

WHEN DO PROMPTING AND PREFIX-TUNING WORK? A THEORY OF CAPABILITIES AND LIMITATIONS

Aleksandar Petrov, Philip H.S. Torr & Adel Bibi

Department of Engineering Science

University of Oxford

Oxford, United Kingdom

{aleksandar.petrov, philip.torr, adel.bibi}@eng.ox.ac.uk

ABSTRACT

Context-based fine-tuning methods, including prompting, in-context learning, soft prompting (also known as prompt tuning), and prefix-tuning, have gained popularity due to their ability to often match the performance of full fine-tuning with a fraction of the parameters. Despite their empirical successes, there is little theoretical understanding of how these techniques influence the internal computation of the model and their expressiveness limitations. We show that despite the continuous embedding space being more expressive than the discrete token space, soft-prompting and prefix-tuning are strictly less expressive than full fine-tuning, even with the same number of learnable parameters. Concretely, context-based fine-tuning cannot change the relative attention pattern over the content and can only bias the outputs of an attention layer in a fixed direction. This suggests that while techniques like prompting, in-context learning, soft prompting, and prefix-tuning can effectively elicit skills present in the pretrained model, they cannot learn novel tasks that require new attention patterns.

1 INTRODUCTION

Language model advances are largely driven by larger models and more training data (Kaplan et al., 2020; Rae et al., 2021). Training cutting-edge models is out of reach for most academic researchers, small enterprises, and individuals, and it has become common to instead fine-tune open-source pretrained models (Devlin et al., 2019; Min et al., 2021). Yet, due to escalating computational demands, even fine-tuning of the larger models has become prohibitively expensive (Lialin et al., 2023).

As a result, there is an acute need for more efficient fine-tuning methods, either by sparsely modifying the parameters of the model or modifying its input context. Examples of the first type include adapter modules which introduce a few trainable layers to modify the behaviour of the frozen pretrained network (Rebuffi et al., 2017; Houlsby et al., 2019; Hu et al., 2023). One can also use low-rank updates, which also results in a reduced number of trainable parameters (Hu et al., 2021).

Context-based fine-tuning has been motivated by the success of few-shot and zero-shot learning (Wei et al., 2021; Kojima et al., 2022). The most popular context-based approach is prompting, where generation is conditioned on either human-crafted or automatically optimized tokens (Shin et al., 2020; Liu et al., 2023). In-context learning —prompting via providing input-label examples— is another widely used technique (Brown et al., 2020). Given the challenges of discrete optimization over tokens, there is a growing interest in methods that optimize real-valued embeddings (Lester et al., 2021). It is widely believed that these *soft prompts* offer greater expressiveness due to the expansive nature of continuous space. Furthermore, beyond only optimizing input embeddings, one can optimize the inputs of every attention layer (Li and Liang, 2021). This technique, *prefix-tuning*, has proven to be very successful and competitive to full fine-tuning (Liu et al., 2022).

While context-based fine-tuning approaches have witnessed impressive empirical successes and widespread adoption, we have little theoretical understanding of how they work. In this work, we analyse the influence of prompts and prefixes on the internal computations of a pretrained model and delineate their limitations. Specifically, we address the following questions:

1. **Soft prompting and prefix-tuning are motivated by the embedding space being larger than the token space. However, can a transformer utilize the additional capacity?** We show that with a careful choice of transformer weights, controlling a single embedding can generate any of the V^N completions of N tokens, while controlling a token can produce only V completions, with V being the vocabulary size. Thus, a transformer can indeed exploit the embedding space.
2. **Since prefix-tuning is more expressive than prompting, is it as expressive as full fine-tuning?** Despite the expressiveness of continuous space, prefix-tuning has structural limitations. A prefix cannot change the relative attention over the content tokens and can only bias the output of the attention block in a constant direction. In contrast, full fine-tuning can learn new attention patterns and arbitrarily modify attention block outputs, making it strictly more powerful.
3. **If context-based fine-tuning methods suffer from such structural limitations, how come they have high empirical performance?** We show that the prefix-induced bias can steer the model towards a pretraining task. Prefix-tuning can also combine skills picked up during pretraining to solve some new tasks similar to pretraining tasks. However, it cannot learn a completely new task. This is not simply because of the small number of learnable parameters: fine-tuning the same number of parameters can be sufficient to learn the novel task. Hence, context-based fine-tuning can elicit or combine pretrained model skills but cannot learn completely new behaviors.

2 BACKGROUND

2.1 THE TRANSFORMER ARCHITECTURE

We outline a formulation of a simplified transformer architecture (Vaswani et al., 2017). Assume that the model has vocabulary size V (also referred to as *number of tokens*). The input is a sequence $(\mathbf{x}_1, \dots, \mathbf{x}_p)$, $\mathbf{x}_i \in \{1, \dots, V\}$, $\forall i$. Each token is mapped to a d_e -dimensional vector that is the \mathbf{x}_i -th column of an embedding matrix $\mathbf{E} \in \mathbb{R}^{d_e \times V}$. The attention mechanism is position-invariant, so typically position encodings are added. For a model with maximum input length N (*context size*), we use a one-hot position encoding $\mathbf{e}_N(i)$ concatenated with the embedding. Therefore, the embedding for the i -th position provided to the first attention block would be $\mathbf{x}_i = [\mathbf{E}_{:, \mathbf{x}_i}^\top, \mathbf{e}_N^\top(i)]^\top$.

