ + \sum_{j=1}^L f_j([\mathbf{A}_j \mathbf{z}_j + \mathbf{p}_j^0; \mathbf{x}_j])$$

Given such an additive form, we propose a **divide-and-conquer (DC)** algorithm to alternately optimize prompt at each layer

# Divide-and-Conquer

- Layer-specific optimizer
- Layer-specific random projection
- Alternate from the bottom to top

---

**Algorithm 1:** DC Algorithm for BBTv2

---

**Require:**  $L$ -layer PTM Inference API  $f$ ,  
Loss function  $\mathcal{L}$ ,  
Budget of API calls  $\mathcal{B}$ ,  
Derivative-free optimizers  $\{\mathcal{M}_j\}_{j=1}^L$

- 1: Initialize random projections  $\mathbf{A}_1, \dots, \mathbf{A}_L$
- 2: Initialize parameters  $\mathbf{z}_1^{(0)}, \dots, \mathbf{z}_L^{(0)}$
- 3: Deep prompts  $\mathbf{p} = \langle \mathbf{A}_1 \mathbf{z}_1^{(0)}, \dots, \mathbf{A}_L \mathbf{z}_L^{(0)} \rangle$
- 4: **for**  $i = 1$  to  $\mathcal{B}/L$  **do**
- 5:   **for**  $j = 1$  to  $L$  **do**
- 6:     Evaluate:  $loss = \mathcal{L}(f(\mathbf{p}))$
- 7:     Update:  $\mathbf{z}_j^{(i)} \leftarrow \mathcal{M}_j(\mathbf{z}_j^{(i-1)}, loss)$
- 8:     Replace:  $\mathbf{p}_j \leftarrow \mathbf{A}_j \mathbf{z}_j^{(i)}$
- 9:   **end for**
- 10: **end for**
- 11: **return** Optimized deep prompts  $\mathbf{p}$

