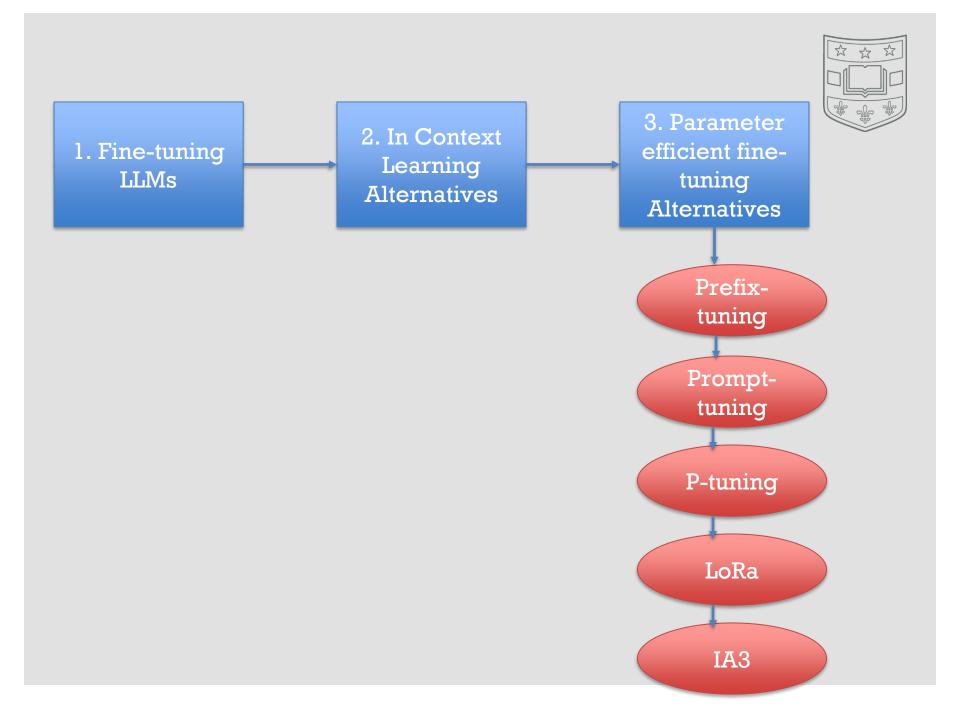
Parameter efficient fine-tuning

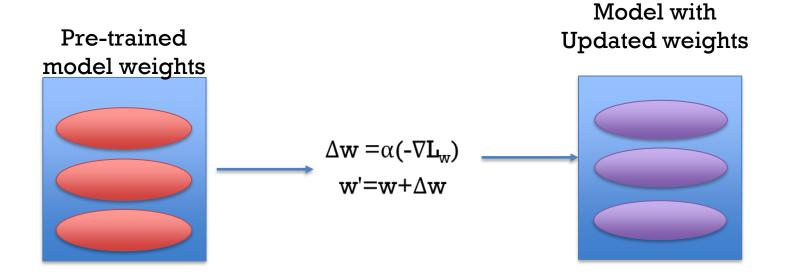






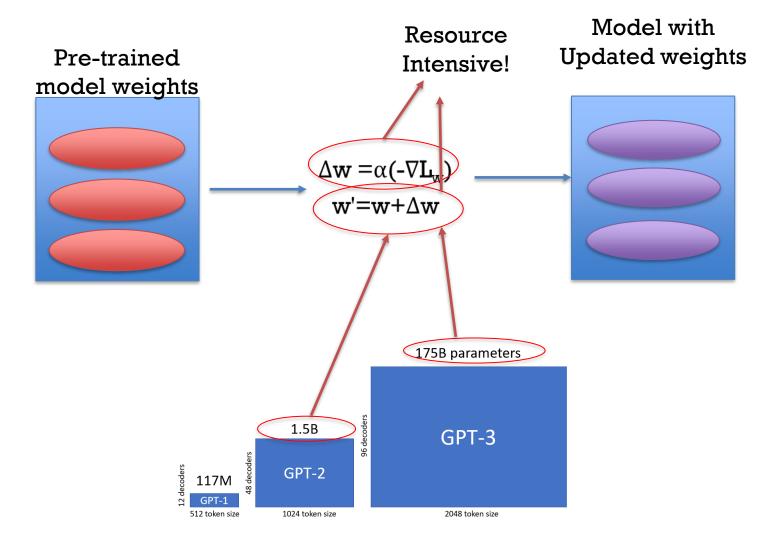
# Fine-tuning LLMs (or any deep learning model)





# Problem 1: As models get larger, fine-tuning becomes more computationally intensive!





### Other problems

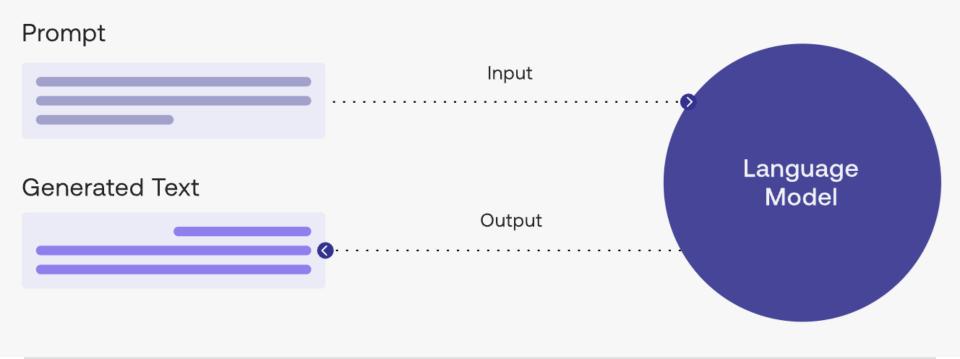


- 1. Specialized for a single domain
  - Have to fine-tune for each specific domain
- 2. Require sufficiently **abundant** and **representative** data
  - Labelling can be expensive!!!

### In context learning



Idea: Feed it input-target examples ("shots")



### In context learning (test)



#### Example of a 4-shot input:

{Please unscramble the letters into a word, and write that word:

- asinoc=casino,
- yfrogg=froggy,
- plesim=simple,
- iggestb=biggest,
- astedro=?}

### In context learning (test)



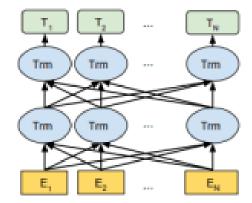
#### Example of a 4-shot input:

{Please unscramble the letters into a word, and write that word:

- asinoc=casino,
- yfrogg=froggy,
- plesim=simple,
- iggestb=biggest,
- astedro=<u>roasted</u>}

## Advantages to ICL

- 1. Enables single model to perform many task immediately without fine-tuning
- 2. Allows for mixed task batches
  - ☐ different examples in a batch of data correspond to different task.

