

# When Do Prompting and Prefix-Tuning Work?

## A Theory of Capabilities and Limitations

Aleksandar Petrov, Philip Torr, Adel Bibi

Department of Engineering Science, University of Oxford, United Kingdom

24.2.29

Minkyu Kim

EffL@POSTECH

# Index

- **Introduction**
- **Background**
- **Analyze & Delineate**
  - Soft prompting has more capacity than prompting
  - Prefix-tuning can only bias the output of an attention head
  - The bias can elicit skills from the pretrained model
  - Effects of prefix-tuning beyond the single attention layer
- **Limitation**
- **Conclusion**

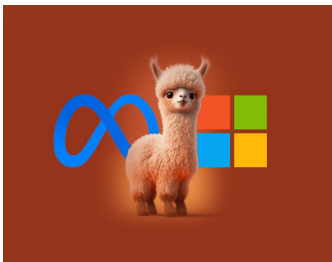
# Introduction

Training the cutting-edge language model. needs **huge cost**

- GPT4. more than \$100M
- LLaMA-2-7b. \$85K

**Note.** Without cost to buy hardware (~ \$30M)

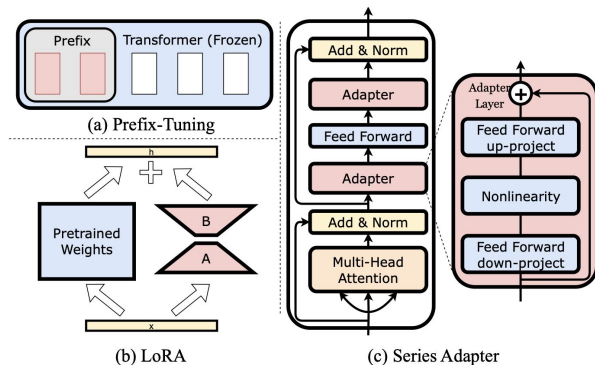
👉 Out of reach for most academic researchers (even fine-tuning)



# Introduction

## Alternative. ‘Efficient’ fine-tuning

- Sparsely modifying the parameters of the model
  - Adapter modules, Low-rank updates
- **Modifying its input context** (Context-based fine-tuning)
  - In-context learning (GPT3)
  - Prompt tuning(Hard prompt tuning, Soft prompt tuning, Prefix-tuning)



# Introduction

## Modifying its input context (Context-based fine-tuning).

- The most popular approach: prompting
  - Generation is conditioned on either **human-crafted** or **automatically optimized** tokens
  - Have witnessed impressive empirical successes and widespread adoption
- Type
  - Tokens(word sequence): **Hard** Prompt 👉 discrete information
  - Vectors(parameters): **Soft** Prompt 👉 continuous information
    - greater expressiveness due to the **expansive nature of continuous space**

# Introduction

However, we have little **theoretical understanding** of how they work

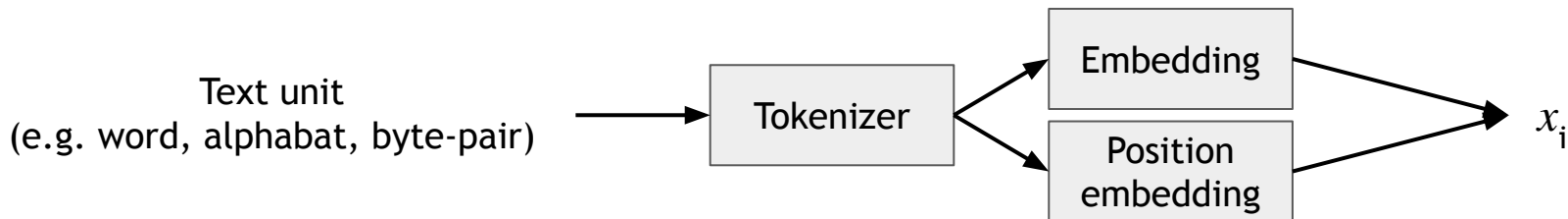
- In this work, we...
  - Analyse the influence of prompts and prefixes on the computations of a pretrained model
  - Delineate their limitations
- Specifically, we address the following questions:
  - Can a transformer utilize the soft prompt or prefix vector?
  - Since prefix-tuning is more expressive than prompting, is it as expressive as full fine-tuning?
  - If context-based fine-tuning methods suffer from such structural limitations, how come they have high empirical performance?