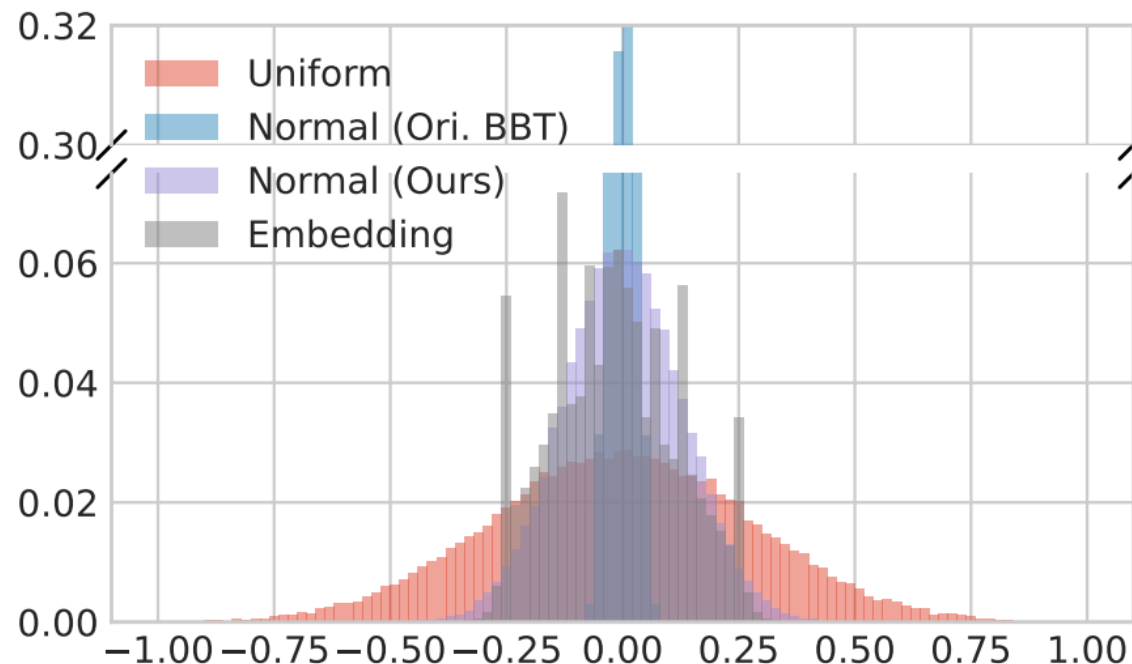
---

# Revisiting Random Projection (Embedding)

Generating random projections from a normal distribution with std dev as

$$\sigma_A = \frac{\hat{\sigma}}{\sqrt{d}\sigma_z}$$

A visualization of generated prompt with RoBERTa-large



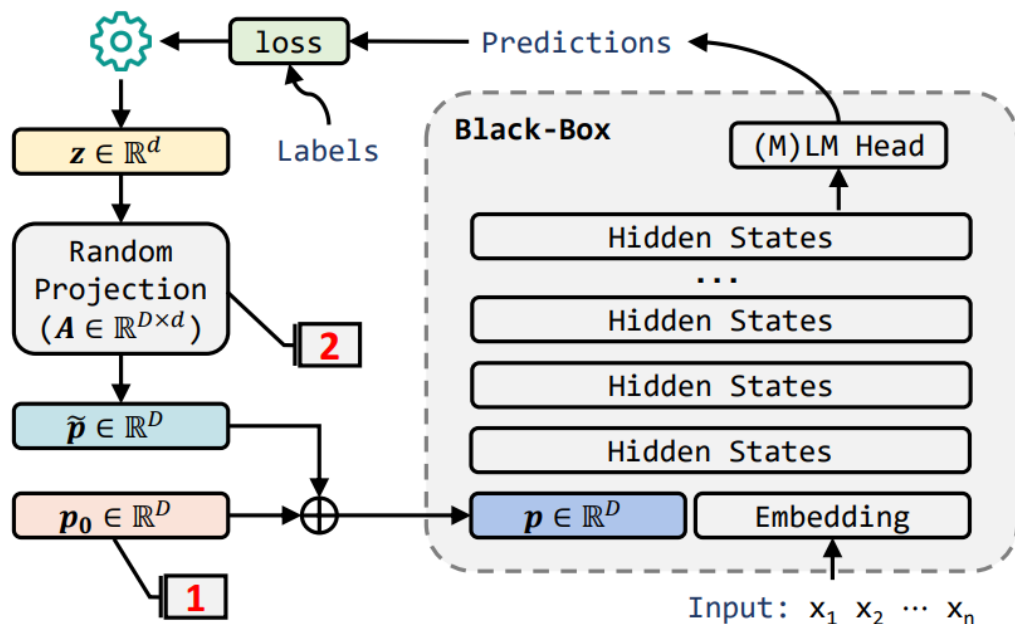
# BBTv2: Towards A Gradient-Free Future

Main improvements of BBTv2

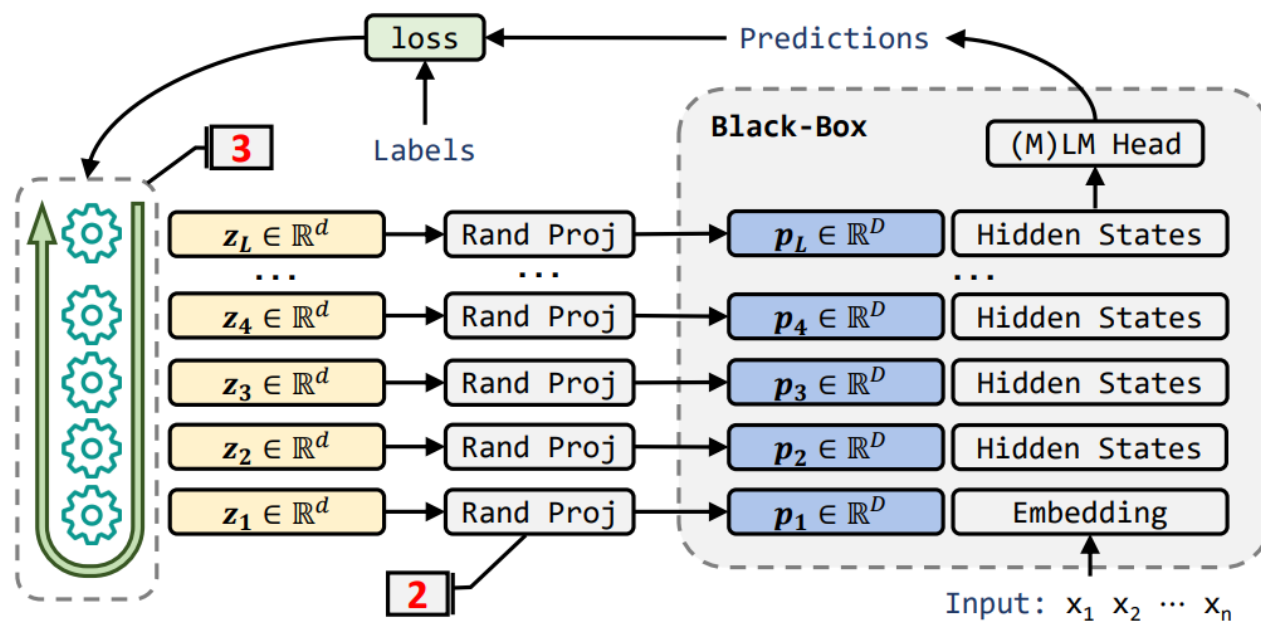
- Get rid of prompt pre-training
- Improved random projection
- Deep prompts