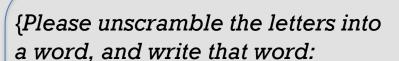




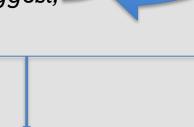


#### Problems with ICL

#### Can display unintuitive behavior.



- asinoc=casino, 🍁
- yfrogg=froggy,
- plesim=simple,
- iggestb=biggest,
- astedro=?}





#### Calibrate Before Use: Improving Few-Shot Performance of Language Models

Tony Z. Zhao \*1 Eric Wallace \*1 Shi Feng 2 Dan Klein 1 Sameer Singh 3

Is it really learning???

### Problems with ICL





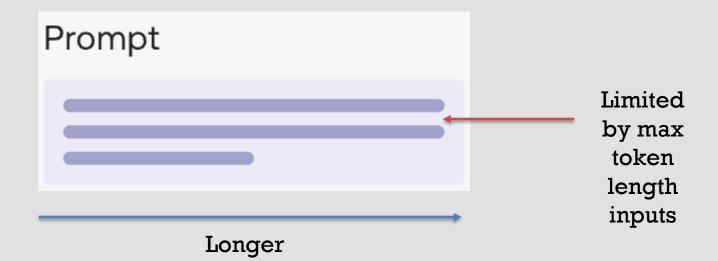




More expensive

Memory







1. Fine-tuning
LLMs

2. In Context
Learning
Alternatives

3. Parameter efficient finetuning
Alternatives

# Welcome to the world of parameter efficient fine-tuning





# What is parameter efficient fine-tuning? Definition by Li & Liang (2021)





Lightweight fine-tuning



Freezes most of the pretrained parameters and modifies the pretrained model with small trainable modules



Goal: identify high-performing architectures of the modules and the subset of pretrained parameters to tune

## Definition by Liu et al (2022)





Provides **strong** performance on any task



Without manual tuning or per-task adjustments



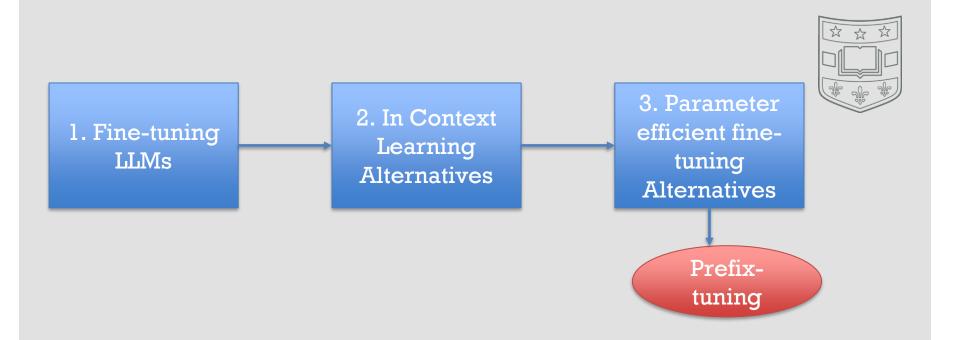
Realistic option in <u>few-shot settings</u> where limited labeled data is available.

#### How?



Exploiting tricks that allow the model to be tuned but adjusting only a small amount of parameters

While keeping the model's remaining weights **frozen** 







ACL Anthology FAQ Corrections Submissions

Search...



#### **Prefix-Tuning: Optimizing Continuous Prompts for Generation**

Xiang Lisa Li, Percy Liang

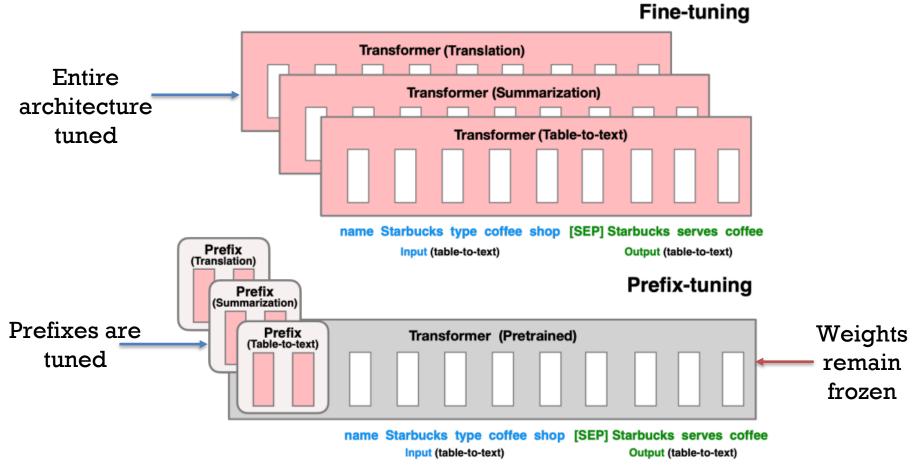
#### **Abstract**

Fine-tuning is the de facto way of leveraging large pretrained language models for downstream tasks. However, fine-tuning modifies all the language model parameters and therefore necessitates storing a full copy for each task. In this paper, we propose prefix-tuning, a lightweight alternative to fine-tuning for natural language generation tasks, which keeps language model parameters frozen and instead optimizes a sequence of continuous task-specific vectors, which we call the prefix. Prefix-tuning draws inspiration from prompting for language models, allowing subsequent tokens to attend to this prefix as if it were "virtual tokens". We apply prefix-tuning to GPT-2 for table-to-text generation and to BART for summarization. We show that by learning only 0.1% of the parameters, prefix-tuning obtains comparable performance in the full data setting, outperforms fine-tuning in low-data settings, and extrapolates better to examples with topics that are unseen during training.

