

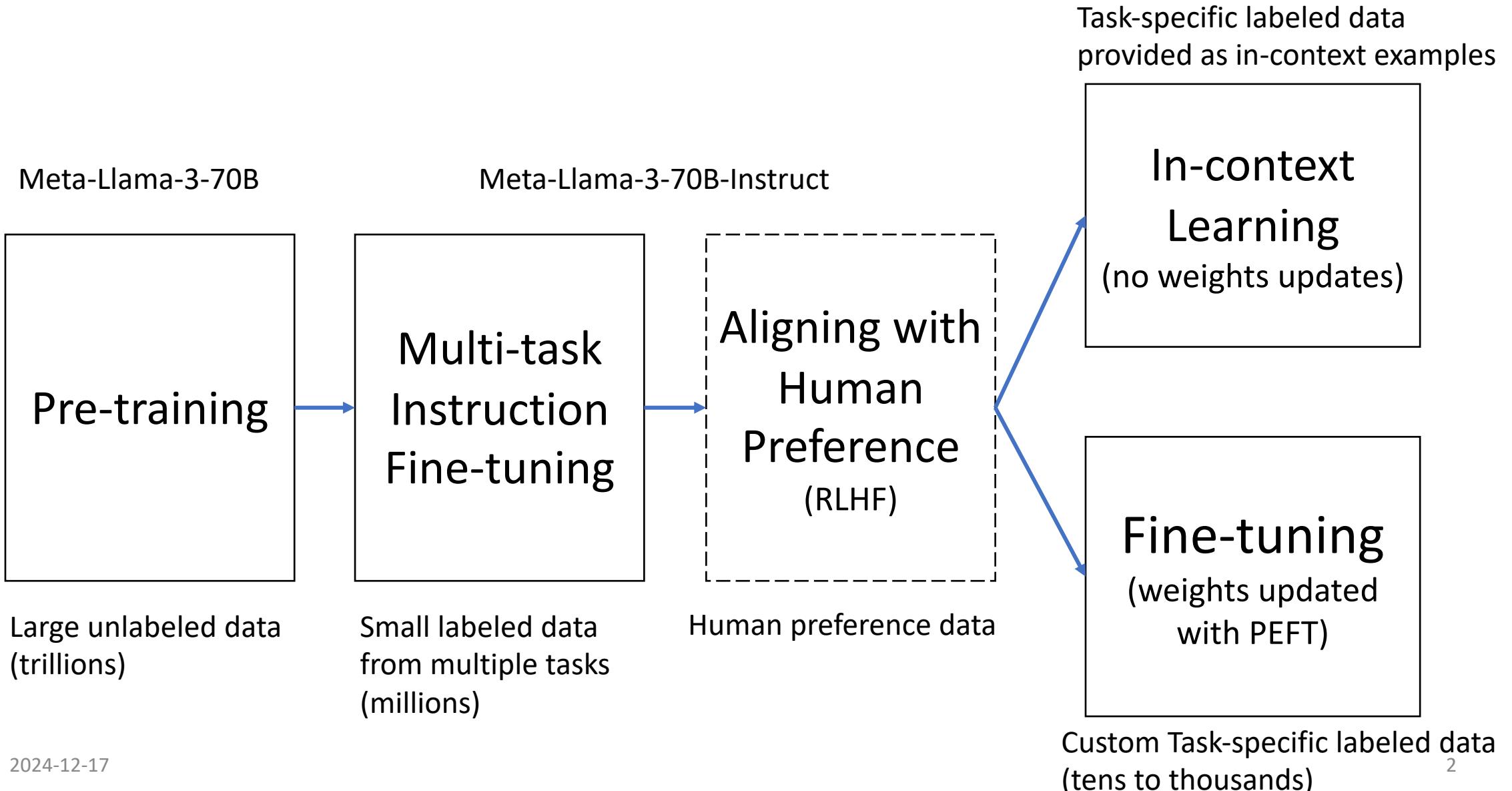
---

**Few-Shot Parameter-Efficient Fine-Tuning is Better  
and Cheaper than In-Context Learning**

# **LLM: Few-Shot Fine-Tuning V.S. In Context Learning**

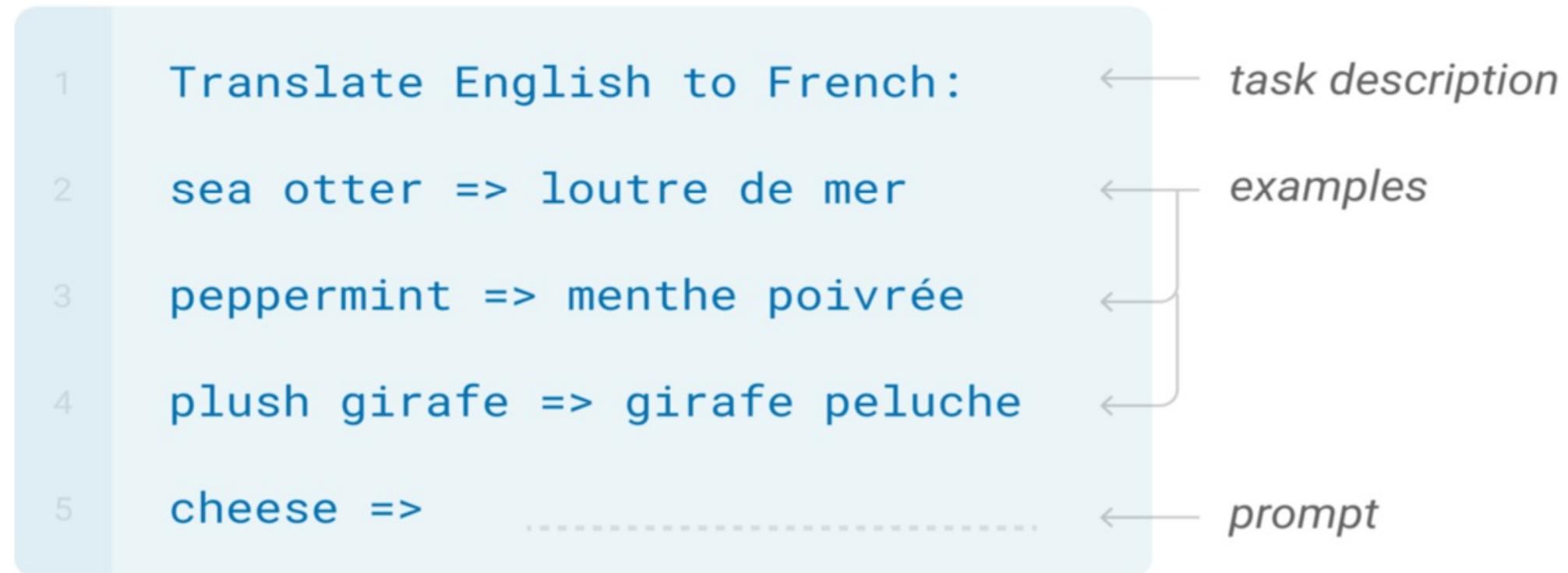
Houston Machine Learning Meetup  
12/13/2024

# LLM Training Paradigm



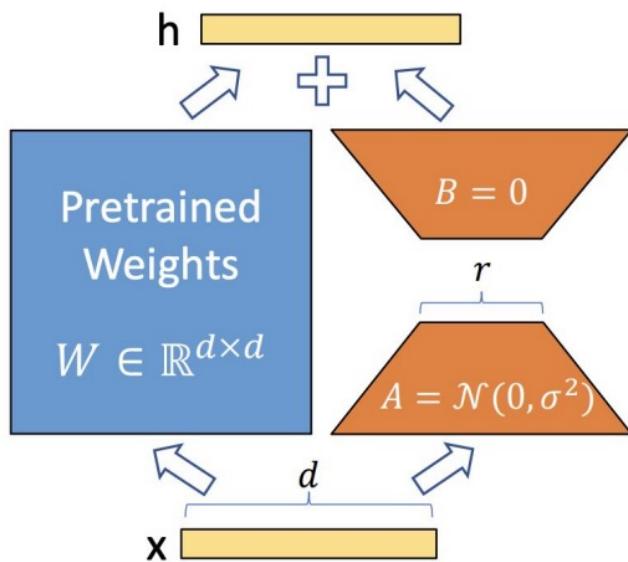
# Few-shot In-Context Learning (ICL)

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



# Few-shot Fine-tuning (PEFT)

- Fine-tuning is a process that involves adjusting a pretrained model's weights to better fit your specific data or task. This is different from in-context learning, as it actually trains a new model on your custom data.



$$h = W_0 x + \Delta W x = W_0 x + BAx$$

- Only update the low-rank matrix
- 10000x less trainable parameter
- 3x less GPU memory requirement
- Apply to any linear layer
- No inference overhead

Figure 1: Our reparametrization. We only train  $A$  and  $B$ .

# Comparison: ICL v.s. PEFT

## ICL

### PROS:

- ICL enables single models to perform many tasks without additional training

### CONS:

- Processing the  $k$  in-context examples increases the inference cost by approximately  $k$  times compared to processing the unlabeled example alone.
- ICL typically produces inferior performance compared to fine-tuning.
- The exact formatting of the prompt (including the wording and ordering of examples) can have significant and unpredictable impact on the model's performance.

## PEFT

### PROS:

- No additional inference cost compared to processing the unlabeled example directly
- Have significant better performance comparing to ICL

### CONS:

- One-time fine-tuning cost to train the model

---

# Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning

---

**Haokun Liu\***   **Derek Tam\***   **Mohammed Muqeeth\***

**Jay Mohta**   **Tenghao Huang**   **Mohit Bansal**   **Colin Raffel**

Department of Computer Science  
University of North Carolina at Chapel Hill

- Published at NIPS 2022, ~760 citations so far.
- This paper introduces  $(IA)^3$ : Infused Adapter by Inhibiting and Amplifying Inner Activations
- Argues that PEFT is better than ICL in few-shot settings
- Among the first to answer:

**whether PEFT methods work well when very little labeled data is available**

# Overview

(IA)<sup>3</sup>: “Infused Adapter by Inhibiting and Amplifying Inner Activations”