# BBTv2: Towards A Gradient-Free Future

## Main improvements of BBTv2



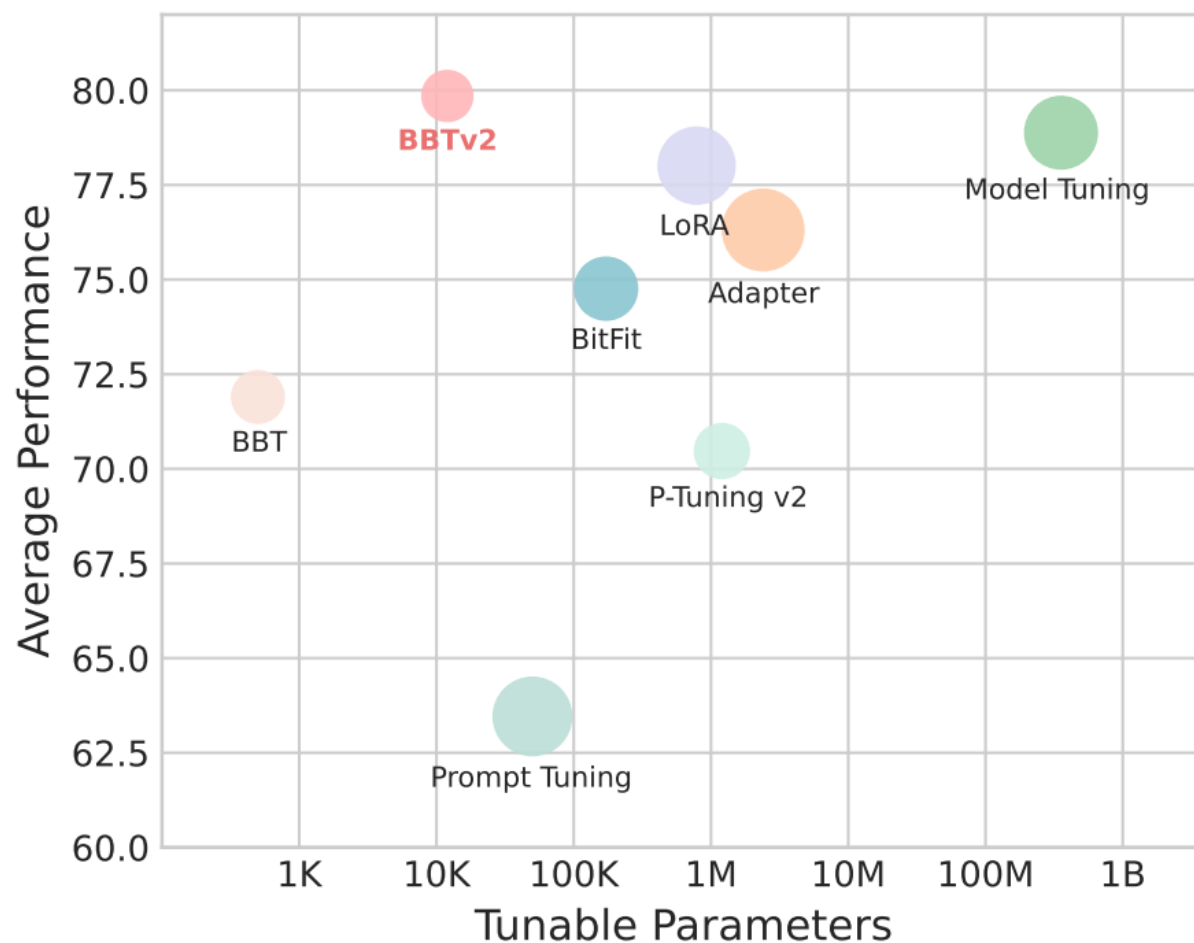
(a) BBT



(b) BBTv2

# Experiments of BBTv2

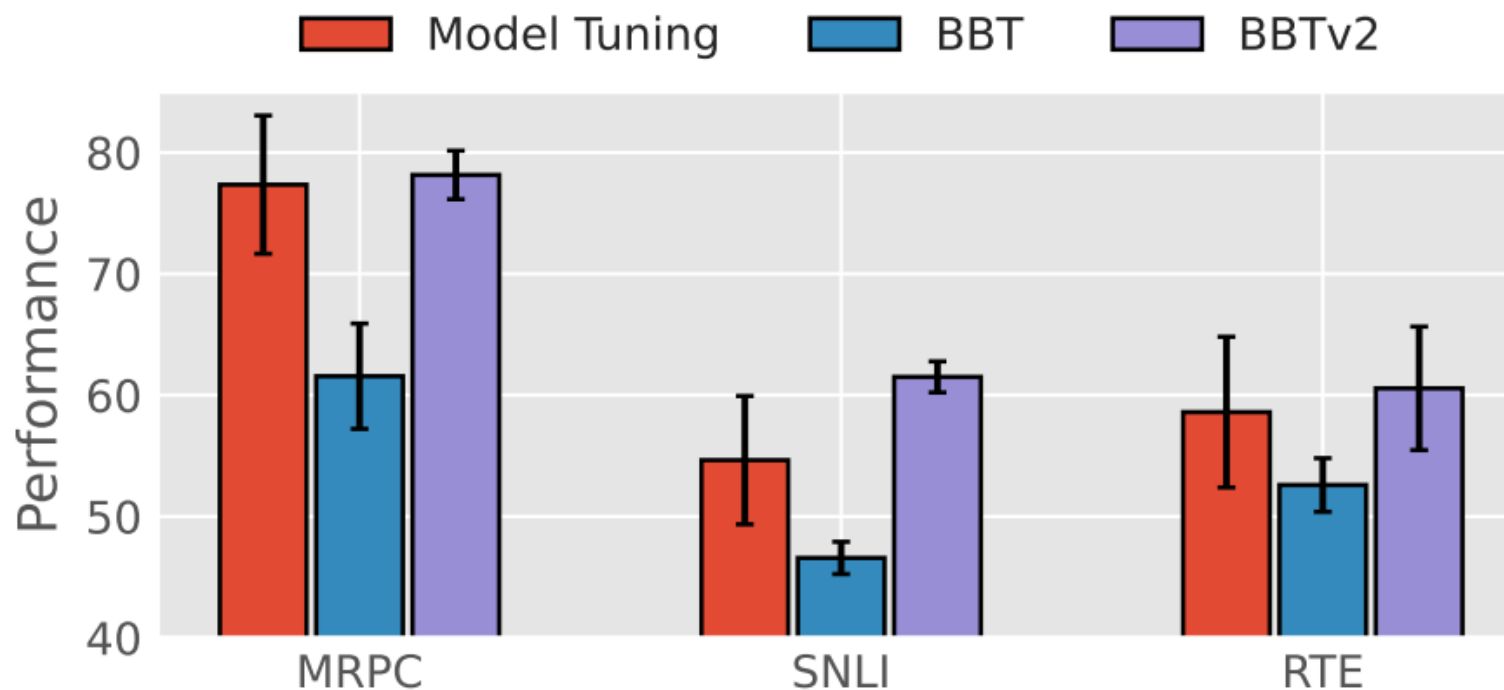
Comparable to full model tuning but merely tuning ~10k parameters