(Authors answered these in next sections)

# Background

## Input data.

- **Input sequences.**  $(x_1, \dots, x_p)$ ,  $x_i \in \{1, \dots, V\}$ 
  - $V$ : size of vocabulary
- **Embedding matrix (token  $\rightarrow$  vector).**  $E \in \mathbb{R}^{d_e \times V}$ 
  - $d_e$ : dimension of embedding
- **One-hot  $i$ -th position encoding.**  $e_N(i)$
- **Embedding for  $i$ -th position.**  $x_i = [E_{:,x_i}^T, e_N^T(i)]^T$



# Background

**Transformer architecture.** Consists of alternating attention blocks

- **Attention block.** consists of H heads
  - Attention matrix for head 'h':  $\mathbf{A}^h \in \mathbb{R}^{p \times p}$
- **Linear.** applied to all positions of the sequence
  - followed by ReLU activation layer

$$\mathbf{A}_{ij}^h = \frac{\exp\left(T/\sqrt{k}(\mathbf{W}_Q^h \mathbf{x}_i)^\top (\mathbf{W}_K^h \mathbf{x}_j)\right)}{\sum_{r=1}^p \exp\left(T/\sqrt{k}(\mathbf{W}_Q^h \mathbf{x}_i)^\top (\mathbf{W}_K^h \mathbf{x}_r)\right)}$$

[Attention matrix]

$$\mathcal{A}[(\mathbf{W}_Q^1, \dots, \mathbf{W}_Q^H), (\mathbf{W}_K^1, \dots, \mathbf{W}_K^H), (\mathbf{W}_V^1, \dots, \mathbf{W}_V^H)](\mathbf{x}_1, \dots, \mathbf{x}_p) = (\mathbf{t}_1, \dots, \mathbf{t}_p)$$
$$\mathbf{t}_i = \sum_{h=1}^H \sum_{j=1}^p \mathbf{A}_{ij}^h \mathbf{W}_V^h \mathbf{x}_j.$$

[Attention block]

$$\mathcal{L}[\mathbf{M}, \mathbf{b}](\mathbf{x}) = \mathbf{M}\mathbf{x} + \mathbf{b}$$
$$\hat{\mathcal{L}}[\mathbf{M}, \mathbf{b}](\mathbf{x}) = \text{ReLU}(\mathbf{M}\mathbf{x} + \mathbf{b})$$

[Linear layer]

$$(\mathbf{y}_1, \dots, \mathbf{y}_p) = \left( \mathcal{A}_1 \circ \hat{\mathcal{L}}_{1,1} \circ \mathcal{L}_{1,2} \circ \mathcal{A}_2 \circ \hat{\mathcal{L}}_{2,1} \circ \mathcal{L}_{2,2} \circ \text{softmax} \right) \left( \left[ \begin{matrix} \mathbf{E}_{:,x_1} \\ \mathbf{e}_N(1) \end{matrix} \right], \dots, \left[ \begin{matrix} \mathbf{E}_{:,x_p} \\ \mathbf{e}_N(p) \end{matrix} \right] \right)$$

[Output of the transformer]