A transformer consists of alternating attention blocks which operate across the whole sequence and Multi-Layer Perceptrons (MLPs) that operate on each individual element. Each attention block consists of H heads. Each head h is parameterized by query, key, and value matrices $\mathbf{W}_Q^h, \mathbf{W}_K^h \in \mathbb{R}^{k \times d_{in}}, \mathbf{W}_V^h \in \mathbb{R}^{d_{out} \times d_{in}}$.¹ The attention matrix $\mathbf{A}^h \in \mathbb{R}^{p \times p}$ for head h then has elements

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

where $p \leq N$ is the current length of the input and $T > 0$ is an inverse temperature parameter.² Equation (1) is the softmax function, hence with high enough T , it will result in approximately one-hot encoding of the maximum j . The output of the attention block \mathcal{A} with H heads is then $(\mathbf{t}_1, \dots, \mathbf{t}_p)$, where each position i is the sum of the attention-weighted values across all heads:

$$\begin{aligned} \mathcal{A}[(\mathbf{W}_Q^1, \dots, \mathbf{W}_Q^H), (\mathbf{W}_K^1, \dots, \mathbf{W}_K^H), (\mathbf{W}_V^1, \dots, \mathbf{W}_V^H)](\mathbf{x}_1, \dots, \mathbf{x}_p) &= (\mathbf{t}_1, \dots, \mathbf{t}_p), \\ \mathbf{t}_i &= \sum_{h=1}^H \sum_{j=1}^p \mathbf{A}_{ij}^h \mathbf{W}_V^h \mathbf{x}_j. \end{aligned} \quad (2)$$

A transformer then applies an MLP to each output of an attention block before passing them to next attention block. We will consider linear layers $\mathcal{L}[\mathbf{M}, \mathbf{b}](\mathbf{x}) = \mathbf{M}\mathbf{x} + \mathbf{b}$ and ReLU-activated linear layers $\hat{\mathcal{L}}[\mathbf{M}, \mathbf{b}](\mathbf{x}) = \text{ReLU}(\mathbf{M}\mathbf{x} + \mathbf{b})$. When we compose attention blocks and linear or softmax layers, we will implicitly assume that the linear layer is applied to all positions of the sequence. Furthermore, we will use the *then* operator \circledast for left-to-right function composition. Therefore, a transformer model predicting confidences over the vocabulary can, for example, be represented as:

$$(\mathbf{y}_1, \dots, \mathbf{y}_p) = \left(\mathcal{A}_1 \circledast \hat{\mathcal{L}}_{1,1} \circledast \mathcal{L}_{1,2} \circledast \mathcal{A}_2 \circledast \hat{\mathcal{L}}_{2,1} \circledast \mathcal{L}_{2,2} \circledast \text{softmax} \right) \left(\left[\begin{array}{c} \mathbf{E}_{:, \mathbf{x}_1} \\ \mathbf{e}_N(1) \end{array} \right], \dots, \left[\begin{array}{c} \mathbf{E}_{:, \mathbf{x}_p} \\ \mathbf{e}_N(p) \end{array} \right] \right), \quad (3)$$

¹For the first block, d_{in} must be $d_e + N$ but may be different for the deeper blocks.

²A causal model has $\mathbf{A}_{ij} = 0$ for $j > i$. This does not affect our results so we will skip the masking step.

where the output dimension of the last layer has to be V . The next token for a deterministic transformer is selected to be the last element’s largest logit: $x_{p+1} = \arg \max_{u \in \{1, \dots, V\}} \mathbf{y}_{p,u}$. Given an input (x_1, \dots, x_p) , the model then autoregressively extends this sequence one token at a time, following Equation (3) either until the sequence reaches a length N or until a special termination token.

A transformer has no separation between the system prompt S , user provided input X and the autoregressively response Y . Thus, a sequence conditional on user input is denoted as $(S_1, \dots, S_{n_S}, X_1, \dots, X_{n_X}, Y_1, \dots, Y_{n_Y})$ and one without user input as $(S_1, \dots, S_{n_S}, Y_1, \dots, Y_{n_Y})$.

2.2 CONTEXT-BASED FINE-TUNING OF A PRETRAINED MODEL

We now define prompting, soft prompting and prefix-tuning with the previously introduced notation.

Prompting. The most frequently used content-based fine-tuning approach is *prompting*: prefixing the input (X_1, \dots, X_{n_X}) with a token sequence $S \in \{1, \dots, V\}^{n_S}$ to guide the model response: $(S_1, \dots, S_{n_S}, X_1, \dots, X_{n_X})$. This is how most people interact with language models such as ChatGPT.

Soft prompting. Soft prompting replaces the embeddings of the system input $\mathbf{E}_{:,S_i}$ with learned vectors $\mathbf{s}_i \in \mathbb{R}^{d_e}$ called *virtual tokens* (Hambardzumyan et al., 2021; Lester et al., 2021; Qin and Eisner, 2021). Hence, the input in Equation (3) is modified to be:

$$\left(\left[\begin{array}{c} \mathbf{s}_1 \\ \mathbf{e}_N(1) \end{array} \right], \dots, \left[\begin{array}{c} \mathbf{s}_{n_S} \\ \mathbf{e}_N(n_S) \end{array} \right], \left[\begin{array}{c} \mathbf{E}_{:,X_1} \\ \mathbf{e}_N(n_S + 1) \end{array} \right], \dots, \left[\begin{array}{c} \mathbf{E}_{:,X_{n_X}} \\ \mathbf{e}_N(n_S + n_X) \end{array} \right] \right) \quad (4)$$

with \mathbf{s}_i chosen to maximize the likelihood of a target response $Y=(Y_1, \dots, Y_{n_Y})$, i.e., $\arg \max_{\mathbf{s}_1, \dots, \mathbf{s}_{n_S} \in \mathbb{R}^{d_e}} \sum_{j=1}^{n_Y} \log \mathbf{y}_{n_S+n_X+j, Y_j}$, where $\mathbf{y}_{n_S+n_X+j}$ are autoregressively generated.

Prefix-tuning. Prefix-tuning applies soft prompting across the depth of the model (Li and Liang, 2021; Liu et al., 2022). The first n_S positions for all attention blocks are learnable parameters, replacing the input $(x_1^l, \dots, x_{n_X}^l)$ for layer l with $(s_1^l, \dots, s_{n_S}^l, x_1^l, \dots, x_{n_X}^l)$, where all \mathbf{s}_i^l constitute the prefix. Hence, prefix-tuning can be formulated as $\arg \max_{\{\mathbf{s}_i^l, \dots, \mathbf{s}_{n_S}^l\}_{l=1}^{n_S}} \sum_{j=1}^{n_Y} \log \mathbf{y}_{n_S+n_X+j, Y_j}$. Prefix-tuning has been successful at fine-tuning models (Vu et al., 2022; Wu and Shi, 2022; Choi and Lee, 2023; Ouyang et al., 2023; Bai et al., 2023), leading to calls for language models provided as a service (La Malfa et al., 2023) to allow providing prefixes instead of prompts (Sun et al., 2022).