# Experiments of BBTv2

Improve BBT on entailment tasks

- Be comparable to full model tuning without pre-trained prompt embedding

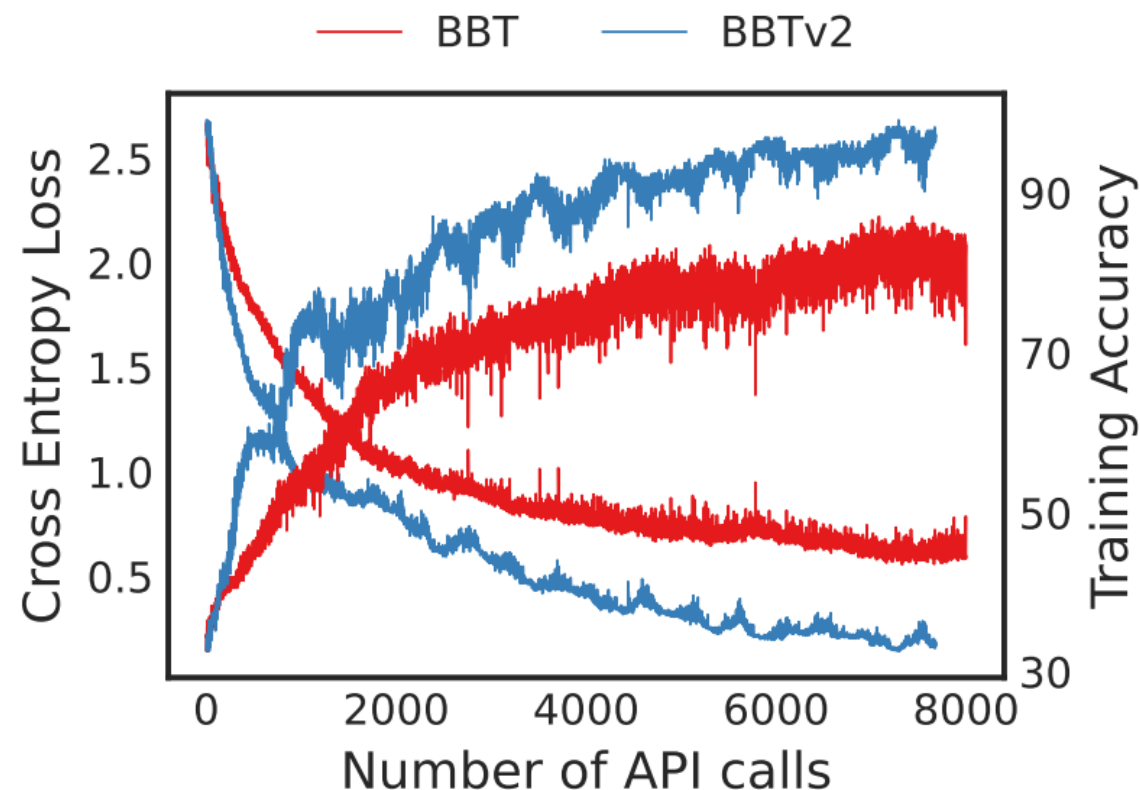




# Experiments of BBTv2

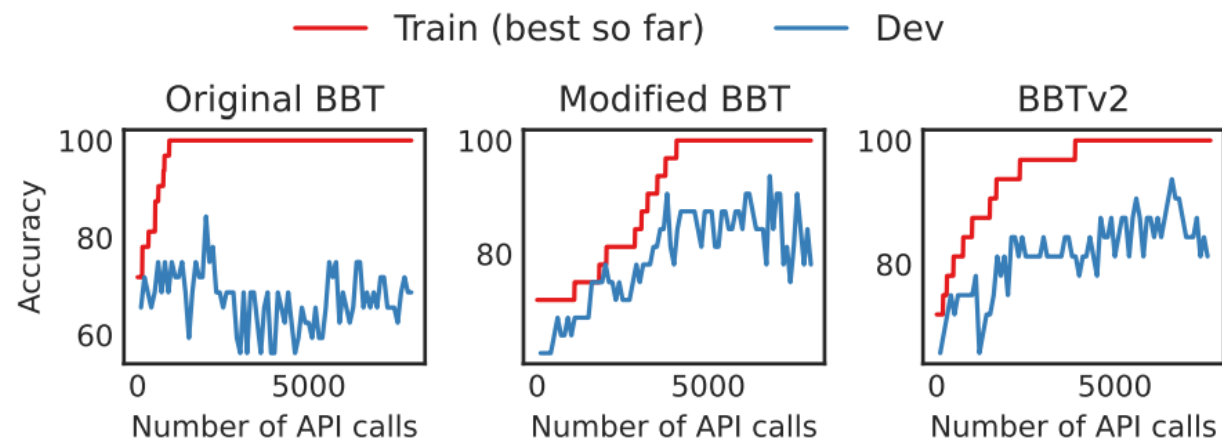
Improve BBT on many-label classification tasks

- Faster convergence than BBT on DBPedia (14 classes)

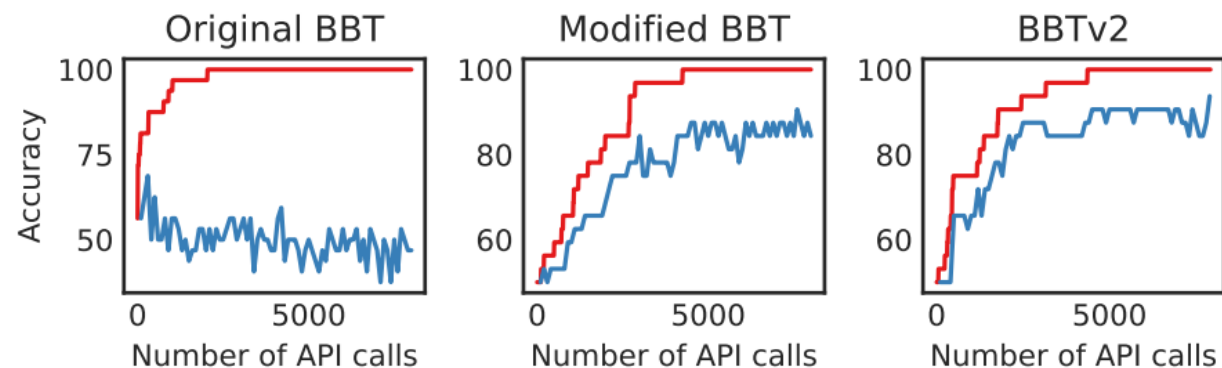


# Experiments of BBTv2

## Generalization across LMs



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



# Experiments of BBTv2

## Overall comparison

| Method                        | Tunable Params | SST-2<br>acc            | Yelp P.<br>acc          | AG's News<br>acc        | DBPedia<br>acc          | MRPC<br>F1              | SNLI<br>acc             | RTE<br>acc              | Avg.         |
|-------------------------------|----------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|--------------|
| <i>Gradient-Based Methods</i> |                |                         |                         |                         |                         |                         |                         |                         |              |
| Model Tuning                  | 355M           | <u>85.39</u> $\pm 2.84$ | <u>91.82</u> $\pm 0.79$ | 86.36 $\pm 1.85$        | <u>97.98</u> $\pm 0.14$ | <b>77.35</b> $\pm 5.70$ | 54.64 $\pm 5.29$        | <b>58.60</b> $\pm 6.21$ | <b>78.88</b> |
| Adapter                       | 2.4M           | 83.91 $\pm 2.90$        | 90.99 $\pm 2.86$        | 86.01 $\pm 2.18$        | <b>97.99</b> $\pm 0.07$ | 69.20 $\pm 3.58$        | <u>57.46</u> $\pm 6.63$ | 48.62 $\pm 4.74$        | 76.31        |
| BitFit                        | 172K           | 81.19 $\pm 6.08$        | 88.63 $\pm 6.69$        | <u>86.83</u> $\pm 0.62$ | 94.42 $\pm 0.94$        | 66.26 $\pm 6.81$        | 53.42 $\pm 10.63$       | 52.59 $\pm 5.31$        | 74.76        |
| LoRA                          | 786K           | <b>88.49</b> $\pm 2.90$ | 90.21 $\pm 4.00$        | <b>87.09</b> $\pm 0.85$ | 97.86 $\pm 0.17$        | <u>72.14</u> $\pm 2.23$ | <b>61.03</b> $\pm 8.55$ | 49.22 $\pm 5.12$        | <u>78.01</u> |