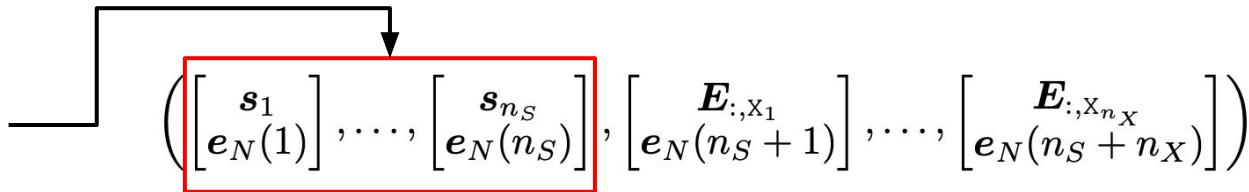


# Background

## Context-based fine-tuning of a pretrained model.

- **Hard Prompting.** prefixing the input with a token sequence
  - Guide the model response
- **Soft prompting.** learnable input embeddings

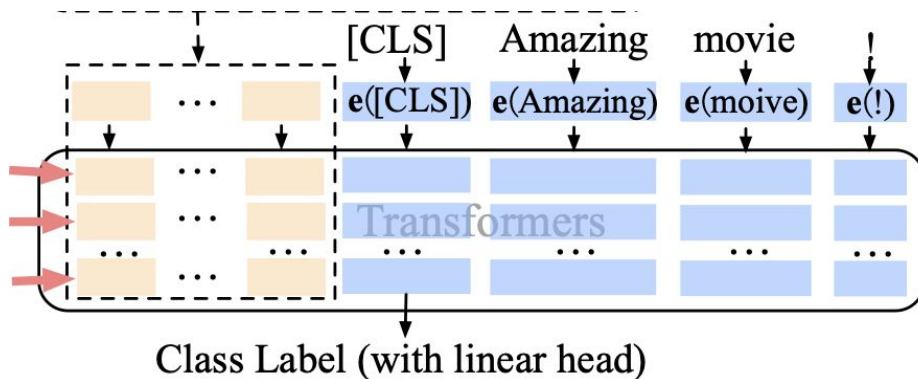
Constant: **hard** prompting  
Learnable: **soft** prompting



# Background

## Context-based fine-tuning of a pretrained model.

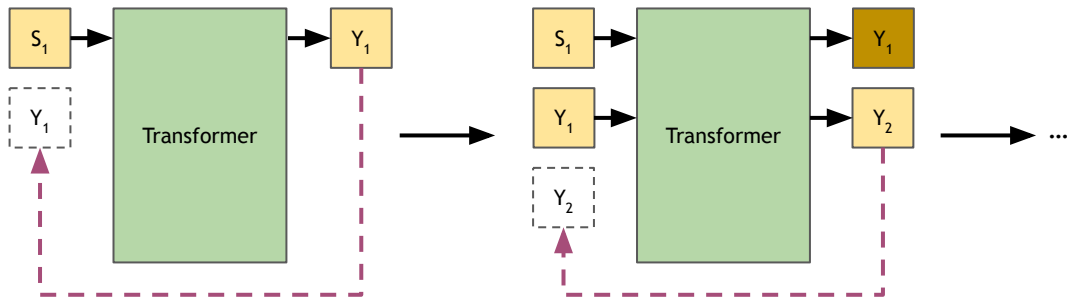
- **Prefix tuning.** applied soft prompt across the depth of the model
  - The first  $n_s$  positions for all attention blocks are learnable parameters
    - Optimize the inputs of every attention layer
  - Successful at fine-tuning models



# Soft prompting has more capacity than prompting

Case: **Unconditional** generation with a single system token

- Vocabulary size:  $V$
- Deterministic autoregressive function
- If we use **hard** prompt. Enable to get number of  $V$  outputs
  - $[f(s_1), f(s_2), \dots, f(s_v)]$  (**Finite** number('V') of case)
- If we use **soft** prompt.
  - Real-value  $\rightarrow$  **Infinite** number of cases



$$(Y_1, \dots, Y_N) = f(S_1) = f_{S_1}$$

[using **hard** prompt]

$$(Y_1, \dots, Y_N) = f(s_1) = f_{s_1}$$

[using **soft** prompt]

# Soft prompting has more capacity than prompting

Guessing the successful result of soft prompt (comparing with hard prompt)