Any token-based prompt (S_1, \dots, S_{n_S}) has a corresponding soft prompt $(\mathbf{s}_i = \mathbf{E}_{:,S_i})$ but the reverse does not hold. Similarly, every soft prompt $(\mathbf{s}_1, \dots, \mathbf{s}_{n_S})$ can be represented as a prefix by setting the deeper prefixes to be the values that the model would compute at these positions $(\mathbf{s}_i^l = (\mathcal{A}_1 \circ \dots \circ \mathcal{L}_{l-1,-1})([\mathbf{s}_1^\top, \mathbf{e}_N^\top(1)]^\top, \dots, [\mathbf{s}_l^\top, \mathbf{e}_N^\top(l)]^\top))$. The reverse also does not hold: there are prefixes that cannot be represented as a soft prompt. A hierarchy emerges: *prompting* < *soft prompting* < *prefix-tuning*, with prefix-tuning the most powerful of the three. Hence, we focus on examining its performance relative to full fine-tuning but our findings also apply to prompting and soft prompting.

3 SOFT PROMPTING HAS MORE CAPACITY THAN PROMPTING

The success of soft prompting (and by extension, prefix-tuning) is commonly attributed to the larger capacity of the uncountably infinite embedding space compared to the finite token space. Yet, increased capacity from this larger domain is beneficial only if the model can utilize it. This section confirms that this is indeed the case by constructing a transformer that can generate exponentially more completions by varying a single virtual token than by varying a hard token.

Consider unconditional generation (representing a function with no inputs) with a single system token: $(Y_1, \dots, Y_N) = f(S_1) = f_{S_1}$. For a deterministic autoregressive function, there are a total of V functions in this family, hence the upper bound on the number of outputs of length N that one can generate by varying the first token S_1 is V : the first token fully determines the rest of the sequence. Generally, if one varies the first N_S tokens, there are at most V^{N_S} unique outputs. What if instead of the token S_1 we vary a single virtual token \mathbf{s}_1 : $(Y_1, \dots, Y_N) = f(\mathbf{s}_1) = f_{\mathbf{s}_1}$? This family of functions is indexed by a real vector and hence is infinite: in principle, one could generate all V^N possible

output sequences by only controlling s_1 .³ Still, that does not mean that the transformer architecture can represent a function that achieves that in practice, i.e., it is not obvious if there is a surjective map from $\{f_{s_1} : s_1 \in \mathbb{R}^{d_e}\}$ to $\{1, \dots, V\}^N$. We demonstrate that, in fact, there is a transformer f for which such a surjective map exists:

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 $s_1 = ((Y_1 - 1)/V, \dots, (Y_N - 1)/V)$.*

However, conditional generation is more interesting: given a user input (X_1, \dots, X_{n_X}) , we want to generate a target response (Y_1, \dots, Y_{n_Y}) . Even in the simple case of one system token, the user provides one token and the model generates one token in response ($Y_1 = f(S_1, X_1) = f_{s_1}(X_1)$), we cannot control response of the model to any user input with the system token. There are V^V maps from X_1 to Y_1 , but S_1 can take on only V values: $|\{f_{s_1} : S_1 \in 1, \dots, V\}| = V < V^V$. Hence, tokens cannot be used to specify an arbitrary map from user input to model output. However, a single virtual token can specify any of the V^V maps, i.e., there exists a transformer $f_{s_1}(X_2)$ for which there is a surjective map from $\{f_{s_1} : s_1 \in \mathbb{R}^{d_e}\}$ to $\{1, \dots, V\}^{\{1, \dots, V\}}$.

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 $s_1 = (m(1)/V, \dots, m(V)/V)$. That is, by selecting s_1 we control the model response to any user input.*

Theorem 2 builds on Theorem 1 by showing that soft prompting is also more expressive for governing the *conditional* behavior of a transformer model. This also holds for longer responses $n_Y > 1$ by increasing the length of the soft prompt, or longer user inputs $n_X > 1$, by increasing the depth of the model. We provide proofs in Appendix A, as well as working Python implementations.

This section showed that soft prompting, and by implication, prefix-tuning, possess greater expressiveness than prompting. As we can fully determine the map from user input to model response using virtual tokens, our findings may appear to suggest that soft prompting is as powerful as full fine-tuning. However, this is not at all the case. There are structural constraints on the capabilities of soft prompting and prefix-tuning; they cannot facilitate the learning of an entirely new task. The following section elucidates this discrepancy and reconciles these seemingly contradictory results.

4 PREFIX-TUNING CAN ONLY BIAS THE OUTPUT OF AN ATTENTION HEAD

We just saw that soft prompting and prefix-tuning can fully control the conditional behavior of a transformer. However, that assumed a specific design for the network weights. Given a fixed pretrained model, as opposed to a manually crafted one, can prefix-tuning be considered equally powerful to full fine-tuning? In this section, we show that a prefix S cannot change the relative attention over the content X, Y and can only bias the attention block outputs in a subspace of rank n_S , the length of the prefix, making it strictly less powerful than full fine-tuning.

While full fine-tuning can alter the attention pattern of an attention head, prefix-tuning cannot. Recall the attention A_{ij} position i gives to position j for a trained transformer (Equation (1)):

$$A_{ij} = \frac{\exp(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_j)}{\sum_{r=1}^p \exp(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_r)} = \frac{\exp(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{H} \mathbf{x}_j)}{\sum_{r=1}^p \exp(T/\sqrt{k} \mathbf{x}_i^\top \mathbf{H} \mathbf{x}_r)}, \quad (5)$$

where $\mathbf{W}_Q^\top \mathbf{W}_K = \mathbf{H}$. Full fine-tuning can enact arbitrary changes to \mathbf{W}_Q and \mathbf{W}_K and hence, assuming the input does not change (e.g., at the first attention layer), we get the following attention:

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

³For example, LLaMA-7B (Touvron et al., 2023) has 24 426 unique completions when prompted with each of its 32 000 tokens and we found a non-exhaustive set of 46 812 unique 10-token-long sequences by controlling the first virtual token. Hence, in practice, one can generate more outputs by soft prompting than by prompting.

where the changes to \mathbf{W}_Q and \mathbf{W}_K are folded into $\Delta\mathbf{H}$. It is clear that by varying $\Delta\mathbf{H}$ full fine-tuning can change the attention patterns arbitrarily. However, let us see how is attention affected by the presence of a prefix. For now, assume we have a prefix of length one (\mathbf{s}_1) at position 0.