| Prompt Tuning                 | 50K            | 68.23 $\pm 3.78$        | 61.02 $\pm 6.65$        | 84.81 $\pm 0.66$        | 87.75 $\pm 1.48$        | 51.61 $\pm 8.67$        | 36.13 $\pm 1.51$        | <u>54.69</u> $\pm 3.79$ | 63.46        |
| P-Tuning v2                   | 1.2M           | 64.33 $\pm 3.05$        | <b>92.63</b> $\pm 1.39$ | 83.46 $\pm 1.01$        | 97.05 $\pm 0.41$        | 68.14 $\pm 3.89$        | 36.89 $\pm 0.79$        | 50.78 $\pm 2.28$        | 70.47        |
| <i>Gradient-Free Methods</i>  |                |                         |                         |                         |                         |                         |                         |                         |              |
| Manual Prompt                 | 0              | 79.82                   | 89.65                   | 76.96                   | 41.33                   | 67.40                   | 31.11                   | 51.62                   | 62.56        |
| In-Context Learning           | 0              | 79.79 $\pm 3.06$        | 85.38 $\pm 3.92$        | 62.21 $\pm 13.46$       | 34.83 $\pm 7.59$        | 45.81 $\pm 6.67$        | <u>47.11</u> $\pm 0.63$ | <u>60.36</u> $\pm 1.56$ | 59.36        |
| Feature-MLP                   | 1M             | 64.80 $\pm 1.78$        | 79.20 $\pm 2.26$        | 70.77 $\pm 0.67$        | 87.78 $\pm 0.61$        | 68.40 $\pm 0.86$        | 42.01 $\pm 0.33$        | 53.43 $\pm 1.57$        | 66.63        |
| Feature-BiLSTM                | 17M            | 65.95 $\pm 0.99$        | 74.68 $\pm 0.10$        | 77.28 $\pm 2.83$        | <u>90.37</u> $\pm 3.10$ | <u>71.55</u> $\pm 7.10$ | 46.02 $\pm 0.38$        | 52.17 $\pm 0.25$        | 68.29        |
| BBT                           | 500            | <u>89.56</u> $\pm 0.25$ | <b>91.50</b> $\pm 0.16$ | <u>81.51</u> $\pm 0.79$ | 79.99* $\pm 2.95$       | 61.56 $\pm 4.34$        | 46.58 $\pm 1.33$        | 52.59 $\pm 2.21$        | <u>71.90</u> |
| <b>BBTv2</b>                  | 12K            | <b>90.41</b> $\pm 0.71$ | <u>90.69</u> $\pm 0.66$ | <b>85.06</b> $\pm 0.49$ | <b>92.59</b> $\pm 0.17$ | <b>78.15</b> $\pm 2.00$ | <b>61.50</b> $\pm 1.28$ | <b>60.56</b> $\pm 5.09$ | <b>79.85</b> |