- **Unconditional** generation with a single system token
  - If we use **soft** prompt.
    - Note: it does not mean that the transformer architecture can represent a function that achieves any type of output **in practice**
      - it is not obvious if there is a surjective map from **soft** prompt to **word** in vocab.
  - **Theorem 1.** Condition that enable to create ‘soft prompt + transformer’  
, which generates infinite unique outputs

**Theorem 1** (Exponential unconditional generation capacity of a single virtual token). *For any  $V, N > 0$ , there exists a transformer with vocabulary size  $V$ , context size  $N$ , embedding size  $d_e = N$ , one attention layer with two heads and a three-layer MLP such that it generates any token sequence  $(Y_1, \dots, Y_N) \in \{1, \dots, V\}^N$  when conditioned on the single virtual token  $\mathbf{s}_1 = ((Y_1 - 1)/V, \dots, (Y_N - 1)/V)$ .*

# Soft prompting has more capacity than prompting

Guessing the successful result of soft prompt (comparing with hard prompt)

- **Conditional** generation with a single system token

- Prompt + User input ( $X_1, \dots, X_{n_x}$ ), Output ( $Y_1, \dots, Y_{n_y}$ )
- If we use **hard** prompt.
  - $X_1 \rightarrow Y_1$ :  $V^V$  cases ( $1 \rightarrow [1, \dots, V], 2 \rightarrow [1, \dots, V], \dots, V \rightarrow [1, \dots, V]$  👉  $[1, \dots, V] \rightarrow [1, \dots, V]$ )
    - $S_1 \in \{1, \dots, V\}$ :  $V$  cases = less than  $V^V$ 
      - $\rightarrow$  Can't control the output by modifying system token  
( $S_1$  cannot be used to specify an arbitrary map from input to output)
- If we use **soft** prompt.
  - $S_1$  is 'real value' = infinite cases
    - Can control the all cases of output( $V^V$ )

$$(Y_1 = f(\boxed{S_1, X_1}) = f_{S_1}(X_1))$$

# Soft prompting has more capacity than prompting

Guessing the successful result of soft prompt (comparing with hard prompt)

- **Conditional** generation with a single system token
  - **Theorem 2.** **Soft** prompting is also more expressive for the conditional behavior of a transformer model

**Theorem 2** (Conditional generation capacity for a single virtual token ( $n_X=n_Y=1$ )). *For any  $V>0$ , there exists a transformer with vocabulary size  $V$ , context size  $N=2$ , embedding size  $d_e=V$ , one attention layer with two heads and a three-layer MLP that reproduces any map  $m:[1, \dots, V] \rightarrow [1, \dots, V]$  from a user input token to a model response token when conditioned on a single virtual token  $\mathbf{s}_1=(m^{(1)}/V, \dots, m^{(V)}/V)$ . That is, by selecting  $\mathbf{s}_1$  we control the model response to any user input.*

# Soft prompting has more capacity than prompting

## Conclusion.

- Soft prompting, Prefix-tuning possesses **greater expressiveness** than prompting
- We can fully determine the map from user input to output using virtual tokens
  - Soft prompting is **as powerful as** full fine-tuning

# Prefix-tuning can only bias the attention head's output

Strength of soft prompting & prefix-tuning.

- Needs specific network's parameters (Theorem 1, 2)

So, we can ask question,

- Given a fixed pre-trained model, Can prefix-tuning be considered equally powerful to full fine-tuning?
  - Answer. prefix **can't change** the relative attention over the content  $X, Y$ 
    - Can **only bias** the attention block outputs in a subspace of rank  $n_s$ 
      - 👉 Making it strictly less powerful than full fine-tuning



# Prefix-tuning can only bias the attention head's output

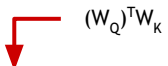
Prefix-tuning **can't alter the attention pattern** of an attention head