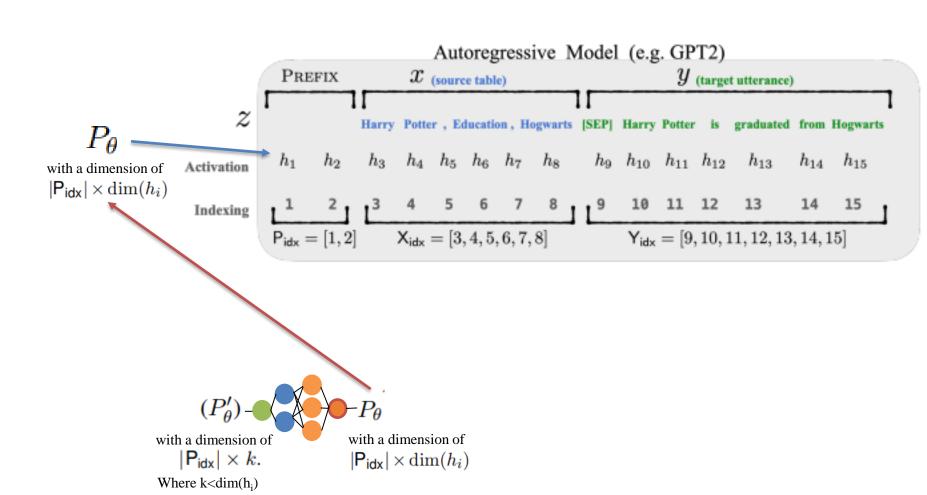


### Prefix Tuning

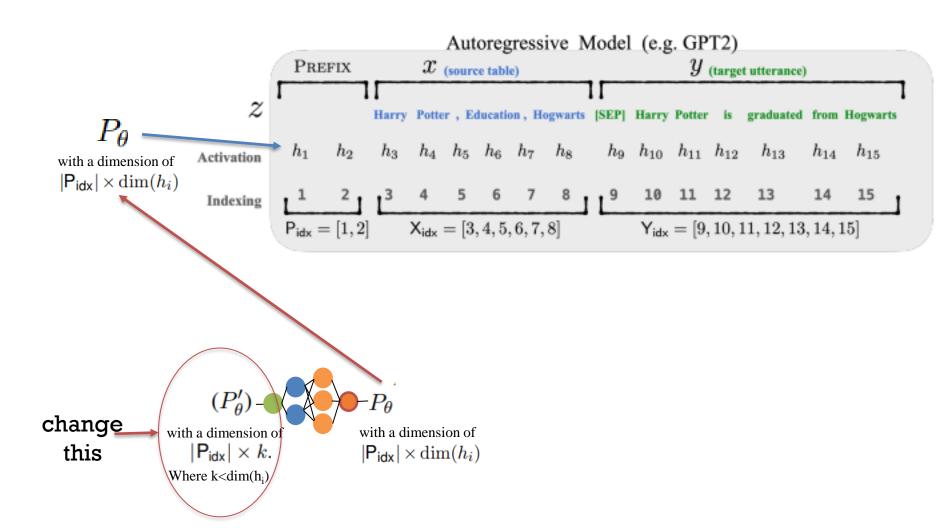








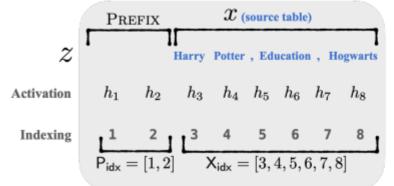


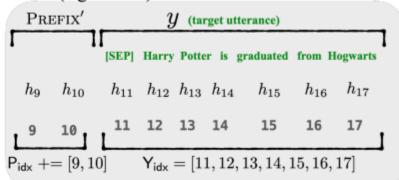




одрини сан

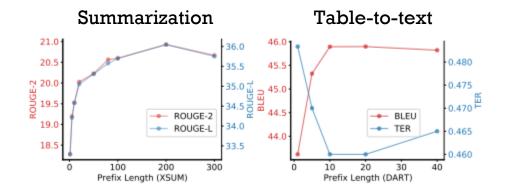






# Ideal prefix length





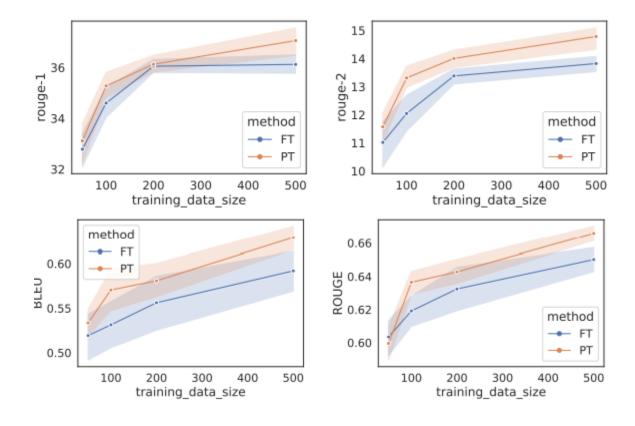
### Performances

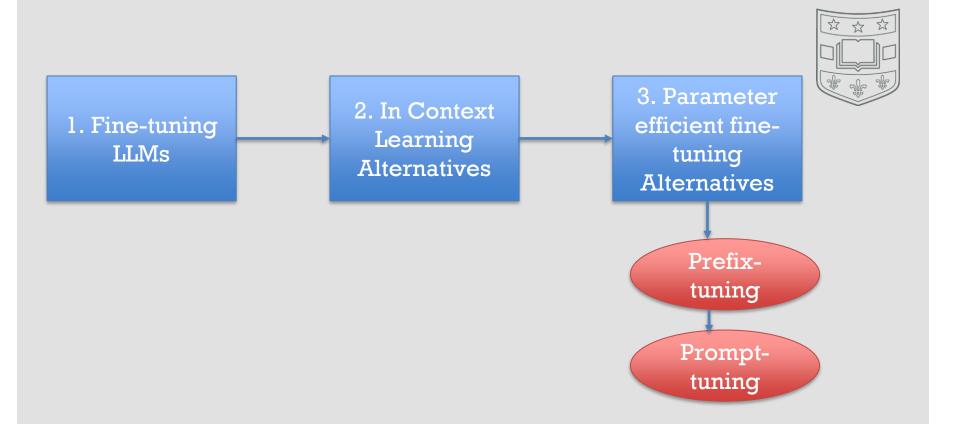


			E2E						V	VebNL	.G				1		Ι	OART		
	BLEU	NIST	MET	R-L	CIDEr		BLEU	J		MET			TER ↓		BLEU	MET	TER $\downarrow$	Mover	BERT	BLEURT
						S	U	Α	S	U	A	S	U	A						
										GP	T-2 <sub>ME</sub>	DIUM								
FINE-TUNE	68.2	8.62	46.2	71.0	2.47	64.2	27.7	46.5	0.45				0.76	0.53	46.2	0.39	0.46	0.50	0.94	0.39
FT-TOP2	68.1	8.59	46.0	70.8	2.41	53.6	18.9	36.0	0.38	0.23	0.31	0.49	0.99	0.72	41.0	0.34	0.56	0.43	0.93	0.21
ADAPTER(3%)	68.9	8.71	46.1	71.3	2.47	60.4	48.3	54.9	0.43	0.38	0.41	0.35	0.45	0.39	45.2	0.38	0.46	0.50	0.94	0.39
ADAPTER(0.1%)	66.3	8.41	45.0	69.8	2.40	54.5	45.1	50.2	0.39	0.36	0.38	0.40	0.46	0.43	42.4	0.36	0.48	0.47	0.94	0.33
Prefix(0.1%)	<b>69.7</b>	8.81	46.1	71.4	2.49	62.9	45.6	55.1	0.44	0.38	0.41	0.35	0.49	0.41	46.4	0.38	0.46	0.50	0.94	0.39
										G	PT-2L	ARGE								
FINE-TUNE	68.5	8.78	46.0	69.9	2.45	65.3	43.1	55.5	0.46	0.38	0.42	0.33	0.53	0.42	47.0	0.39	0.46	0.51	0.94	0.40
Prefix	70.3	8.85	46.2	71.7	2.47	63.4	47.7	56.3	0.45	0.39	0.42	0.34	0.48	0.40	46.7	0.39	0.45	0.51	0.94	0.40
SOTA	68.6	8.70	45.3	70.8	2.37	63.9	52.8	57.1	0.46	0.41	0.44	-	-	-	-	-	-	-	-	-

# Performances under low data settings







Search...

#### The Power of Scale for Parameter-Efficient Prompt Tuning

Brian Lester, Rami Al-Rfou, Noah Constant