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

The numerator of $\mathbf{A}_{ij}^{\text{pt}}$ is the same as in Equation (5), i.e., the prefix does not affect it. It only adds the term $\exp(\frac{T}{\sqrt{k}} \mathbf{x}_i^\top \mathbf{H} \mathbf{s}_1)$ to the denominator. Therefore, the attention position i gives to the content positions $j \geq 1$ is simply scaled down by the attention it now gives to the prefix. If *tomato* attends the most to *salad* in a particular context, no prefix can change that. This becomes evident by rewriting $\mathbf{A}_{ij}^{\text{pt}}$ as the attention of the pretrained model scaled by the attention ‘‘stolen’’ by the prefix:

$$\mathbf{A}_{ij}^{\text{pt}} = \mathbf{A}_{ij} \sum_{r=1}^p \mathbf{A}_{ir}^{\text{pt}} = \mathbf{A}_{ij} (1 - \mathbf{A}_{i0}^{\text{pt}}). \quad (6)$$

Hence, prefix-tuning cannot affect the relative attention patterns across the content, it will only scale them down. In other words, one cannot modify what an attention head attends to via prefix-tuning.

Prefix-tuning only adds a bias to the attention block output. Let us see how this attention scaling down affects the output of the attention block. Following Equation (2), the output at position i for the pretrained (\mathbf{t}_i), the fully fine-tuned (\mathbf{t}_i^{ft}) and the prefix-tuned (\mathbf{t}_i^{pt}) models are as follows:⁴

$$\begin{aligned} \mathbf{t}_i &= \sum_{j=1}^p \mathbf{A}_{ij} \mathbf{W}_V \mathbf{x}_j, & \mathbf{t}_i^{\text{ft}} &= \sum_{j=1}^p \mathbf{A}_{ij}^{\text{ft}} (\mathbf{W}_V + \Delta \mathbf{W}_V) \mathbf{x}_j, \\ \mathbf{t}_i^{\text{pt}} &= \mathbf{A}_{i0}^{\text{pt}} \mathbf{W}_V \mathbf{s}_1 + \sum_{j=1}^p \mathbf{A}_{ij}^{\text{pt}} \mathbf{W}_V \mathbf{x}_j \stackrel{(6)}{=} \mathbf{A}_{i0}^{\text{pt}} \mathbf{W}_V \mathbf{s}_1 + \sum_{j=1}^p \mathbf{A}_{ij} (1 - \mathbf{A}_{i0}^{\text{pt}}) \mathbf{W}_V \mathbf{x}_j = \mathbf{A}_{i0}^{\text{pt}} \mathbf{W}_V \mathbf{s}_1 + (1 - \mathbf{A}_{i0}^{\text{pt}}) \mathbf{t}_i. \end{aligned} \quad (7)$$

Hence, prefix-tuning only biases the pretrained attention block value at each position i towards the constant vector $\mathbf{W}_V \mathbf{s}_1$, which is independent of the content ($\mathbf{x}_1, \dots, \mathbf{x}_p$). The content only affects the scale $\mathbf{A}_{i0}^{\text{pt}}$ of the bias via the amount of attention on the prefix. This is in contrast with full fine-tuning where $\Delta \mathbf{W}_Q$, $\Delta \mathbf{W}_K$ and $\Delta \mathbf{W}_V$ allow for a content-dependent change of the attention and value computation. Note that these results also hold for suffix-tuning (placing the prefix after the input) but *not* for suffix soft-prompting.

We validate that this indeed is the case when prefix-tuning real-world transformers. In Figures 6 and 7, we show that a prefix applied to LLaMA’s first layer does not change the relative attention distribution over the content positions X and results in a bias with a constant direction.

Longer prefixes define larger subspaces for the bias but are not fully utilized in practice. In the case of a longer prefix ($\mathbf{s}_1, \dots, \mathbf{s}_{n_S}$), the bias vector is in a subspace of dimensionality n_S : $\mathbf{t}_i^{\text{pt}} = \sum_{j=1}^{n_S} \mathbf{A}_{i,S_j}^{\text{pt}} \mathbf{W}_V \mathbf{s}_j + (1 - \sum_{j=1}^{n_S} \mathbf{A}_{i,S_j}^{\text{pt}}) \mathbf{t}_i$, where i goes over the content and j over the prefix positions. Larger prefixes thus have a larger subspace to modify the attention block output. The specific direction is determined by the relative distribution of attention across the prefix positions. However, when we examine the distribution of attention across the prefix positions for various inputs as in Appendix B, it appears that the prefixes do not span this subspace. Regardless of the input, the attention $\mathbf{A}_{i,S_j}^{\text{pt}}$ over the prefix positions remains nearly constant. Thus, prefix-tuning does not seem to make full use of the space that the vectors $\mathbf{W}_V \mathbf{s}_j$ span. We hypothesise that this is due to the two competing optimization goals for the vectors \mathbf{s}_j : at the same time they need to ‘‘grab attention’’ when interacting with \mathbf{W}_K and determine the bias direction when multiplied with \mathbf{W}_V .

So, is prefix-tuning equivalent to full fine-tuning or is it less powerful than full fine-tuning? In Section 3, we showed that prefix-tuning, in principle, has a large capacity to influence the behavior of the model. But then, in this section, we showed that it has some severe limitations, including not being able to affect the attention pattern and only biasing the attention layer activations. These two results seem to be contradicting one another, so how do we reconcile them?

The constructions for the results in Section 3 (described in Appendix A) are simply an algorithm that extracts the completion from a lookup table encoded in the virtual tokens. The attention patterns

⁴He et al. (2021a) show a similar analysis but do not study the expressiveness of prefix-tuning.