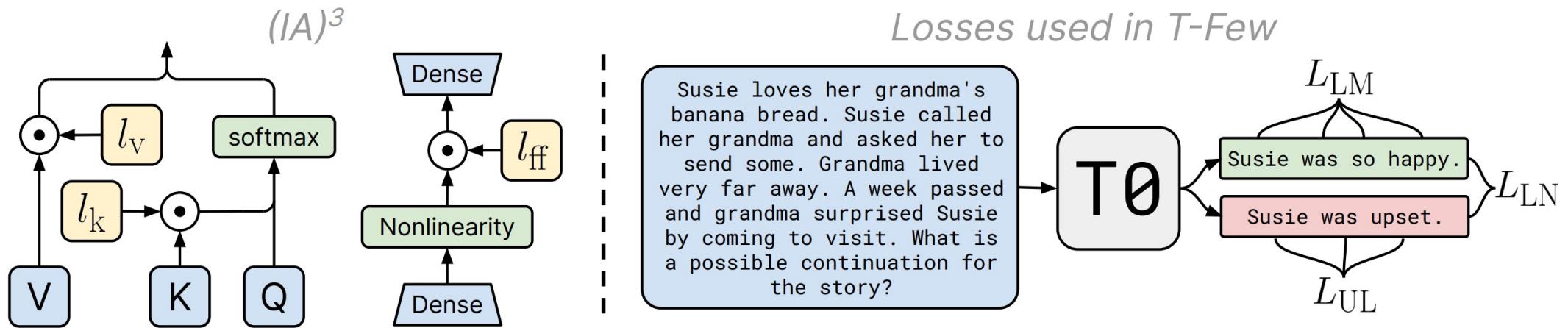
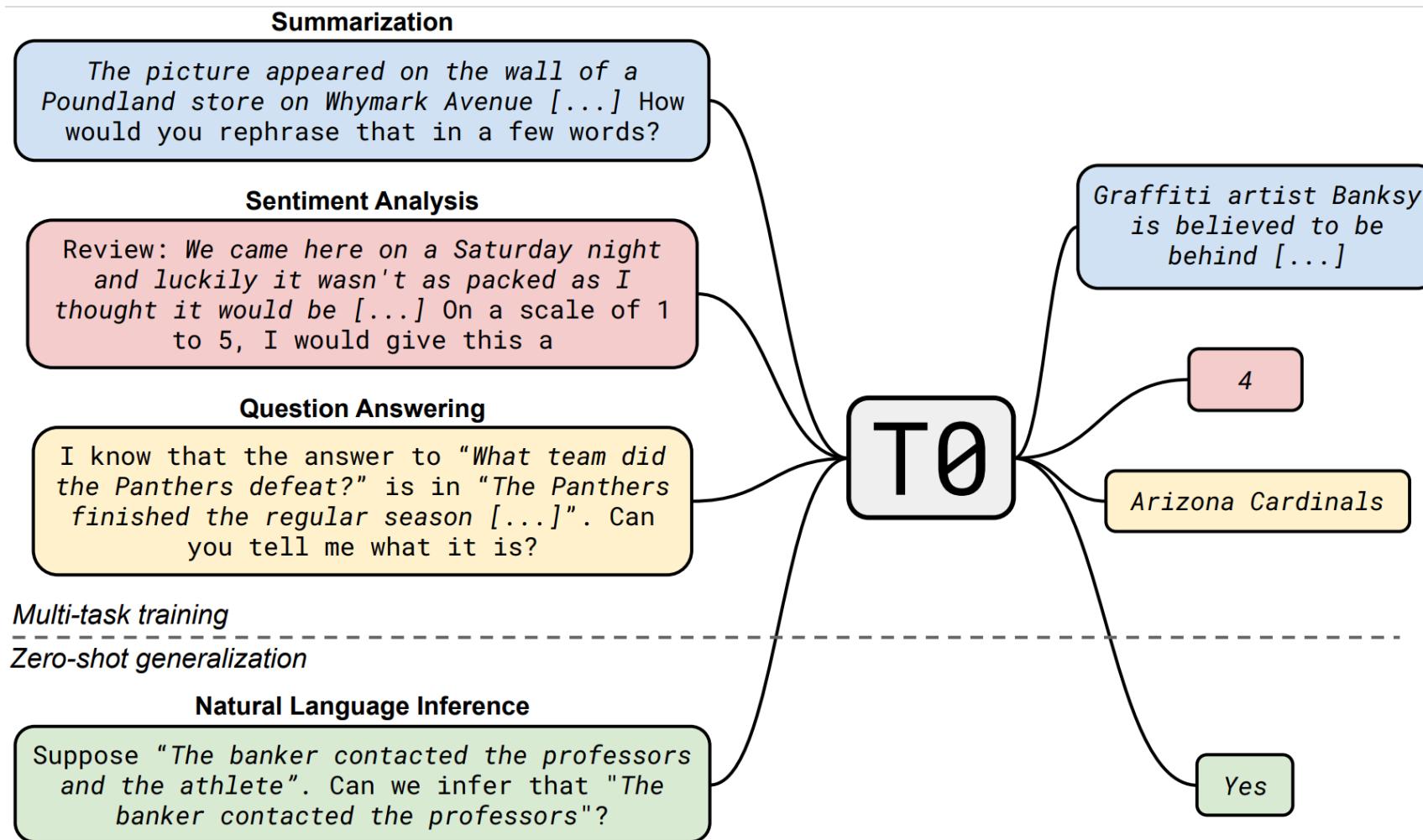


Figure 1: Diagram of (IA)<sup>3</sup> and the loss terms used in the T-Few recipe. *Left:* (IA)<sup>3</sup> introduces the learned vectors  $l_k$ ,  $l_v$ , and  $l_{ff}$  which respectively rescale (via element-wise multiplication, visualized as  $\odot$ ) the keys and values in attention mechanisms and the inner activations in position-wise feed-forward networks. *Right:* In addition to a standard cross-entropy loss  $L_{LM}$ , we introduce an unlikelihood loss  $L_{UL}$  that lowers the probability of incorrect outputs and a length-normalized loss  $L_{LN}$  that applies a standard softmax cross-entropy loss to length-normalized log-probabilities of all output choices.