#### Abstract

In this work, we explore "prompt tuning," a simple yet effective mechanism for learning "soft prompts" to condition frozen language models to perform specific downstream tasks. Unlike the discrete text prompts used by GPT-3, soft prompts are learned through backpropagation and can be tuned to incorporate signals from any number of labeled examples. Our end-to-end learned approach outperforms GPT-3's few-shot learning by a large margin. More remarkably, through ablations on model size using T5, we show that prompt tuning becomes more competitive with scale: as models exceed billions of parameters, our method "closes the gap" and matches the strong performance of model tuning (where all model weights are tuned). This finding is especially relevant because large models are costly to share and serve and the ability to reuse one frozen model for multiple downstream tasks can ease this burden. Our method can be seen as a simplification of the recently proposed "prefix tuning" of Li and Liang (2021) and we provide a comparison to this and other similar approaches. Finally, we show that conditioning a frozen model with soft prompts confers benefits in robustness to domain transfer and enables efficient "prompt ensembling." We release code and model checkpoints to reproduce our experiments.





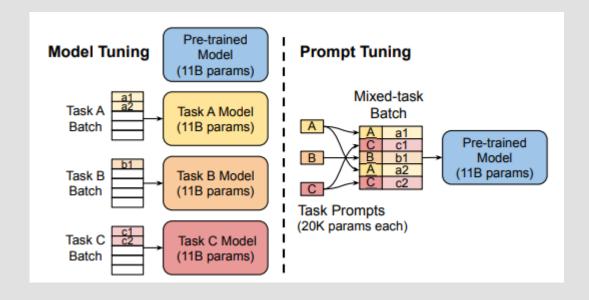






### Prompt tuning

#### Same idea as Prefix-tuning!



## Ways to initialize prompts



#### Random Uniform

Train from scratch, using random initialization

#### Sampled vocab

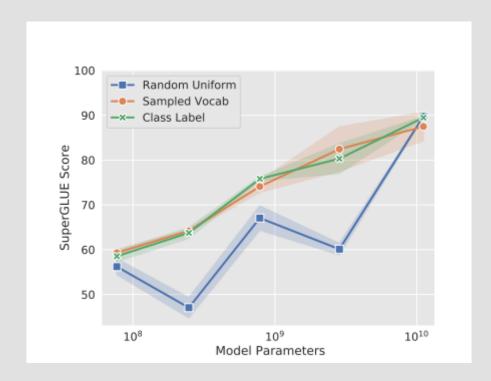
initialize each prompt token to an embedding drawn from the model's vocabulary.