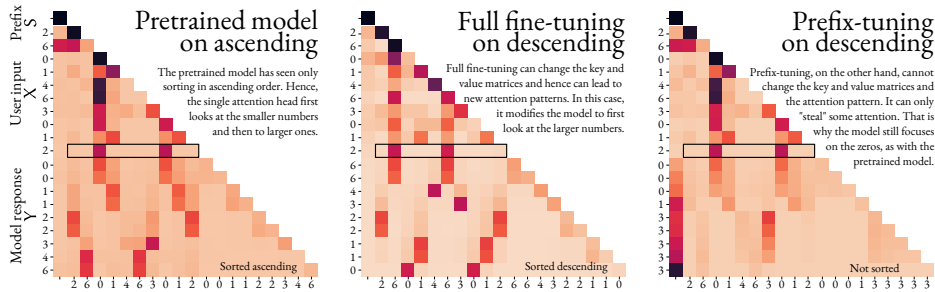


Figure 1: Attention patterns of a small transformer pretrained on sorting in ascending order. The model is given the prefix S and user input X and generates Y autoregressively. We have highlighted the attention when the first response Y_1 is being generated. Full fine-tuning sorts in descending order but prefix-tuning cannot as it cannot update the learned attention. Note how the relative attention of X to X in the left and right plots is exactly the same: the prefix cannot change the attention pattern for the same inputs. The relative attention of X to X in the center plot is very different because full fine-tuning can arbitrarily change W_Q and W_K .

are simply extracting the current position embedding and the virtual token and hence the attention does not depend on the actual content in the tokens. There is no need to learn a new attention pattern to learn a different map from input to output. Furthermore, the virtual token designates the map precisely by acting as a bias. Therefore, the observations in these two sections do not contradict one another. Soft prompting and prefix-tuning can be on par with full fine-tuning but only in very limited circumstances: when all the knowledge is represented in the virtual token as a lookup table and the model simply extracts the relevant entry. Transformers do not behave like this in practice. Models are typically trained with token inputs rather than virtual tokens. Moreover, if we had a lookup table of the responses to each input we would not need a learning algorithm in the first place.

Therefore, the limitations from this section hold for real-world pretrained transformers. Then how come prefix-tuning has been reported to achieve high accuracy and often to be competitive to full fine-tuning? The next section aims to explain when and why prefix-tuning can work in practice.

5 THE BIAS CAN ELICIT SKILLS FROM THE PRETRAINED MODEL

Pretraining potentially exposes a model to different classes of completions for the same token sequence. For a string like *I didn't enjoy the movie*, the model may have seen completions such as *I found the acting to be sub par*, *This is negative sentiment* or *Je n'ai pas aimé le film*. Hence, a pretrained model could do text completion, sentiment analysis, or translation. However, the input does not fully determine the desired completion type and the model can generate any one of them. Hence, following our results from Section 4, we hypothesise that prefix-tuning cannot be used to gain *new knowledge* but can bring to the surface *latent knowledge* present in the pretrained model.⁵ We test this hypothesis by constructing small transformers trained on one or few tasks. We use a minimal transformer model (Karpathy, 2020) to show that prefix-tuning cannot learn a new task that full fine-tuning can, then that it can easily elicit a latent skill from pretraining, and finally how it can even learn *some* new tasks, provided they can be solved by combining pretraining skills.

Prefix-tuning cannot learn a new task requiring a different attention pattern.

To check if prefix-tuning can learn a new task, we train a 1-layer, 1-head transformer to sort numbers into ascending order and then fine-tune it to sort in descending order. During training, the model sees random sequences of 10 digits from 0 to 7 followed by their ascending sorted order. The pre-trained accuracy (fully matching the sorted sequence) is 91%. Full fine-tuning on the descending task leads to 85% test accuracy, hence full fine-tuning successfully learns the new task. However, prefix-tuning with a prefix size $n_S=1$ results in 0% accuracy, hence prefix-tuning fails to learn the new task at all.

Table 1: A transformer pretrained on sorting in ascending order cannot be prefix-tuned to sort in descending order. 10 random seeds.

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

⁵A similar hypothesis has also been proposed by Reynolds and McDonnell (2021) for fine-tuning in general.

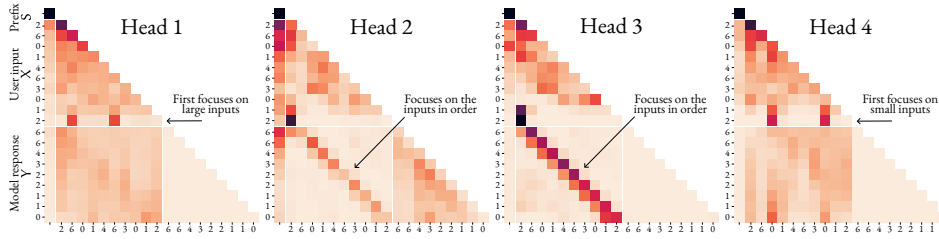
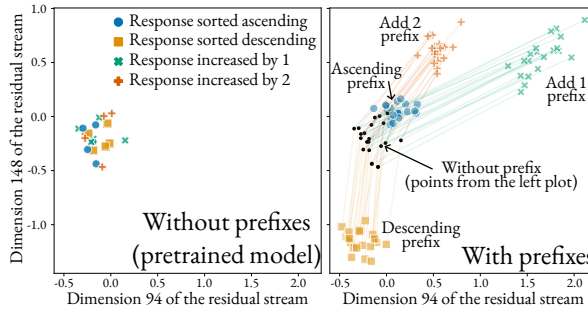


Figure 2: Model pretrained on the four tasks. The four attention heads specialize in the skills necessary to solve these tasks: look at the elements in order, look first at the smallest elements or first at the largest elements.

Figure 3: Attention block activations for ten sequences at the last input position (10) when pretrained on the four tasks. The left plot shows the pretrained activations t_{10} are not predictive of the completion. The right plot shows prefixes cluster the activations t_{10}^{pt} . Connecting the pretrained and prefixed activations highlights the bias. No dimensionality reduction is used; the clustering is solely due to the prefixes.



The attention patterns in Figure 1 show why this is the case: the pretrained model learns to attend first to the smallest numbers and then to the larger ones. When fully fine-tuned, the attention patterns are reversed: they now first attend to the largest values. However, following Section 4, prefix-tuning cannot change the attention pattern over the input sequence and will still attend to the smallest values. Therefore, prefix-tuning indeed cannot learn a new task if it requires new attention patterns.

Prefix-tuning can elicit a skill from the pretrained model. The second part of our hypothesis was that prefix-tuning can elicit latent skills in the pretrained model. To test that, we pretrain a 1-layer, 4-head model with solutions sorted in ascending (\nearrow) or descending (\searrow) order, or adding one (+1) or two (+2) to each element of the input sequence. Each solution is shown with 25% probability. The model has no indication of what the task is, hence, it assigns equal probability to all tasks, as shown in the first row in Table 2. Full fine-tuning for each task naturally results in high accuracy. However, prefix-tuning can also reach very high accuracy for all tasks: above 90%. Compared to the previous case, prefix-tuning is more successful here because the pretrained model contains the attention mechanisms for solving the four tasks, as shown in Figure 2.