# Choose a pre-trained model: T0



# Choose a pre-trained model: T0

- T0 is based on T5, an encoder-decoder Transformer model that was pre-trained via a masked language modeling objective on a large corpus of unlabeled text data.
- T0 was created by fine-tuning T5 on a multitask mixture of datasets in order to enable zero-shot generalization, i.e. the ability to perform tasks without any additional gradient-based training.
- T0 was released in 3B and 11B parameter variants, referred to as “T0-3B” and simply “T0” respectively
- T0 shows zero-shot task generalization on English natural language prompts, outperforming GPT-3 (175B) on many tasks, while being 16x smaller.

# Model loss: Unlikelihood loss

cross-entropy loss  $L_{LM} = -\frac{1}{T} \sum_t \log p(y_t | \mathbf{x}, y_{<t})$

$$L_{UL} = -\frac{\sum_{n=1}^N \sum_{t=1}^{T^{(n)}} \log(1 - p(\hat{y}_i^{(n)} | \mathbf{x}, \hat{y}_{<t}^{(n)}))}{\sum_{n=1}^N T^{(n)}}$$

We discourage the model from predicting tokens from incorrect target sequences, where  $\hat{y}^{(n)} = (\hat{y}_1^{(n)}, \hat{y}_2^{(n)}, \dots, \hat{y}_{T^{(n)}}^{(n)})$  is the  $n$ -th of  $N$  incorrect target sequences.

# Model loss: Length Normalization

$$\beta(\mathbf{x}, \mathbf{y}) = \frac{1}{T} \sum_{t=1}^T \log p(y_t | \mathbf{x}, y_{<t}).$$

$$L_{\text{LN}} = -\log \frac{\exp(\beta(\mathbf{x}, \mathbf{y}))}{\exp(\beta(\mathbf{x}, \mathbf{y})) + \sum_{n=1}^N \exp(\beta(\mathbf{x}, \hat{\mathbf{y}}^{(n)}))}$$

we maximize the length-normalized log probability of the correct answer choice

# Model PEFT: Recalining vectors

We found it was sufficient to introduce **rescaling vectors on the keys and values in self-attention and encoder-decoder attention mechanisms and on the intermediate activation of the position-wise feed-forward networks.**

the notation from Vaswani et al. [33], we introduce three learned vectors  $l_k \in \mathbb{R}^{\bar{d}_k}$ ,  $l_v \in \mathbb{R}^{\bar{d}_v}$ , and  $l_{ff} \in \mathbb{R}^{d_{ff}}$ , which are introduced into the attention mechanisms as:

$$\text{softmax}\left(\frac{Q(l_k \odot K^T)}{\sqrt{d_k}}\right) (l_v \odot V)$$

and in the position-wise feed-forward networks as  $(l_{ff} \odot \gamma(W_1 x))W_2$ , where  $\gamma$  is the feed-forward network nonlinearity. We introduce a separate set of  $l_k$ ,  $l_v$ , and  $l_{ff}$  vectors in each Transformer layer

•  $\odot$  represents element-wise multiplication

# Comparisons with other PEFT methods

We compare with 9 strong PEFT methods:

- **BitFit**: updates only the bias parameters
- **Adapters**: introduce task-specific layers after the self-attention and position-wise feed-forward network
- **Compacter and Compacter++**: improve upon adapters by using low-rank matrices and hypercomplex multiplication
- **Prompt tuning**: learns task-specific prompt embeddings that are concatenated to the model's input
- **FISH Mask**: chooses a subset of parameters to update based on their approximate Fisher information
- **Intrinsic SAID (Structure-aware intrinsic dimension)**: performs optimization in a low-dimensional subspace
- **Prefix-tuning**: learns task-specific vectors that are concatenated to the model's activations
- **LoRA**: assigns low-rank updates to parameter matrices.

Additionally, we include the baselines of **full-model fine-tuning**.

# Evaluation

- Evaluation is done using the T0 Model
- For each dataset tested, the number of few-shot examples used varies from 20 to 70, matching the same numbers as the evaluation from the paper: “Language models are few-shot learners” focusing on few-shot in context learning.

# Evaluation Datasets

- Datasets used for testing are sentence completion (COPA, H-SWAG, and Story Cloze), natural language inference (ANLI, CB, and RTE), coreference resolution (WSC and Winogrande), and word sense disambiguation (WiC).