- $A_{ij}$ : Attention that i-th position token gives to j-th position token

$$A_{ij} = \frac{\exp\left(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_j\right)}{\sum_{r=1}^p \exp\left(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_r\right)} = \frac{\exp\left(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{H} \mathbf{x}_j\right)}{\sum_{r=1}^p \exp\left(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{H} \mathbf{x}_r\right)}$$

- **Full fine-tuning.** Can enact arbitrary changes to  $\mathbf{W}_Q$  and  $\mathbf{W}_K$ 
  - Change the attention patterns arbitrarily

$$A_{ij}^{\text{ft}} = \frac{\exp\left(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{H} \mathbf{x}_j + T/\sqrt{k} \mathbf{x}_i^\top \Delta \mathbf{H} \mathbf{x}_j\right)}{\sum_{r=1}^p \exp\left(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{H} \mathbf{x}_r + T/\sqrt{k} \mathbf{x}_i^\top \Delta \mathbf{H} \mathbf{x}_r\right)}$$

  $(\mathbf{W}_Q)^\top \mathbf{W}_K$

# Prefix-tuning can only bias the attention head's output

Prefix-tuning **can't alter the attention pattern** of an attention head

- **Prefix-tuning.** If we have a prefix of length one at position 0,
  - Only increase denominator of attention where to get attention not at 0  
👉 prefix can't change the attention pattern (can't modify what an attention head attends to)
  - Re-writing the  $A_{ij}^{\text{pt}}$  makes this more evident

$$A_{ij}^{\text{pt}} = A_{ij} \sum_{r=1}^p A_{ir}^{\text{pt}} = A_{ij} (1 - A_{i0}^{\text{pt}}) \leftarrow \text{just scaling}$$

$$A_{i0}^{\text{pt}} = \frac{\exp\left(\frac{T}{\sqrt{k}} \mathbf{x}_i^\top \mathbf{H} \mathbf{s}_1\right)}{\exp\left(\frac{T}{\sqrt{k}} \mathbf{x}_i^\top \mathbf{H} \mathbf{s}_1\right) + \sum_{r=1}^p \exp\left(\frac{T}{\sqrt{k}} \mathbf{x}_i^\top \mathbf{H} \mathbf{x}_r\right)}, \quad A_{ij}^{\text{pt}} = \frac{\exp\left(\frac{T}{\sqrt{k}} \mathbf{x}_i^\top \mathbf{H} \mathbf{x}_j\right)}{\exp\left(\frac{T}{\sqrt{k}} \mathbf{x}_i^\top \mathbf{H} \mathbf{s}_1\right) + \sum_{r=1}^p \exp\left(\frac{T}{\sqrt{k}} \mathbf{x}_i^\top \mathbf{H} \mathbf{x}_r\right)} \quad \text{for } j \geq 1.$$

Attention to prefix      Attention to other tokens

(Added value to denominator)

# Prefix-tuning can only bias the attention head's output

Prefix-tuning **only adds a bias** to the attention block output

- **Term containing prefix.** Independent of the content ( $x_1, \dots, x_p$ )
  - Act as bias
- In contrast with full fine-tuning
  - Allow for **content-dependent change** of attention and value computation.

$$t_i = \sum_{j=1}^p A_{ij} W_V x_j, \quad t_i^{\text{ft}} = \sum_{j=1}^p A_{ij}^{\text{ft}} (W_V + \boxed{\Delta W_V}) x_j,$$

$$t_i^{\text{pt}} = A_{i0}^{\text{pt}} W_V s_1 + \sum_{j=1}^p A_{ij}^{\text{pt}} W_V x_j \stackrel{(6)}{=} A_{i0}^{\text{pt}} W_V s_1 + \sum_{j=1}^p A_{ij} (1 - A_{i0}^{\text{pt}}) W_V x_j = \boxed{A_{i0}^{\text{pt}} W_V s_1} + (1 - A_{i0}^{\text{pt}}) t_i.$$