Table 2: A transformer pretrained on several tasks can be prefix-tuned for one of them. 10 random seeds.

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	$0 \pm 0\%$	$0 \pm 0\%$	$1 \pm 2\%$	$98 \pm 5\%$

If all a prefix does is bias the attention layer activations, how can it steer the model to collapse its distribution onto one task? This is likely due to the attention block solving all tasks in parallel and placing their solutions in different subspaces of the residual stream (intermediate representation, Elhage et al., 2021). As the MLP needs to select one solution to generate, a further indicator on the selected task (or lack of selection thereof) should also be represented. The bias induced by the prefix then acts on this “selection subspace” to nudge the MLP to select the desired solution.

This can be clearly seen from the activations of the attention layer at the last input position (X_{n_X}). This is the position where the task selection happens as the first output element fully describes the task. Figure 3 shows plots of randomly selected dimensions of the residual stream with and without a prefix. The attention block activations of the pretrained model (without prefix) show no correlation with the output it is about to generate, demonstrating that the choice of completion is indeed random and is not determined by the activation block. However, the prefix-tuned activations for the same

inputs are clustered as a result of the prefix-induced bias. This shows that the bias induced by the prefix acts as a “task selector” of the subspace of the residual stream specializing in the desired task.

Prefix-tuning can combine knowledge from pretraining tasks to solve new tasks. Prefix-tuning eliciting one type of completion learned in pretraining starts to explain its practical utility. Still, prefix-tuning seems to be successful also at tasks that the pretrained model has not seen. As we showed above, a model trained to sort in one order cannot be prefix-tuned to sort in the other. Then how is it possible for prefix-tuning to learn a new task? We posit that this can happen, as long as the “skill” required to solve the new task is a combination of “skills” the pretrained model has seen.

We test this by pretraining a 40-layer 4-head model with the same four tasks. We prefix-tune ($n_S=10$) for two new tasks: incrementing the ascending sorted sequence ($\nearrow+1$) and double histogram (mapping each element to the number of elements in the sequence with the same value, \mathbb{H}). The pretrained model has not seen either task. Prefix-tuning results in 93% accuracy for $\nearrow+1$ which is a combination of the \nearrow and $+1$ pretraining tasks and just 1% for the \mathbb{H} task which requires different skills: finding other instances of the same token and counting. Note that \mathbb{H} is not a hard task: it requires 2 layers and 2 heads to be solved exactly (Weiss et al., 2021). Therefore, prefix-tuning is can indeed combine skills that the model has learned in order to solve a novel task but cannot learn a completely new task requiring new skills.

Table 3: Prefix tuning can learn a new task requiring only pretraining skills ($\nearrow+1$) but cannot learn a completely new task (\mathbb{H}). Average accuracy over 3 seeds.

Accuracy on:	\nearrow	\searrow	+1	+2	$\nearrow+1$	\mathbb{H}
Pretrained	17%	23%	34%	25%	0%	0%
Prefix-tune on \nearrow	100%	0%	0%	0%	0%	0%
Prefix-tune on \searrow	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 $\nearrow+1$	0%	0%	0%	0%	93%	0%
Prefix-tune on \mathbb{H}	0%	0%	0%	0%	0%	1%

6 EFFECTS OF PREFIX-TUNING BEYOND THE SINGLE ATTENTION LAYER

Section 4 focused exclusively on a single attention layer. However, even if a prefix can only induce a bias on its output, this bias, when passed through the following MLP and subsequent attention layers, can exhibit increasingly complex behaviors. While it is not immediately obvious that the limitations discussed above translate to deeper networks, we provide evidence to this end. Section 6.1 shows that a prefix can change the attention pattern of the following attention layer but only in a linear fashion while full fine-tuning can also enact bilinear effects. Section 6.2 shows that while prefix-tuning can be thought of as neural network training, the representational capacity of the resulting architecture is severely limited. Therefore, prefix-tuning appears to be less expressive than full fine-tuning, even when both methods are given the same number of learnable parameters.

6.1 PREFIX-TUNING CAN CHANGE THE ATTENTION, ALBEIT THE ONE OF THE NEXT LAYER

Let us examine how the prefix of one attention layer affects the following attention layer. For simplicity, assume there is no MLP between the attention layers, residual connections or layer norms: the output $\mathbf{t}_i^{(1)}$ of the first is the input $\mathbf{x}_i^{(2)}$ of the second. The pretrained outputs are $\mathbf{t}_i^{(1)} = \sum_{j=1}^p \mathbf{A}_{ij}^{(1)} \mathbf{W}_V^{(1)} \mathbf{x}_j^{(1)}$, resulting in the second layer attention $\tilde{\mathbf{A}}_{ij}^{(2)} = \tau/\sqrt{k} \mathbf{t}_i^{(1)\top} \mathbf{H}^{(2)} \mathbf{t}_j^{(1)}$. Here $\tilde{\mathbf{A}}_{ij}$ is the pre-softmax attention, i.e., $\mathbf{A}_{ij} = \exp \tilde{\mathbf{A}}_{ij} / \sum_{r=1}^p \exp \tilde{\mathbf{A}}_{ir}$. For prefix-tuning we then have:

$$\begin{aligned} \mathbf{t}_i^{\text{pt}(1)} &= \mathbf{A}_{i0}^{\text{pt}(1)} \mathbf{W}_V \mathbf{s}_1^{(1)} + \sum_{j=1}^p \mathbf{A}_{ij}^{\text{pt}(1)} \mathbf{W}_V^{(1)} \mathbf{x}_j^{(1)} \stackrel{(7)}{=} \underbrace{\mathbf{A}_{i0}^{\text{pt}(1)}}_{\alpha_i} \underbrace{\mathbf{W}_V \mathbf{s}_1^{(1)}}_{\boldsymbol{\mu}} + (1 - \mathbf{A}_{i0}^{\text{pt}(1)}) \mathbf{t}_i^{(1)}, \\ \tilde{\mathbf{A}}_{ij}^{\text{pt}(2)} &= \frac{\tau}{\sqrt{k}} \mathbf{t}_i^{\text{pt}(1)\top} \mathbf{H}^{(2)} \mathbf{t}_j^{\text{pt}(1)}, \\ &= \frac{\tau}{\sqrt{k}} (\alpha_i \alpha_j \underbrace{\boldsymbol{\mu}^\top \mathbf{H}^{(2)} \boldsymbol{\mu}}_{\text{constant}} + \alpha_j (1 - \alpha_i) \underbrace{\mathbf{t}_i^{(1)\top} \mathbf{H}^{(2)} \boldsymbol{\mu}}_{\text{depends only on } \mathbf{t}_i^{(1)}} + \alpha_i (1 - \alpha_j) \underbrace{\boldsymbol{\mu}^\top \mathbf{H}^{(2)} \mathbf{t}_j^{(1)}}_{\text{depends only on } \mathbf{t}_j^{(1)}} + (1 - \alpha_i)(1 - \alpha_j) \underbrace{\mathbf{t}_i^{(1)\top} \mathbf{H}^{(2)} \mathbf{t}_j^{(1)}}_{\text{pretrained attention } \tilde{\mathbf{A}}_{ij}^{(2)}}). \end{aligned}$$

The presence of $\boldsymbol{\mu}$ shows that the prefix of layer 1 can change the attention pattern of the following layer. This change is content-specific: the second and the third terms depend on the inputs, hence a simple bias can affect the attention when passed through MLPs and further attention blocks. Compare with Equation (6), which showed a prefix cannot change the attention of the same layer.

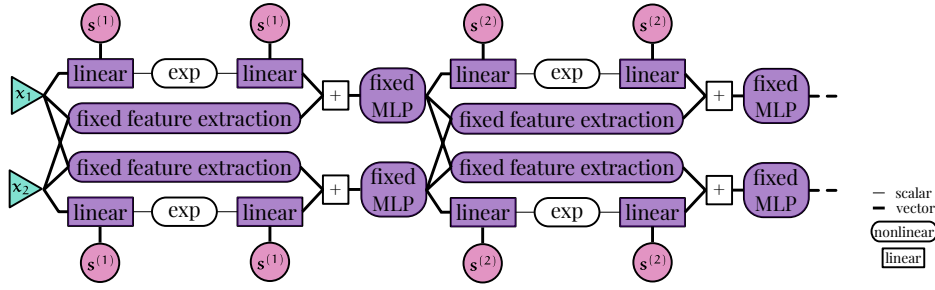


Figure 4: Prefix-tuning as a neural network architecture. While linearities and non-linearities are present, the only learnable parameters $s^{(1)}, s^{(2)}, \dots$ have limited interaction with the inputs x_1 and x_2 . The interaction of the prefix parameters with each input is only via the scalar attention, shown here with a light connection. The mixing of information between the inputs happens via residual connections with the pretrained fixed feature extraction and hence is not learnable. The MLP is also fixed and hence only acts as a multivariate activation function. This limited interaction explains why prefix-tuning struggles to learn new tasks even in deeper models.

Still, even considering this cross-layer effect, prefix-tuning is more limited in its expressiveness than full fine-tuning. While the second and the third terms are input-dependent, each depends on one input position only. The prefix does not change the bilinear dependency on both the query and key. This is something that the full fine-tuning can achieve: $\hat{A}_{ij}^{\text{ft}(2)} = T/\sqrt{k} \mathbf{t}_i^{\text{ft}(1)\top} (\mathbf{H}^{(2)} + \Delta \mathbf{H}^{(2)}) \mathbf{t}_j^{\text{ft}(1)}$.

6.2 EXPRESSIVITY OF PREFIX-TUNING ACROSS DEEPER MODELS

We considered the effect of the presence of the prefix in the first attention layer on the attention of the second. However, the effects become more complex as one adds more attention attention layers. We also ignored the MLPs between, but in practice they can play an important role. Here, instead, we analyse prefix-tuning as learning a neural network. We argue that while the resulting architecture includes both linear operation and non-linear activations, the structure is unlikely to learn efficiently.

For simplicity, we will consider two inputs x_1 and x_2 and a single prefix s . The output of the attention head, parameterized by s is then:

$$\begin{aligned} \mathcal{A}_s(x_1, x_2) &= \langle y_1, y_2 \rangle \\ y_1 &= \frac{\exp(x_1^\top \mathbf{H} s) \mathbf{W}_V s + \exp(x_1^\top \mathbf{H} x_1) \mathbf{W}_V x_1 + \exp(x_1^\top \mathbf{H} x_2) \mathbf{W}_V x_2}{C_1} \\ y_2 &= \frac{\exp(x_2^\top \mathbf{H} s) \mathbf{W}_V s + \exp(x_2^\top \mathbf{H} x_1) \mathbf{W}_V x_1 + \exp(x_2^\top \mathbf{H} x_2) \mathbf{W}_V x_2}{C_2} \end{aligned} \quad (8)$$

where we have omitted the T/\sqrt{k} factors and have folded the softmax normalization into C_1 and C_2 . The **layer inputs**, **pretrained parameters** and the **learnable parameters** are correspondingly highlighted. The attention head is clearly a non-linear operation. However, the learnable parameter s participates only in the left term. It interacts with only one of the inputs at a time and only by computing a single scalar value $x_i^\top \mathbf{H} s$. As we discussed above, the interaction between x_1 and x_2 is non-trainable and can be thought as a hard-coded feature extraction. Each of the outputs is then passed to a pretrained MLP which can be thought of as an activation function. This would be a *multivariate* activation function, which while unusual in the contemporary practice has been studied before (Solazzi and Uncini, 2004). Figure 4 illustrates the computation graph of the resulting neural network and shows that the only learnable interaction between the inputs is indirect. Therefore, prefix-tuning can be considered as learning a neural network where the only interaction between the inputs happens via non-learnable residual connections. Nevertheless, the alternating linear and nonlinear operations are reminiscent of the standard neural network architecture and their universal approximation properties (Hassoun, 1995). That begs the question if the prefix-tuning architecture can be an universal approximator and whether it would be a parameter-efficient one.