## Sentence Completion

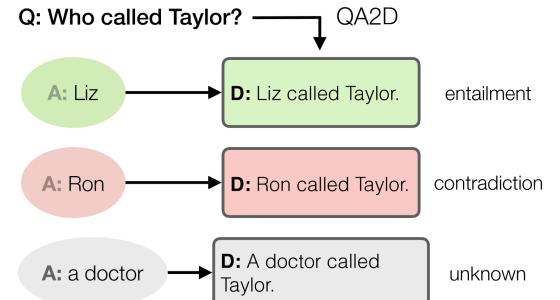
Premise: I tipped the bottle. What happened as a RESULT?  
Alternative 1: The liquid in the bottle froze.  
Alternative 2: The liquid in the bottle poured out.

## Coreference Resolution

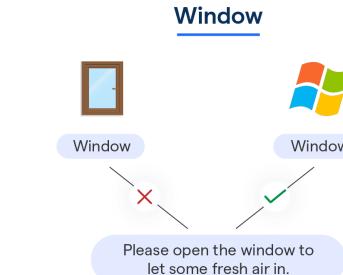
*"I voted for Nader because he was most aligned with my values," she said.*

## Natural Language Inference

**P:** Taylor is a journalist [...]. She was playing golf with Ron when her phone rang. It was Liz, her mother's friend. [...]



## Word Sense Disambiguation



# Results

	Situations With Adversarial Generations		
	COPA	H-Swag	Story
T-Few	93.0 <sub>2.0</sub>	67.1 <sub>6.0</sub>	97.1 <sub>1.0</sub>
T0	90.8	33.7	9.1 <sub>0.0</sub>
T5+LM	68.0	60.95	6.1 <sub>0.0</sub>
GPT-3 (175B)	92.0	79.3	8.1 <sub>0.0</sub>
GPT-3 (13B)	86.0	71.3	8.1 <sub>0.0</sub>
GPT-3 (6.7B)	83.0	67.3	8.1 <sub>0.0</sub>

	RTE	CB	SST-2
T-Few	85.6 <sub>2.9</sub>	87.5 <sub>3.6</sub>	90.1 <sub>1.0</sub>
T0	81.2	78.6	9.1 <sub>0.0</sub>
T5 + LM	53.4	32.1	6.1 <sub>0.0</sub>
GPT-3 (175B)	72.9	82.1	8.1 <sub>0.0</sub>
GPT-3 (13B)	60.6	66.1	8.1 <sub>0.0</sub>
GPT-3 (6.7B)	49.5	60.7	8.1 <sub>0.0</sub>

On stage, a woman takes a seat at the piano. She

- a) sits on a bench as her sister plays with the doll.
- b) smiles with someone as the music plays.
- c) is in the crowd, watching the dancers.
- d) nervously sets her fingers on the keys.**

A girl is going across a set of monkey bars. She

- a) jumps up across the monkey bars.
- b) struggles onto the monkey bars to grab her head.
- c) gets to the end and stands on a wooden plank.**
- d) jumps up and does a back flip.

The woman is now blow drying the dog. The dog

- a) is placed in the kennel next to a woman's feet.**
- b) washes her face with the shampoo.
- c) walks into frame and walks towards the dog.
- d) tried to cut her face, so she is trying to do something very close to her face.

Table 9: Comparing T-Few with few-shot ICL methods. All GPT-3 numbers are from Brown et al. [4] and all T0 numbers are from Sanh et al. [1]. Subscripts are IQR.

# T-Few shine in Inference Cost

- The T-Few method shines in terms of Inference FLOPS
- Inference FLOPs listed are FLOPS to make a prediction for a single example

Method	Inference FLOPs	Training FLOPs	Disk space	Acc.
T-Few	1.1e12	2.7e16	4.2 MB	72.4%
T0 [1]	1.1e12	0	0 B	66.9%
T5+LM [14]	4.5e13	0	16 kB	49.6%
GPT-3 6.7B [4]	5.4e13	0	16 kB	57.2%
GPT-3 13B [4]	1.0e14	0	16 kB	60.3%
GPT-3 175B [4]	1.4e15	0	16 kB	66.6%

**T-Few training costs is as much as using GPT-3 175B to process 20 inferences with few-shot ICL.**

# Training cost is not as big as you thought

- Training an eleven billion parameter encoder-decoder model for 1,000 steps with a batch size of 8 length-103 sequences requires approximately  $3 \times 11e9 \times 1,000 \times 8 \times 103 = 2.7e16$  FLOPs.
- While not insignificant, this is only about 20 times larger than the FLOPs required to process a single example with few-shot ICL using GPT-3 175B.
- **Training T-Few costs as much as using GPT-3 175B to process 20 examples with few-shot ICL.**
- We also found that fine-tuning T0 with T-Few on a single dataset only takes about half an hour on a single NVIDIA A100 GPU. As of writing, this would cost about \$2 USD using Microsoft Azure.