#### Class Label

initialize the prompt with embeddings that enumerate the output classes

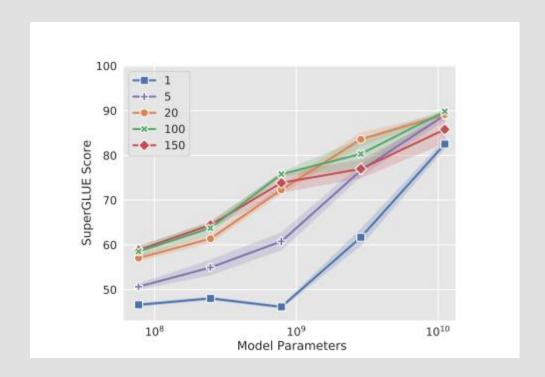
# Which one is the best strategy





# Ideal prompt length?

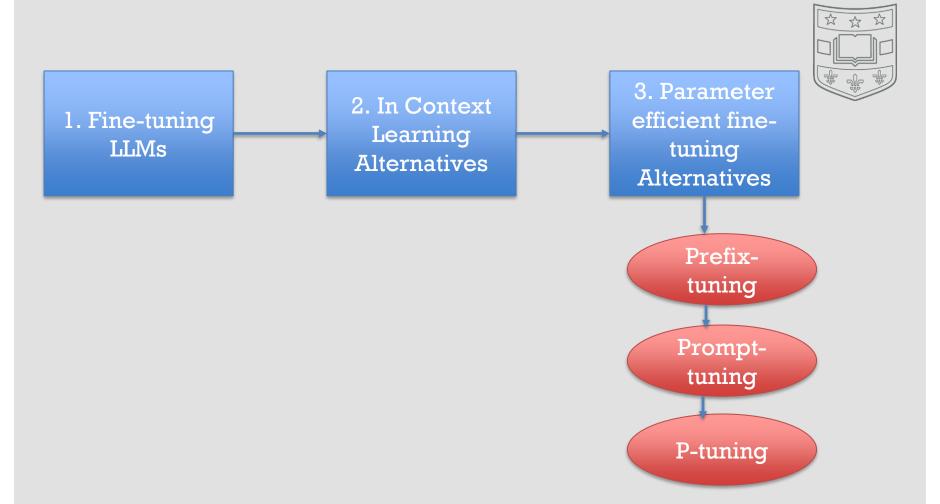




### Performance



Dataset	Domain	Model	Prompt	$\Delta$
SQuAD	Wiki	94.9 ±0.2	$94.8 \pm 0.1$	-0.1
TextbookQA	Book	54.3 ±3.7	<b>66.8</b> ±2.9	+12.5
BioASQ	Bio	$77.9 \pm 0.4$	<b>79.1</b> $\pm 0.3$	+1.2
RACE	Exam	59.8 ±0.6	<b>60.7</b> $\pm 0.5$	+0.9
RE	Wiki	$88.4 \pm 0.1$	<b>88.8</b> $\pm 0.2$	+0.4
DuoRC	Movie	<b>68.9</b> $\pm 0.7$	$67.7 \pm 1.1$	-1.2
DROP	Wiki	<b>68.9</b> ±1.7	$67.1 \pm 1.9$	-1.8







#### AI Open

Available online 26 August 2023

In Press, Journal Pre-proof  $\ \$  What's this?  $\ \$ 



#### GPT understands, too

Xiao Liu <sup>a 1</sup>, Yanan Zheng <sup>a 1</sup>, Zhengxiao Du <sup>a</sup>, Ming Ding <sup>a</sup>, Yujie Qian <sup>b</sup>, Zhilin Yang <sup>a</sup> ∠ ⋈, Jie Tang <sup>a</sup> ⋈

Show more V

+ Add to Mendeley 📽 Share 🗦 Cite

https://doi.org/10.1016/j.aiopen.2023.08.012 7

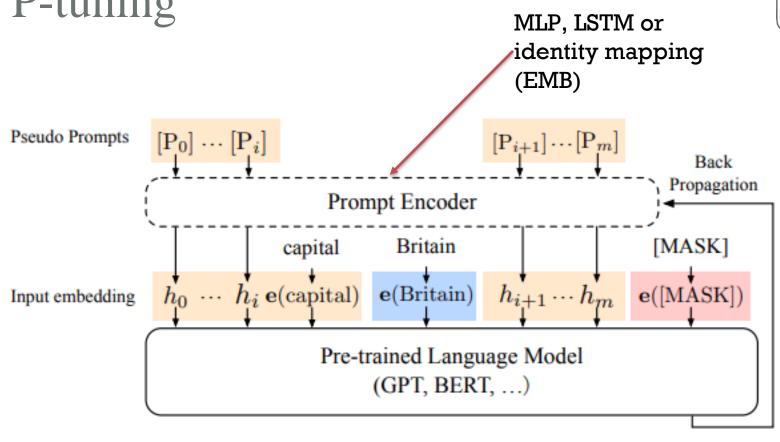
Get rights and content 7

Under a Creative Commons license →

open access







(b) P-tuning

# MLP, LSTM or Identity mapping (EMB)



Task	LSTM	MLP	EMB
WiC-ACC	56.27 ±1.54	$55.25 \pm 3.09$	53.96 ±3.23
CB-ACC. CB-F1.	81.70 ±7.49 77.41 ±9.15	88.39 ±3.72 84.24±5.15	82.59±3.69 67.27±6.78
BoolQ-ACC.	75.41±3.09	76.46±2.84	76.87±1.69

## Performances



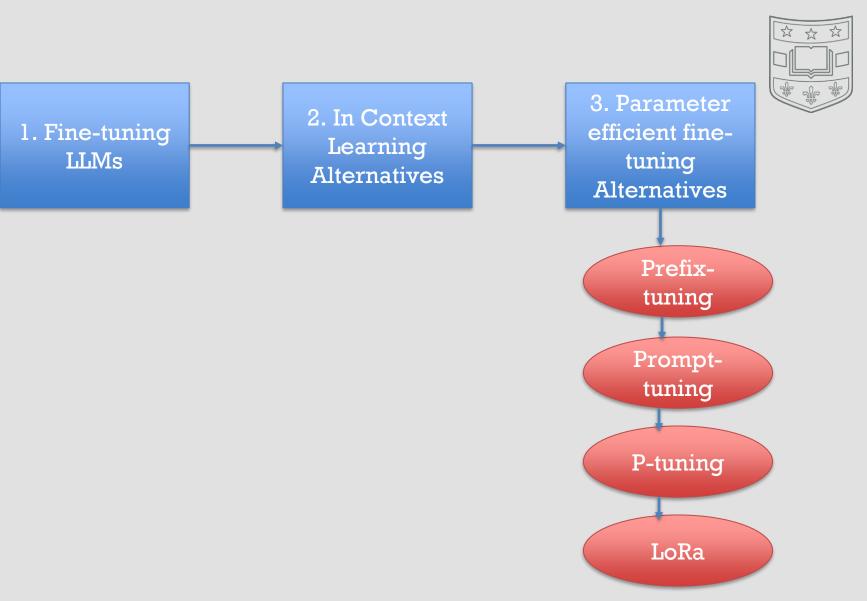
### (a) Fully-supervised performance with base-scale models.

Method		BoolQ CB (Acc.) (Acc.) (F1)			WiC	RTE	MultiRC		WSC COPA	Avg.	
	Method		(Acc.)	(F1)	(Acc.)	(Acc.)	(EM)	(F1a)	(Acc.)	(Acc.)	11.8.
BERT-Base	CLS-FT	72.9	85.1	73.9	71.1	68.4	16.2	66.3	63.5	67.0	66.2
(109M)	PET-FT	73.7	87.5	90.8	67.9	70.4	13.7	62.5	60.6	70.0	67.1
(10911)	P-tuning	73.9	89.2	92.1	68.8	71.1	14.8	63.3	63.5	72.0	68.4
GPT2-Base	CLS-FT	71.2	78.6	55.8	65.5	67.8	17.4	65.8	63.0	64.4	63.0
(117M)	PET-FT	74.8	87.5	88.1	68.0	70.0	23.5	69.7	66.3	78.0	70.2
(11/1/1)	P-tuning	75.0	91.1	93.2	68.3	70.8	23.5	69.8	63.5	76.0	70.4

#### (b) Fully-supervised performance with large-scale models.

Method		BoolQ	CB		WiC	RTE	MultiRC		WSC	COPA	Avg.
	Wicthod	(Acc.)	(Acc.)	(F1)	(Acc.)	(Acc.)	(EM)	(F1a)	(Acc.)	(Acc.)	Avg.
BERT-Large	CLS-FT <sup>1</sup>	77.7	94.6	93.7	74.9	75.8	24.7	70.5	68.3	69.0	72.5
(335M)	PET-FT	77.2	91.1	93.5	70.5	73.6	17.7	67.0	80.8	75.0	73.1
(333141)	P-tuning	77.8	96.4	97.4	72.7	75.5	17.1	65.6	81.7	76.0	74.6
GPT2-Med.	CLS-FT	71.0	73.2	51.2	65.2	72.2	19.2	65.8	62.5	66.0	63.1
(345M)	PET-FT	78.3	96.4	97.4	70.4	72.6	32.1	74.4	73.0	80.0	74.9
(343141)	P-tuning	78.9	98.2	98.7	69.4	75.5	29.3	74.2	74.0	81.0	75.6

.