# Prefix-tuning can only bias the attention head's output

## Conclusion,

- In real-world scenario, prefix-tuning has some limitations,
  - It can be competitive with full fine-tuning **only in very limited circumstances**
  - Transformers do not behave like we want
    - Models are typically trained with token inputs rather than soft prompts

# The bias can elicit skills from the pre-trained model

**Target.** Explain when and why prefix-tuning can work in practice.

- **Hypothesis.** prefix-tuning cannot be used to gain **new knowledge** but can bring to the surface **latent knowledge present in the pretrained model**
- **Test.** Constructing small transformers trained on one or few tasks
  - Task: Sort numbers into ascending/descending order
    - Training strategy 1: Pre-training task = Fine-tuning task
    - Training strategy 2: Pre-training task != Fine-tuning task

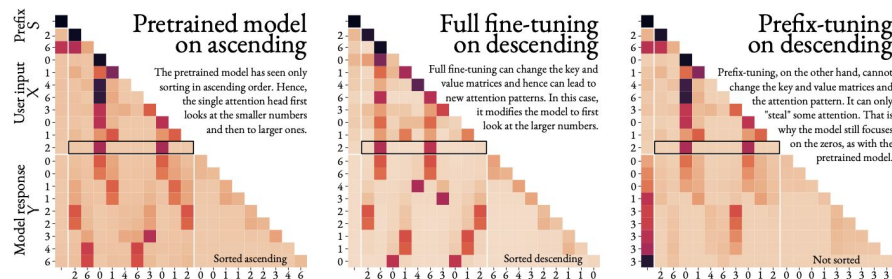
# The bias can elicit skills from the pre-trained model


Test results. Prefix-tuning **can't learn a new task** requiring a different attention pattern

- 1-layer, 1-head transformer

	Ascending	Descending
Pretrain on asc.	$91 \pm 5\%$	$0 \pm 0\%$
Full fine-tune on desc.	$0 \pm 0\%$	$85 \pm 5\%$
Prefix-tune on desc.	$0 \pm 0\%$	$0 \pm 0\%$

Accuracy (10 random seeds)



: the attention when the first response  $Y_1$  is being generated

# The bias can elicit skills from the pre-trained model

**Test results.** Prefix-tuning **can elicit a skill** from the pre-trained model

- Pretrain a 1-layer, 4-head model with solutions...

- Sorted in ascending ( $\nearrow$ ) order (1st task)
- Sorted in descending ( $\searrow$ ) order (2nd task)
- Adding one (+1) to each element (3th task)
- Adding two (+2) to each element (4th task)

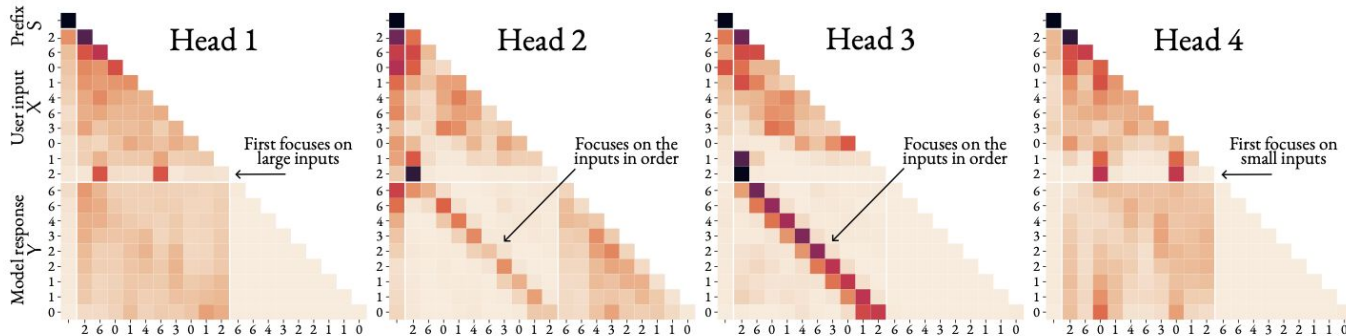
 (Each head learns a different task)