# Ablation study on the effects of pre-training (PT) and losses

	COPA	H-Swag	StoryCloze	Winogrande	WSC	WiC	
T-Few	93.0 <sub>2.0</sub>	67.1 <sub>6.0</sub>	97.9 <sub>0.3</sub>	74.3 <sub>1.5</sub>	75.0 <sub>5.5</sub>	62.1 <sub>5.8</sub>	PT: simply <b>pre-train</b> the new parameters introduced by (IA) <sup>3</sup> on the same multitask mixture used to train T0
- PT	92.0 <sub>2.0</sub>	64.5 <sub>6.6</sub>	97.8 <sub>0.8</sub>	72.7 <sub>1.0</sub>	73.1 <sub>6.3</sub>	60.8 <sub>6.4</sub>	
- $L_{UL}$ - $L_{LN}$	91.0 <sub>2.0</sub>	52.1 <sub>2.7</sub>	97.4 <sub>0.5</sub>	71.9 <sub>1.1</sub>	71.2 <sub>1.0</sub>	62.2 <sub>2.4</sub>	
- PT - $L_{UL}$ - $L_{LN}$	94.0 <sub>2.3</sub>	52.7 <sub>4.9</sub>	98.0 <sub>0.3</sub>	74.0 <sub>1.1</sub>	72.6 <sub>4.8</sub>	62.6 <sub>5.0</sub>	
	RTE	CB	ANLI-R1	ANLI-R2	ANLI-R3	Acc.	
T-Few	85.6 <sub>2.9</sub>	87.5 <sub>3.6</sub>	59.3 <sub>3.6</sub>	49.8 <sub>2.6</sub>	44.8 <sub>8.0</sub>	72.4	
- PT	84.5 <sub>2.8</sub>	83.9 <sub>5.4</sub>	57.9 <sub>3.2</sub>	48.6 <sub>3.0</sub>	43.1 <sub>5.7</sub>	70.8	
- $L_{UL}$ - $L_{LN}$	82.0 <sub>0.7</sub>	82.1 <sub>3.6</sub>	54.8 <sub>0.4</sub>	46.1 <sub>0.6</sub>	40.8 <sub>5.2</sub>	68.3	
- PT - $L_{UL}$ - $L_{LN}$	84.5 <sub>2.9</sub>	80.4 <sub>3.6</sub>	57.1 <sub>3.1</sub>	47.1 <sub>2.4</sub>	43.8 <sub>5.9</sub>	69.7	

Table 10: T-Few ablation results when omitting (IA)<sup>3</sup> pre-training (PT) and/or the  $L_{UL}$  and  $L_{LN}$  losses. Subscripts are IQR.

# RAFT benchmark: Real-world Annotated Few-shot Tasks



Task Specification in Natural Language

Organization	Paper Title	Label
North Carolina State Univ., Raleigh, NC, USA	3Gb/s AC-coupled chip-to-chip communication using a low-swing pulse receiver	university
Advanced LCD Technology Development Center...	Sub-Micron CMOS / MOS-Bipolar Hybrid TFTs for System Displays	company
imec, Heverlee, Belgium	24.4 A 680nA fully integrated implantable ECG-acquisition IC with analog...	research institute
...	...	...

Partially-labeled Dataset

The dataset is a list of institutions that have contributed papers to semiconductor conferences in the last 25 years, as catalogued by IEEE and sampled randomly. The goal is to classify the institutions into one of three categories: "university", "company" or "research institute". 50 labeled examples are provided.

Label the sentence based on whether it is related to an adverse drug effect (ADE)

The following is a banking customer service query. Classify it as one of the following 77 classes.

Label the impact statement as "mentions a harmful application" or "doesn't mention a harmful application" based on whether it mentions a harmful application of the research done in the paper.

Identify whether this paper should be included in a meta-review which includes the findings of systematic reviews on interventions designed to promote charitable donations.

...

Other Tasks

Figure 1: RAFT includes naturally occurring classification datasets, mimicking work that is usually given to human research assistants. Each task comes with natural language instructions and labels in addition to 50 training examples.

# RAFT benchmark: Real-world Annotated Few-shot Tasks

<b>Dataset Name</b>	<i>Long inputs</i>	<i>Domain expertise</i>	<i>Detailed instructions</i>	<i>Number of classes</i>	<i>Test set size</i>
ADE Corpus V2 ( <i>ADE</i> )	–	✓	–	2	5,000
Banking77 ( <i>B77</i> )	–	–	–	77	5,000
NeurIPS impact statement risks ( <i>NIS</i> )	✓	–	–	2	150
OneStopEnglish ( <i>OSE</i> )	✓	–	–	3	516
Overruling ( <i>Over</i> )	–	✓	–	2	2,350
Semiconductor org types ( <i>SOT</i> )	–	–	–	3	449
Systematic review inclusion ( <i>SRI</i> )	✓	–	✓	2	2,243
TAI safety research ( <i>TAI</i> )	✓	✓	✓	2	1,639
Terms of Service ( <i>ToS</i> )	–	✓	✓	2	5,000
TweetEval Hate ( <i>TEH</i> )	–	–	✓	2	2,966
Twitter complaints ( <i>TC</i> )	–	–	–	2	3,399

Table 1: Overview of the tasks in RAFT. *Long inputs*, *Domain expertise*, and *Detailed instructions* are some of the real-world challenges posed by RAFT.

# RAFT results

To collect human baselines, we use the Surge crowdsourcing platform. We randomly select 100 data points from each test set and use a 2-step labeling process: qualification then annotation. The crowdsourced label is the plurality vote of 5 labelers.