#### Computer Science > Computation and Language

[Submitted on 17 Jun 2021 (v1), last revised 16 Oct 2021 (this version, v2)]

### LoRA: Low-Rank Adaptation of Large Language Models

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen

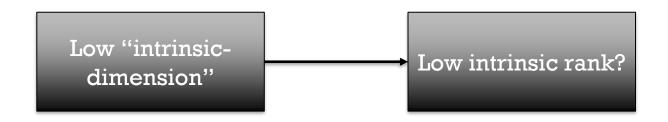
An important paradigm of natural language processing consists of large-scale pre-training on general domain data and adaptation feasible. Using GPT-3 175B as an example -- deploying independent instances of fine-tuned models, each with 175B parameters, trainable rank decomposition matrices into each layer of the Transformer architecture, greatly reducing the number of trainable par by 10,000 times and the GPU memory requirement by 3 times. LoRA performs on-par or better than fine-tuning in model quality on adapters, no additional inference latency. We also provide an empirical investigation into rank-deficiency in language model adapta provide our implementations and model checkpoints for Roberta, Deberta, and GPT-2 at this https URL.

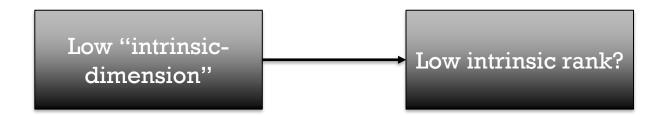
Submitted on 22 Dec 2020]

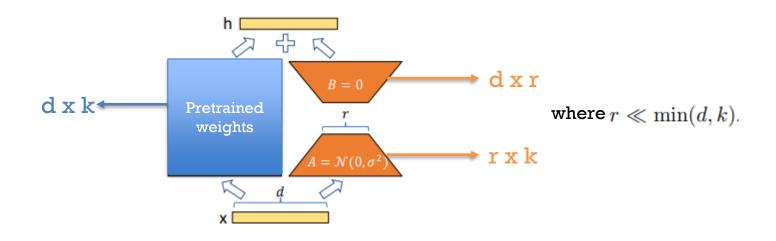
#### ntrinsic Dimensionality Explains the Effectiveness of Language Model Fine-Tuning

rmen Aghajanyan, Luke Zettlemoyer, Sonal Gupta

Although pretrained language models can be fine-tuned to produce state-of-the-art results for a very wide range of language understanding tasks, the dynamics of this process are not well understood, especially in the low data regime. Why can we use relatively vanilla gradient descent algorithms (e.g., without strong regularization) to tune a model with hundreds of millions of parameters on datasets with only hundreds or thousands of labeled examples? In this paper, we argue that analyzing fine-tuning through the lens of intrinsic dimension provides us with empirical and theoretical intuitions to explain this remarkable phenomenon. We empirically show that common pre-trained models have a very low intrinsic dimension; in other words, there exists a low dimension reparameterization that is as effective for fine-tuning as the full parameter space. For example, by optimizing only 200 trainable parameters randomly projected back into the full space, we can tune a RoBERTa model to achieve 90% of the full parameter performance levels on MRPC. Furthermore, we empirically show that pre-training implicitly minimizes intrinsic dimension and, perhaps surprisingly, larger models tend to have lower intrinsic dimension after a fixed number of pre-training updates, at least in part explaining their extreme effectiveness. Lastly, we connect intrinsic dimensionality with low dimensional task representations and compression based generalization bounds to provide intrinsic-dimension-based generalization bounds that are independent of the full parameter count.







$$\Delta W = \alpha AB$$
 Scaling factor 
$$h = W_0 x + \Delta W x = W_0 x + \alpha AB x$$