# The bias can elicit skills from the pre-trained model

Test results. Prefix-tuning **can elicit a skill** from the pre-trained model

- Compared to the previous case, prefix-tuning is more successful
  - because the pretrained model **contains the attention mechanisms for solving the four tasks**

Accuracy on:	$\nearrow$	$\searrow$	+1	+2
Pretrained	$25 \pm 13\%$	$25 \pm 12\%$	$24 \pm 11\%$	$22 \pm 7\%$
Prefix-tune on $\nearrow$	$95 \pm 2\%$	$0 \pm 0\%$	$0 \pm 0\%$	$0 \pm 0\%$
Prefix-tune on $\searrow$	$0 \pm 0\%$	$90 \pm 3\%$	$1 \pm 1\%$	$1 \pm 1\%$
Prefix-tune on +1	$0 \pm 0\%$	$1 \pm 3\%$	$95 \pm 6\%$	$0 \pm 1\%$
Prefix-tune on +2	$0 \pm 0\%$	$0 \pm 0\%$	$1 \pm 2\%$	$98 \pm 5\%$





# The bias can elicit skills from the pre-trained model

**Conclusion.** Prefix-tuning **can combine knowledge** from pre-training tasks to solve new tasks

- “Skill” required to solve the new task = “Skills” the pre-trained model has seen
- **Test again.** 40-layer 4-head model with the same four tasks
  - length of prefixes: 10 ( $n_s = 10$ )
  - New task H: mapping each element to the number of elements in the sequence with the same value

Accuracy on:	↗	↘	+1	+2	↗+1	H
Pretrained	17%	23%	34%	25%	0%	0%
Prefix-tune on ↗	100%	0%	0%	0%	0%	0%
Prefix-tune on ↘	0%	100%	0%	0%	0%	0%
Prefix-tune on +1	0%	0%	100%	0%	0%	0%
Prefix-tune on +2	0%	0%	0%	100%	0%	0%
Prefix-tune on ↗+1	0%	0%	0%	0%	93%	0%
Prefix-tune on H	0%	0%	0%	0%	0%	1%

Combine the knowledge

# Effects of prefix-tuning beyond the single attention layer

**Multi layer case.** Bias made by prefix can exhibit increasingly complex behaviors

- **Strength.**
  - Prefix-tuning can change the attention pattern
- **Weakness.**
  - The representational capacity of the prefix-tuning is severely limited



Prefix-tuning appears to be less expressive than full fine-tuning  
even when both methods are given the same number of learnable parameters

# Effects of prefix-tuning beyond the single attention layer

## Prefix-tuning can change the attention

- Output of first attention(including bias): affect the attention with content
- **Limitation.** Each depends on one input position(i or j) only
  - Full fine-tuning can achieve

$$\tilde{\mathbf{A}}_{ij}^{\text{ft}(2)} = T/\sqrt{k} \mathbf{t}_i^{\text{ft}(1)\top} (\mathbf{H}^{(2)} + \Delta \mathbf{H}^{(2)}) \mathbf{t}_j^{\text{ft}(1)}$$

$$\mathbf{t}_i^{\text{pt}(1)} = \mathbf{A}_{i0}^{\text{pt}(1)} \mathbf{W}_V \mathbf{s}_1^{(1)} + \sum_{j=1}^p \mathbf{A}_{ij}^{\text{pt}(1)} \mathbf{W}_V^{(1)} \mathbf{x}_j^{(1)} \stackrel{(7)}{=} \underbrace{\mathbf{A}_{i0}^{\text{pt}(1)}}_{\alpha_i} \underbrace{\mathbf{W}_V \mathbf{s}_1^{(1)}}_{\boldsymbol{\mu}} + (1 - \mathbf{A}_{i0}^{\text{pt}(1)}) \mathbf{t}_i^{(1)},$$