Method	Ade Corpus V2	Banking 77	Neurips Impact Statement Risks	One Stop English	OVERRULING	Semiconductor Org Types	Systematic Review Inclusion	Tai Safety Research	Terms Of Service	Tweet Eval Hate	Twitter Complaints
T-Few	80.4	69.5	83.3	67.6	95.0	91.5	50.8	73.6	75.0	58.6	87.9
Human baseline [2]	83.0	60.7	85.7	64.6	91.7	90.8	46.8	60.9	62.7	72.2	89.7
PET [50]	82.2	59.3	85.7	64.6	90.8	81.6	49.3	63.8	57.6	48.3	82.4
SetFit [51]	72.6	53.8	87.2	52.1	90.7	68.2	49.3	62.8	62.0	53.2	83.7
GPT-3 [4]	68.6	29.9	67.9	43.1	93.7	76.9	51.6	65.6	57.4	52.6	82.1

Table 11: Detailed per-dataset results for T-Few and the other top-5 methods on RAFT.

Metric: Macro-average F1

Method	Acc.
T-Few	75.8%
Human baseline [2]	73.5%
PET [50]	69.6%
SetFit [51]	66.9%
GPT-3 [4]	62.7%

Table 2: Top-5 best methods on RAFT as of writing. T-Few is the first method to outperform the human baseline and achieves over 6% higher accuracy than the next-best method.

# RAFT: Comparison to other PEFT methods

	# of Param	COPA	H-Swag	StoryCloze	Winogrande1	WSC	WiC	RTE	CB	ANLI-R1	ANLI-R2	ANLI-R3
Full Model Fine-tuning	3B	81.0 <sub>11.0</sub>	46.4 <sub>8.8</sub>	93.8 <sub>2.7</sub>	56.5 <sub>1.5</sub>	65.4 <sub>7.7</sub>	57.7 <sub>3.9</sub>	79.8 <sub>3.6</sub>	87.5 <sub>5.4</sub>	46.6 <sub>2.5</sub>	41.3 <sub>0.9</sub>	40.2 <sub>5.3</sub>
BitFit (with LayerNorm)	1.3M	75.0 <sub>2.0</sub>	29.5 <sub>3.6</sub>	88.6 <sub>0.7</sub>	49.6 <sub>1.3</sub>	61.5 <sub>11.5</sub>	51.7 <sub>2.2</sub>	72.2 <sub>1.1</sub>	57.1 <sub>1.8</sub>	36.5 <sub>0.8</sub>	35.3 <sub>2.2</sub>	36.6 <sub>0.8</sub>
LayerNorm	250K	76.0 <sub>2.0</sub>	29.6 <sub>3.4</sub>	88.7 <sub>0.9</sub>	49.4 <sub>1.4</sub>	63.5 <sub>12.5</sub>	52.2 <sub>1.6</sub>	71.8 <sub>0.4</sub>	57.1 <sub>1.8</sub>	36.5 <sub>0.7</sub>	35.1 <sub>2.6</sub>	36.3 <sub>1.0</sub>
Adapter	12.9M	84.0 <sub>3.0</sub>	41.9 <sub>3.8</sub>	91.7 <sub>3.7</sub>	54.7 <sub>3.6</sub>	65.4 <sub>1.0</sub>	55.5 <sub>2.7</sub>	76.2 <sub>3.6</sub>	87.5 <sub>3.6</sub>	45.1 <sub>2.6</sub>	40.4 <sub>1.2</sub>	35.3 <sub>1.3</sub>
Compacter	807K	84.0 <sub>5.0</sub>	46.4 <sub>2.5</sub>	93.5 <sub>2.2</sub>	55.5 <sub>2.9</sub>	64.4 <sub>6.7</sub>	55.2 <sub>3.8</sub>	75.8 <sub>6.1</sub>	82.1 <sub>3.6</sub>	40.8 <sub>3.3</sub>	37.4 <sub>0.2</sub>	35.8 <sub>3.3</sub>
Compacter++	540K	86.0 <sub>3.0</sub>	46.3 <sub>3.0</sub>	93.5 <sub>1.2</sub>	55.1 <sub>1.1</sub>	65.4 <sub>3.9</sub>	54.1 <sub>2.2</sub>	76.9 <sub>0.4</sub>	82.1 <sub>3.6</sub>	41.7 <sub>0.4</sub>	38.3 <sub>1.8</sub>	36.9 <sub>1.5</sub>
Prompt tuning (10)	41K	67.0 <sub>5.0</sub>	29.9 <sub>0.6</sub>	84.2 <sub>0.8</sub>	51.9 <sub>1.6</sub>	54.8 <sub>10.6</sub>	51.6 <sub>2.0</sub>	52.7 <sub>5.4</sub>	66.1 <sub>1.8</sub>	34.2 <sub>1.9</sub>	33.5 <sub>1.1</sub>	33.5 <sub>1.3</sub>
Prompt tuning (100)	409K	60.0 <sub>19.0</sub>	26.8 <sub>0.6</sub>	74.0 <sub>3.4</sub>	51.1 <sub>0.8</sub>	60.6 <sub>4.8</sub>	50.0 <sub>1.1</sub>	48.0 <sub>2.9</sub>	53.6 <sub>17.9</sub>	33.4 <sub>1.2</sub>	33.8 <sub>0.5</sub>	33.3 <sub>0.8</sub>
Prefix tuning	576K	71.0 <sub>8.0</sub>	42.1 <sub>4.0</sub>	90.2 <sub>3.1</sub>	52.0 <sub>1.3</sub>	56.7 <sub>3.3</sub>	54.2 <sub>3.3</sub>	68.6 <sub>3.3</sub>	84.0 <sub>1.8</sub>	43.3 <sub>4.1</sub>	37.5 <sub>1.2</sub>	36.5 <sub>1.5</sub>
FishMask (0.2%)	6M	82.0 <sub>5.0</sub>	44.1 <sub>4.2</sub>	94.2 <sub>1.8</sub>	54.5 <sub>2.1</sub>	63.5 <sub>4.8</sub>	52.5 <sub>3.3</sub>	76.9 <sub>4.7</sub>	83.9 <sub>3.6</sub>	43.7 <sub>0.3</sub>	39.7 <sub>1.4</sub>	37.2 <sub>1.1</sub>
FishMask (0.02%)	600K	84.0 <sub>6.0</sub>	38.2 <sub>3.6</sub>	93.6 <sub>0.7</sub>	53.9 <sub>2.8</sub>	61.5 <sub>1.0</sub>	53.5 <sub>1.3</sub>	75.5 <sub>5.4</sub>	76.8 <sub>3.6</sub>	39.9 <sub>0.9</sub>	38.1 <sub>2.0</sub>	36.2 <sub>1.8</sub>
Intrinsic SAID	500K	77.0 <sub>4.0</sub>	36.7 <sub>4.5</sub>	89.3 <sub>2.3</sub>	52.7 <sub>2.1</sub>	61.5 <sub>8.7</sub>	55.0 <sub>2.7</sub>	69.0 <sub>7.6</sub>	80.4 <sub>0.0</sub>	40.4 <sub>3.3</sub>	35.4 <sub>4.1</sub>	35.5 <sub>1.6</sub>
Intrinsic SAID	20K	76.0 <sub>4.0</sub>	38.3 <sub>6.4</sub>	89.7 <sub>2.7</sub>	50.9 <sub>1.0</sub>	55.8 <sub>6.7</sub>	55.3 <sub>0.5</sub>	66.1 <sub>5.4</sub>	83.9 <sub>1.8</sub>	41.3 <sub>1.3</sub>	38.5 <sub>1.8</sub>	35.8 <sub>2.0</sub>
LoRA	9.1M	88.0 <sub>5.0</sub>	47.1 <sub>3.2</sub>	93.6 <sub>2.1</sub>	56.8 <sub>3.3</sub>	60.6 <sub>5.8</sub>	55.2 <sub>5.0</sub>	78.3 <sub>7.6</sub>	85.7 <sub>1.8</sub>	45.1 <sub>2.5</sub>	41.0 <sub>1.4</sub>	39.5 <sub>4.8</sub>
(IA) <sup>3</sup>	540K	87.0 <sub>3.0</sub>	49.4 <sub>4.6</sub>	94.7 <sub>2.7</sub>	59.8 <sub>0.6</sub>	68.3 <sub>6.7</sub>	56.0 <sub>4.6</sub>	78.0 <sub>2.5</sub>	87.5 <sub>1.8</sub>	48.6 <sub>2.0</sub>	40.8 <sub>1.5</sub>	40.8 <sub>2.3</sub>