Model & Method	# Trainable Parameters		SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB <sub>base</sub> (FT)*	125.0M	87.6	94.8	90.2	63.6	92.8	91.9	78.7	91.2	86.4
RoB <sub>base</sub> (BitFit)*	0.1M	84.7	93.7	92.7	62.0	91.8	84.0	81.5	90.8	85.2
$RoB_{base}$ ( $Adpt^{D}$ )*	0.3M	87.1 <sub>±.0</sub>	$94.2 \scriptstyle{\pm .1}$	$88.5_{\pm 1.1}$	$60.8_{\pm.4}$	$93.1 \scriptstyle{\pm .1}$	$90.2 \scriptstyle{\pm .0}$	$71.5 {\scriptstyle \pm 2.7}$	$89.7_{\pm .3}$	84.4
$RoB_{base} (Adpt^{D})*$	0.9M	$87.3_{\pm .1}$	$94.7 \scriptstyle{\pm .3}$	$88.4_{\pm.1}$	$62.6_{\pm.9}$	$93.0_{\pm.2}$	$90.6 \scriptstyle{\pm .0}$	$75.9_{\pm 2.2}$	$90.3_{\pm .1}$	85.4
RoB <sub>base</sub> (LoRA)	0.3M	$87.5_{\pm .3}$	$\textbf{95.1}_{\pm .2}$	$89.7 \scriptstyle{\pm .7}$	$63.4 {\scriptstyle \pm 1.2}$	$93.3{\scriptstyle\pm.3}$	$90.8 \scriptstyle{\pm .1}$	$86.6_{\pm .7}$	$\textbf{91.5}_{\pm .2}$	87.2
RoB <sub>large</sub> (FT)*	355.0M	90.2	96.4	90.9	68.0	94.7	92.2	86.6	92.4	88.9
RoB <sub>large</sub> (LoRA)	0.8M	90.6 <sub>±.2</sub>	$96.2 \scriptstyle{\pm .5}$	$\textbf{90.9}_{\pm 1.2}$	$\textbf{68.2}_{\pm 1.9}$	$\textbf{94.9}_{\pm.3}$	$91.6 \scriptstyle{\pm .1}$	$\textbf{87.4}_{\pm 2.5}$	$\textbf{92.6}_{\pm .2}$	89.0
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	3.0M	90.2±.3	96.1 <sub>±.3</sub>	90.2 <sub>±.7</sub>	<b>68.3</b> ±1.0	94.8 <sub>±.2</sub>	<b>91.9</b> <sub>±.1</sub>	83.8 <sub>±2.9</sub>	92.1 <sub>±.7</sub>	88.4
$RoB_{large} (Adpt^{P})^{\dagger}$	0.8M	90.5 <sub>±.3</sub>	$\textbf{96.6}_{\pm.2}$	$89.7_{\pm 1.2}$	$67.8_{\pm 2.5}$	$\textbf{94.8}_{\pm.3}$	$91.7 \scriptstyle{\pm .2}$	$80.1_{\pm 2.9}$	$91.9_{\pm .4}$	87.9
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	6.0M	89.9 <sub>±.5</sub>	$96.2 \scriptstyle{\pm .3}$	$88.7_{\pm 2.9}$	$66.5_{\pm 4.4}$	$94.7 \scriptstyle{\pm .2}$	$92.1_{\pm .1}$	$83.4_{\pm 1.1}$	$91.0{\scriptstyle\pm1.7}$	87.8
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	0.8M	90.3 <sub>±.3</sub>	$96.3_{\pm .5}$	$87.7_{\pm 1.7}$	$66.3_{\pm 2.0}$	$94.7_{\pm .2}$	$91.5{\scriptstyle \pm .1}$	$72.9_{\pm 2.9}$	$91.5_{\pm .5}$	86.4
RoB <sub>large</sub> (LoRA)†	0.8M	90.6 <sub>±.2</sub>	$96.2 \scriptstyle{\pm .5}$	$\textbf{90.2}_{\pm 1.0}$	$68.2{\scriptstyle\pm1.9}$	$\textbf{94.8}_{\pm.3}$	$91.6 \scriptstyle{\pm .2}$	$\textbf{85.2}_{\pm 1.1}$	$\textbf{92.3}_{\pm.5}$	88.6
DeB <sub>XXL</sub> (FT)*	1500.0M	91.8	97.2	92.0	72.0	96.0	92.7	93.9	92.9	91.1
DeB <sub>XXL</sub> (LoRA)	4.7M	$91.9_{\pm .2}$	$96.9_{\pm.2}$	$\textbf{92.6}_{\pm.6}$	$\textbf{72.4}_{\pm 1.1}$	$\textbf{96.0}_{\pm.1}$	$\textbf{92.9}_{\pm.1}$	$\textbf{94.9}_{\pm.4}$	$\textbf{93.0}_{\pm.2}$	91.3

## What is the right *r*? Which weights should LoRa apply to?

	Weight Type	r = 1	r = 2	r = 4	r = 8	r = 64
WikiSQL(±0.5%)	$\left \begin{array}{c} W_q \\ W_q, W_v \\ W_q, W_k, W_v, W_o \end{array}\right $	68.8 73.4 74.1	69.6 73.3 73.7	70.5 73.7 74.0	70.4 73.8 74.0	70.0 73.5 73.9
MultiNLI ( $\pm 0.1\%$ )	$\left \begin{array}{c} W_q \\ W_q, W_v \\ W_q, W_k, W_v, W_o \end{array}\right $	90.7 91.3 91.2	90.9 91.4 91.7	91.1 91.3 91.7	90.7 91.6 91.5	90.7 91.4 91.4

Idea: with a small r even just on  $\mathcal{W}_q$  and  $\mathcal{W}_v$  , you can achieve great performances

# How much more computational resources conserved?



- For GPT3-175B
  - VRAM consumption: 1.2TB --> 350GB
  - 25% speed up in time



NeurIPS Proceedings 🔹 🖼

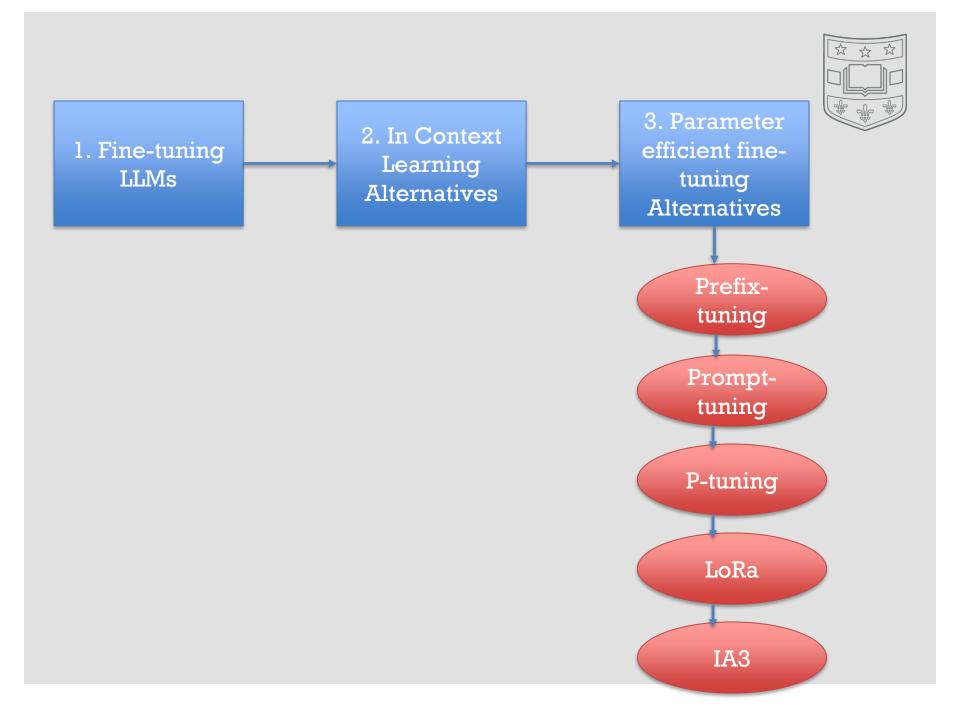
### Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning

Part of Advances in Neural Information Processing Systems 35 (NeurIPS 2022) Main Conference Track