# Experiments of BBTv2

## Versatility across different language models

| LM                          | Method | SST-2            | AG's News        | DBPedia          |
|-----------------------------|--------|------------------|------------------|------------------|
| <i>Encoder-only PTMs</i>    |        |                  |                  |                  |
| BERT                        | BBT    | 76.26 $\pm$ 2.64 | 76.67 $\pm$ 1.12 | 89.58 $\pm$ 0.51 |
|                             | BBTv2  | 79.32 $\pm$ 0.29 | 79.58 $\pm$ 1.15 | 93.74 $\pm$ 0.50 |
| RoBERTa                     | BBT    | 89.56 $\pm$ 0.25 | 81.51 $\pm$ 0.79 | 79.99 $\pm$ 2.95 |
|                             | BBTv2  | 90.41 $\pm$ 0.71 | 85.06 $\pm$ 0.49 | 92.59 $\pm$ 0.17 |
| <i>Decoder-only PTMs</i>    |        |                  |                  |                  |
| GPT-2                       | BBT    | 75.53 $\pm$ 1.98 | 77.63 $\pm$ 1.89 | 77.46 $\pm$ 0.69 |
|                             | BBTv2  | 80.13 $\pm$ 3.28 | 82.18 $\pm$ 1.07 | 91.36 $\pm$ 0.73 |
| <i>Encoder-Decoder PTMs</i> |        |                  |                  |                  |
| BART                        | BBT    | 77.87 $\pm$ 2.57 | 77.70 $\pm$ 2.46 | 79.64 $\pm$ 1.55 |
|                             | BBTv2  | 89.53 $\pm$ 2.02 | 81.30 $\pm$ 2.58 | 87.10 $\pm$ 2.01 |
| T5                          | BBT    | 89.15 $\pm$ 2.01 | 83.98 $\pm$ 1.87 | 92.76 $\pm$ 0.83 |
|                             | BBTv2  | 91.08 $\pm$ 1.49 | 84.32 $\pm$ 1.29 | 92.76 $\pm$ 0.85 |