# Conclusion

- We introduced T-Few, a parameter-efficient few-shot learning recipe that attains **higher accuracy than few-shot ICL** at a lower computational cost and provide a new perspective on how best to **perform few-shot learning** with large language models
- When applying T-Few to the RAFT benchmark, we **attained super-human performance** for the first time and outperformed prior submissions by a large margin.
- T-Few **reduces inference cost with not-so-big one-time training cost**: uses over 1,000× fewer FLOPs during inference than few-shot ICL with GPT-3 and only requires 30 minutes to train on a single NVIDIA A100 GPU (equivalent to 20 GPT-3 inference calls)
- Caveat: all of the experiments were on multiple choices, classification tasks, we are interested in applying T-Few to generative tasks like as summarization and question answering in future work.

# Thank you



## YanAI Talk

@yanaitalk · 2.2K subscribers · 61 videos

Make machine learning easy to understand! ..

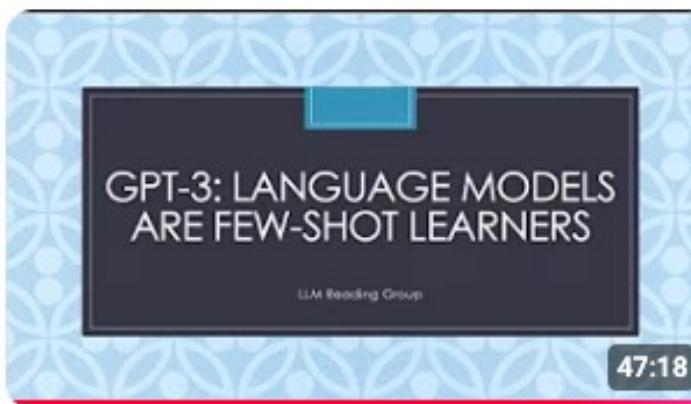
Customize channel

Manage videos



LLM: Exploring the Limits of Transfer Learning with a unified Text-to-Text...

673 views · 1 year ago



LLM: GPT-3

161 views · 1 year ago



LLM: Pretraining, Instruction fine-tuning and RLHF

5.8K views · 1 year ago

**LORA and QLORA:**  
Parameter-Efficient Fine-tuning of LLMs

Low-Rank Adaption

PEFT

Quantized LLMs

48:25

Parameter-efficient Fine-tuning of LLMs with LoRA

308 views · 3 weeks ago