Bibtex Paper Supplemental

#### **Authors**

Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, Colin A. Raffel



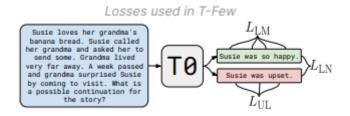
### Introduce 2 losses

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

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

 Length normalized Loss^

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



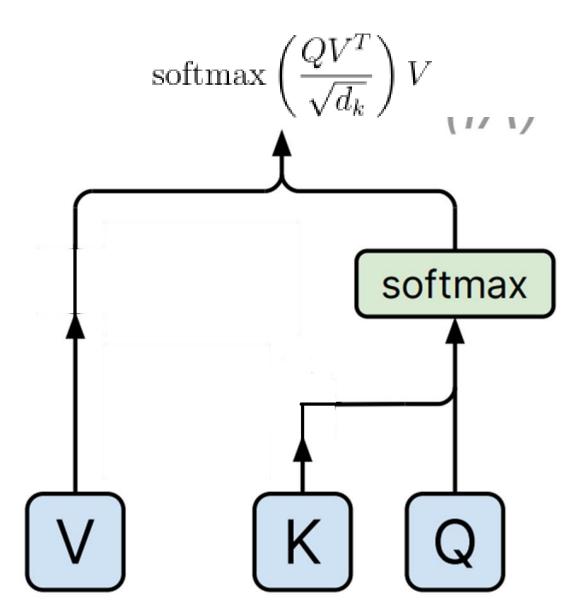
^ basically the log-probability relative to other tokens

## Introduce IA3 training

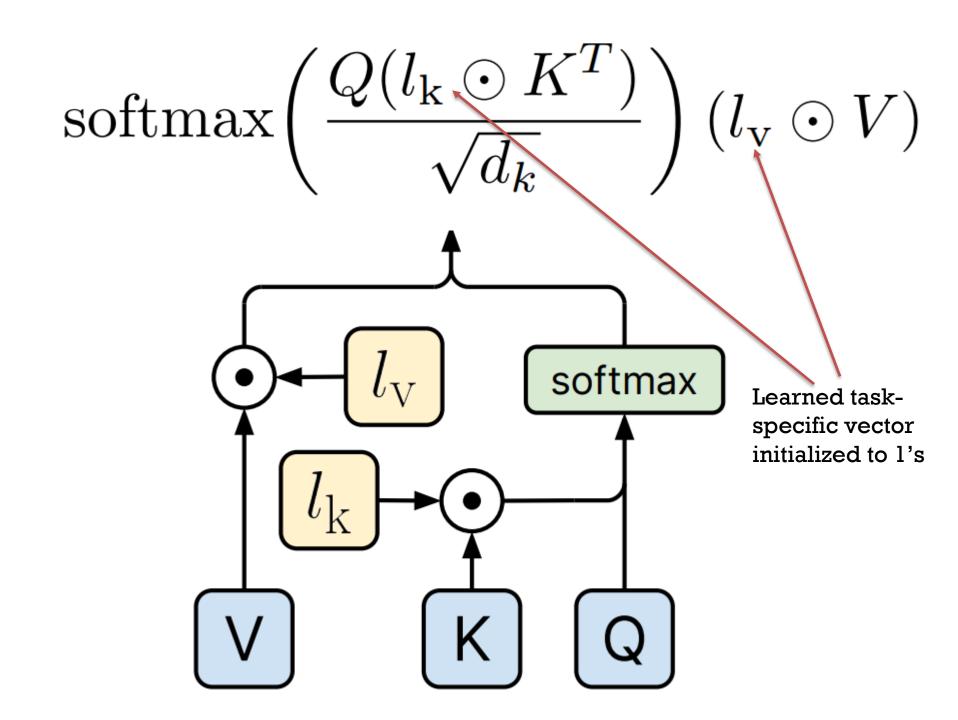


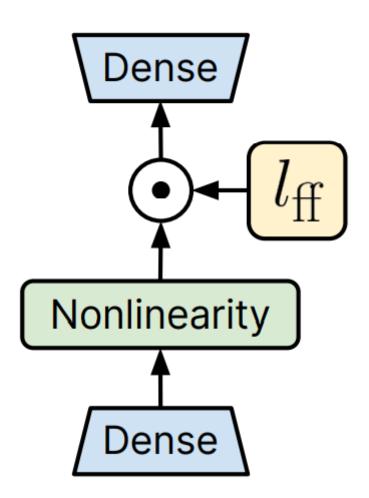
- Short for "Infused Adapter by Inhibiting and Amplifying Inner Activations"
- Introduce 3 learned vectors





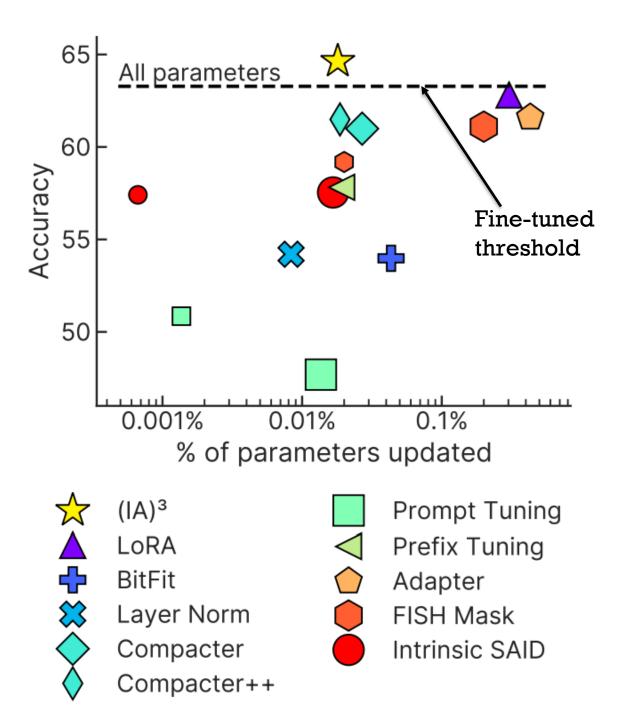
$$\operatorname{softmax}\left(\frac{Q(l_{\mathbf{k}}\odot K^T)}{\sqrt{d_{k}}}\right)(l_{\mathbf{V}}\odot V)$$





# Number of parameters required decreases

From the encoder	From the decoder
$L(d_k + d_v + \bar{d}_{\mathrm{ff}})$	$L(2d_k + 2d_v + d_{\rm ff})$



# Performs better than humans as well with few-shot learning!

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.

## Cheaper as well for prompting