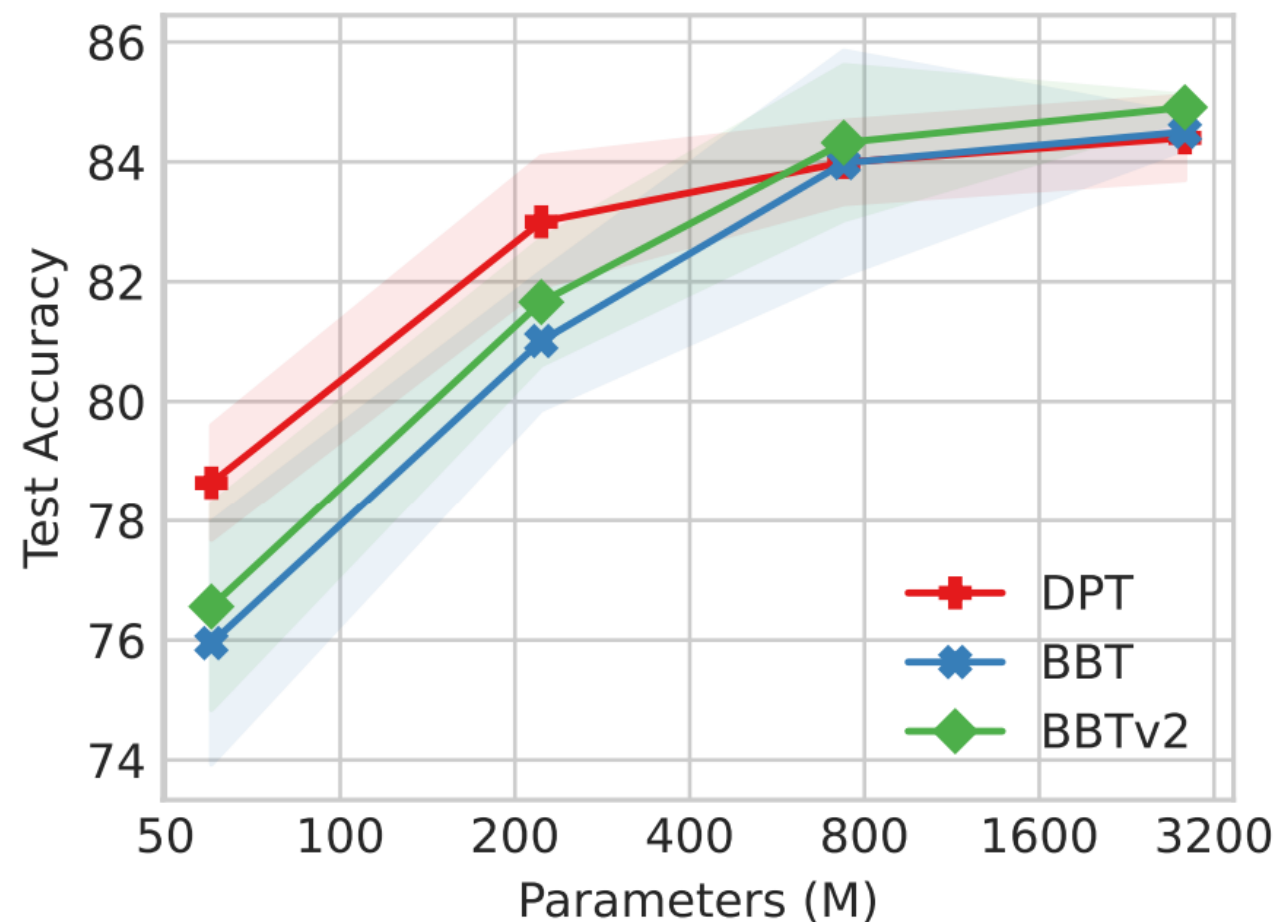
## Comparison on CPM-2 (11B)

| Method       | Tunable Params | ChnSent acc           | LCQMC acc             |
|--------------|----------------|-----------------------|-----------------------|
| Model Tuning | 11B            | <u>86.1</u> $\pm$ 1.8 | <u>58.8</u> $\pm$ 1.8 |
| Vanilla PT   | 410K           | 62.1 $\pm$ 3.1        | 51.5 $\pm$ 3.4        |
| Hybrid PT    | 410K           | 79.2 $\pm$ 4.0        | 54.6 $\pm$ 2.3        |
| LM Adaption  | 410K           | 74.3 $\pm$ 5.2        | 51.4 $\pm$ 2.9        |
| <b>BBTv2</b> | 4.8K           | <b>86.4</b> $\pm$ 0.8 | <b>59.1</b> $\pm$ 2.5 |