$$\begin{aligned} \tilde{\mathbf{A}}_{ij}^{\text{pt}(2)} &= \frac{T}{\sqrt{k}} \mathbf{t}_i^{\text{pt}(1)\top} \mathbf{H}^{(2)} \mathbf{t}_j^{\text{pt}(1)}, \\ &= \frac{T}{\sqrt{k}} (\alpha_i \alpha_j \underbrace{\boldsymbol{\mu}^\top \mathbf{H}^{(2)} \boldsymbol{\mu}}_{\text{constant}} + \alpha_j (1 - \alpha_i) \underbrace{\mathbf{t}_i^{(1)\top} \mathbf{H}^{(2)} \boldsymbol{\mu}}_{\text{depends only on } \mathbf{t}_i^{(1)}} + \alpha_i (1 - \alpha_j) \underbrace{\boldsymbol{\mu}^\top \mathbf{H}^{(2)} \mathbf{t}_j^{(1)}}_{\text{depends only on } \mathbf{t}_j^{(1)}} + (1 - \alpha_i)(1 - \alpha_j) \underbrace{\mathbf{t}_i^{(1)\top} \mathbf{H}^{(2)} \mathbf{t}_j^{(1)}}_{\text{pretrained attention } \tilde{\mathbf{A}}_{ij}^{(2)}). \end{aligned}$$

# Effects of prefix-tuning beyond the single attention layer

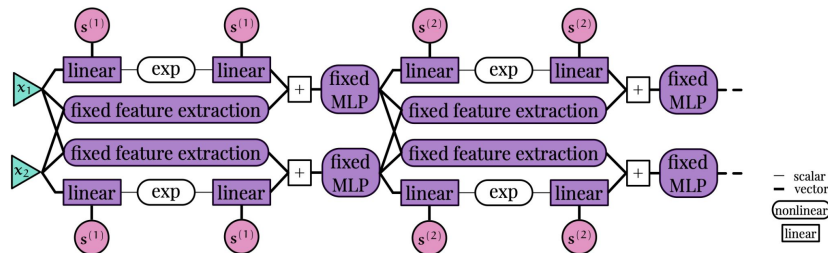
The representational capacity of the prefix-tuning is severely limited

- Consider two inputs  $x_1, x_2$  and a single prefix  $s$  (learnable parameters)
  - $H, W$ : pre-trained parameters
  - prefix interacts with only one of the inputs at a time and only by computing a single scalar value
- Computational graph.
  - Only learnable interaction between the inputs is **indirect**
    - Interaction between the inputs happens via **non-learnable residual connections**

$$\mathcal{A}_s(x_1, x_2) = \langle y_1, y_2 \rangle$$

$$y_1 = \frac{\exp(x_1^\top H s) W_V s + \exp(x_1^\top H x_1) W_V x_1 + \exp(x_1^\top H x_2) W_V x_2}{C_1}$$
$$y_2 = \frac{\exp(x_2^\top H s) W_V s + \exp(x_2^\top H x_1) W_V x_1 + \exp(x_2^\top H x_2) W_V x_2}{C_2}$$

learnable parameter participates



Computational graph

# Limitation

Not analyze about Suffix-tuning.

- Representations for prompt and soft prompt suffixes would depend on the previous positions
  - Whether suffixing is more expressive than prefixing remains an open question

Open question.

- Formal analysis of the conditions under which they may be universal approximators

Training scale. only conducting toy experiments

- In practice, language models are pre-trained with very large datasets
  - can pick up very complex behaviors
  - the limitations when using large-scale pretrained transformers with soft prompt = future work

# Conclusion.

- Soft prompting are strictly more expressive than hard prompting
- Despite this expressivity, prompt tuning suffers from structural limitations
  - Have difficult learning new attention patterns
  - Only bias the output of the attention layer
- Prefix-tuning can easily elicit a skill the pretrained model already seen
  - can even learn a new task
- The effect of the prefix-induced bias is more complicated and powerful
  - still strictly less expressive than full fine-tuning