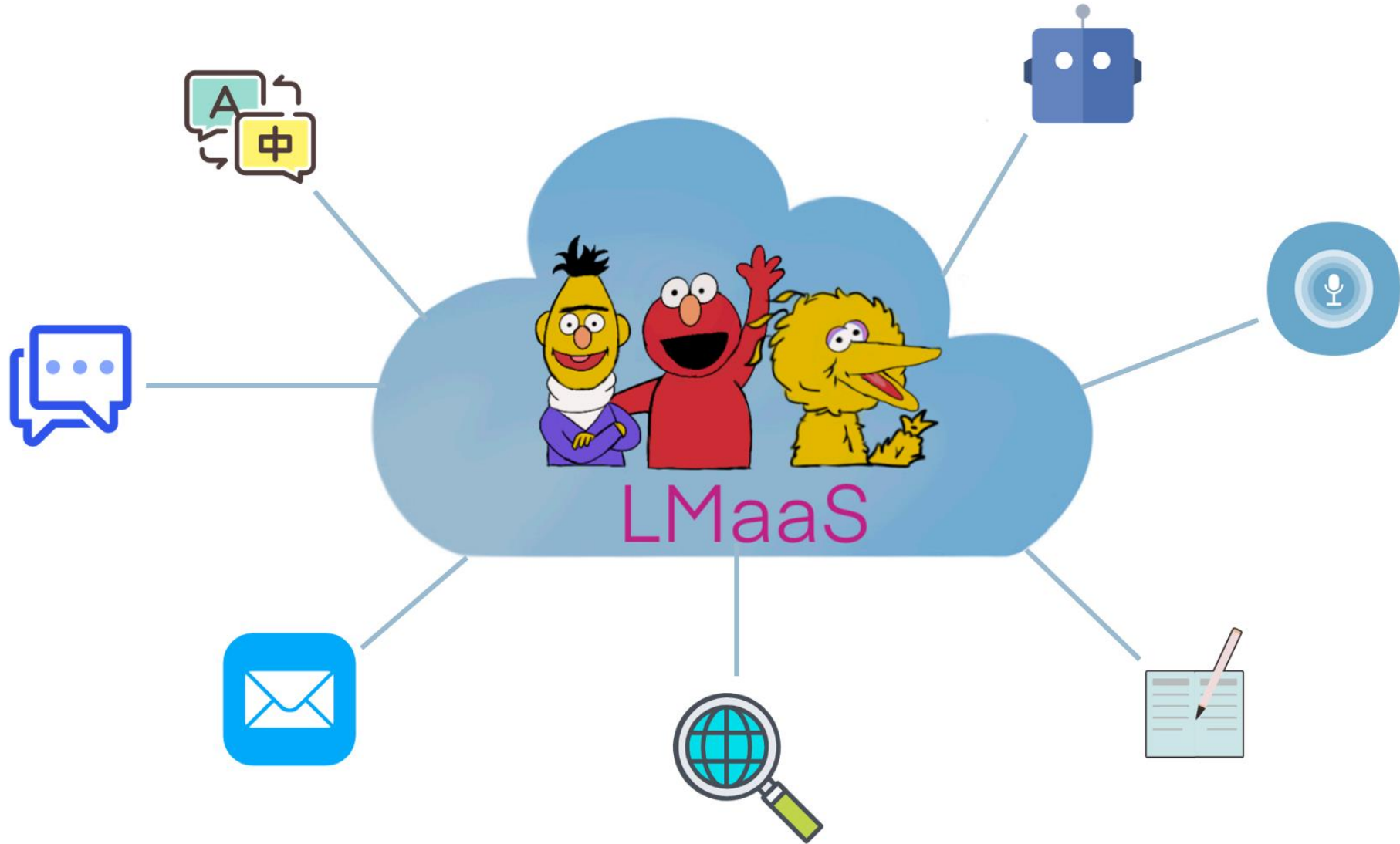
# Experiments of BBTv2

The power of scale (with T5)

- Outperform gradient descent when model size becomes large



# Other Solutions for LMaaS



# Other Solutions for LMaaS

- **Text prompt**: Manually or automatically design task-specific text prompts
- **In-context learning**: Include a few examples in the input at inference time
- **Black-box optimization**: Tuning a small portion of parameters with only the access of the LLM's output probability via black-box optimization (BBO)
- **Feature-based learning**: LLMs can serve as a feature extractor, on which users can build some lightweight learnable model to solve the task
- **Data generation**: Use LLMs to generate a dataset of labeled text pairs, which is then used to locally train a much smaller model

# Check Our Paper List!

☰ README.md ✎

## Language Model as a Service (LMaaS)

last commit today PaperNumber 43

This is a curated list of "Language-Model-as-a-Service (LMaaS)" papers, which is mainly maintained by [Tianxiang Sun](#). We strongly encourage the NLP researchers who are interested in this topic to make pull request to add or update the papers (See [Contributing](#)). Watch this repository for the latest updates!

### Updates

- 2022/7/7: Write a [blog](#) (in Chinese)
- 2022/7/4: Create this paper list

### Contents

- [Introduction](#)
  - [Scope](#)
  - [Advantages](#)
- [Keywords](#)
- [Papers](#)
  - [Text Prompt](#)
  - [In-Context Learning](#)
  - [Black-Box Optimization](#)
  - [Feature-based Learning](#)
  - [Data Generation](#)
- [Contributing](#)

📖 Readme

📄 779 KB

📄 MIT license

★ 87 stars

👁 3 watching

🍴 9 forks

### Releases

No releases published

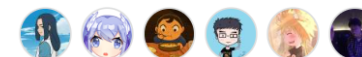
[Create a new release](#)

### Packages

No packages published

[Publish your first package](#)

### Contributors 6





# Resources

- LMaaS paper list: <https://github.com/txsun1997/LMaaS-Papers>
- Code of BBT and BBTv2: <https://github.com/txsun1997/Black-Box-Tuning>
- BBT paper (ICML 2022): <https://arxiv.org/abs/2201.03514>
- BBTv2 paper: <https://arxiv.org/abs/2205.11200>
- Blog for BBT (in Chinese): <https://zhuanlan.zhihu.com/p/455915295>
- Blog for LMaaS (in Chinese): <https://zhuanlan.zhihu.com/p/538857729>

Feel free to reach out if you have any questions or suggestions about our papers, code, or the paper list!



**ICML**  
International Conference  
On Machine Learning

# Thanks!

**Tianxiang Sun**

School of Computer Science, Fudan University

<https://txsun1997.github.io/>