An example of prefix-tuning failing to be an universal approximator. While we leave the formal analysis of the representational capacity of prefix-tuning as future work, we provide an example of pretrained parameters for which the architecture is *not* a universal approximator. As can be seen in Figure 4, all information must pass through the non-learnable MLPs. Thus, if the MLPs destroy all input information, there is no value for the prefixes $s^{(1)}, s^{(2)}, \dots$ that can change that fact. The MLPs can destroy all information, if, e.g., one of their linear layers has a zero weight matrix. While this is unlikely to be learned with typical pretraining, this demonstrates that if prefix-tuning could be a universal approximator, that would pose specific requirements on the pretrained MLPs and it is not clear whether pretraining would satisfy these requirements.

Even if prefix-tuning could be an universal approximator, it would not be parameter-efficient. Whether prefix-tuning is an universal approximator is not of practical relevance as even if it was, it would be much less parameter-efficient than other comparable approaches. While Equation (8) and Figure 4 already hint at that, we designed an experiment to this end. Recall that our pretrained model in Section 5 failed to learn the double histogram task (\mathbb{H}) as it required novel attention patterns. A rank-1 Low Rank Adaptation (LoRA, Hu et al., 2021) applied only to the MLPs in a 4-layer 4-head model pretrained in the exact same way results in 92% accuracy on the \mathbb{H} task. The number of parameters for the LoRA fine-tuning is exactly the same as for a prefix of size 12. However, as can be expected from the results in Section 5, training this prefix results in 0% accuracy. Therefore, prefix-tuning fails to learn a task that LoRA with the same number of parameters can learn.

In conclusion, while prefix-tuning, prompting, soft prompting and in-context learning have complex effects in deeper models, the interaction of the learnable parameters and the model inputs likely still results in very limited expressiveness. In particular, we demonstrated that LoRA can be used to learn a completely new task while prefix-tuning with the exact same number of parameters cannot.

7 DISCUSSION AND RELATED WORKS

Understanding fine-tuning and prefix-tuning. Prior works show that prefixes have low intrinsic dimension allowing transfer to similar tasks and initialization of prefixes for new tasks (Qin et al., 2021; Su et al., 2022; Zhong et al., 2022; Wang et al., 2022b; Zheng et al., 2023). In this work, we offered theoretical insights into their results: this subspace is the span of the prefix-induced bias. Another line of work shows that skills can be localized in the parameter space of pretrained models (Wang et al., 2022a; Panigrahi et al., 2023). Here, we showed that it is also possible to identify subspaces of the residual stream corresponding to individual tasks and select them via prefix-tuning.

Prompting and in-context learning. Prompting and in-context learning are a special case of prefix-tuning. Therefore, the limitations and mechanisms discussed in this work apply to prompting as well: prompts cannot change the distribution of attention of the first attention layer over the content following it and can only induce a bias on the output of this layer (Section 4). Even considering the cross-layer effects, a prompt is strictly less expressive than full fine-tuning (Section 6) and prompting is unlikely to enable the model to solve a completely new task. Our theory thus explains why Kossen et al. (2023) observed that in-context examples cannot overcome pre-training skills.

While context-based fine-tuning approaches cannot learn *arbitrary* new tasks, as shown in Section 5, they can leverage pre-trained skills. For instance, transformers can learn linear models in-context by mimicking gradient descent (Von Oswald et al., 2023) or by approximating matrix inversion (Akyürek et al., 2022). This is consistent with our theory: the prediction updates are enacted as biases in the activations of the attention block. Hence, despite the limitations discussed in this work, context-based methods can result in powerful fine-tuning if the pretrained model has “transferable skills” such as algorithmic fundamentals. Still, there are non-algorithmic applications where our theory predicts that in-context learning will fail, for example, translating to a language that the model has never seen before, even if large number of translation pairs are provided in-context.

Implications for catastrophic forgetting and model alignment. The lack of expressiveness of context-based fine-tuning can be a feature: desirable properties will be maintained. Full fine-tuning can result in catastrophic forgetting (He et al., 2021b; Luo et al., 2023; Mukhoti et al., 2023). Our theory shows that context-based methods will not destroy pretrained skills. Model alignment poses the reverse problem: ensuring that the model cannot pick up undesirable skills during fine-tuning.

Our results show that prompting and prefix-tuning cannot steer the model towards new adversarial behaviors. Hence, the recent successes in adversarial prompting (Zou et al., 2023) indicate that current model alignment methods just mask the undesirable skills rather than removing them.

Implications for model interpretability. An open question for language model interpretability is whether attention is sufficient for explainability (Jain and Wallace, 2019; Wiegrefe and Pinter, 2019). Section 5 points toward the negative: by interfering in the *output* of the attention layer with the bias induced by a prefix, we can change the behavior of the model, without changing its attention. On the flip side, prefix-tuning can be used to understand what “skills” a model has: if prefix-tuning for a task fails, then the model likely lacks one of the key “skills” for that task.

Limitations. The present analysis is largely limited to prefixing with prompts, soft prompts and for prefix-tuning. While our theoretical results hold for suffix-tuning, they do not necessarily apply to suffixing with prompts or soft prompts. That is because the deeper representations for prompt and soft prompt suffixes would depend on the previous positions. This does not apply to suffix-tuning as it fixes all intermediate representations. Therefore, whether suffixing is more expressive than prefixing remains an open question. Separately, while we provided evidence towards context-based fine-tuning methods being parameter inefficient learners, the formal analysis of the conditions under which they may be universal approximators remain an open question. Finally, we mostly considered simple toy problems. In practice, however, language models are pretrained with very large datasets and can pick up very complex behaviors. Hence, the extent to which the limitations we demonstrated apply to large-scale pretrained transformers also remains for future work.

8 CONCLUSION

This paper formally showed that fine-tuning techniques working in embedding space, such as soft prompting and prefix-tuning, are strictly more expressive than prompting which operates in the discrete token space. However, we then demonstrated that despite this larger expressivity, prefix-tuning suffers from structural limitations that prevent it from learning new attention patterns. As a result, it can only bias the output of the attention layer in a direction from a subspace of rank equal to the size of the prefix. We showed that this results in practical limitations by constructing minimal transformers where prefix tuning fails to solve a simple task. This result seems to be at odds with the empirical success of prefix-tuning. We provided explanations towards that. First, we showed that prefix-tuning can easily elicit a skill the pretrained model already has and can even learn a new task, if it has picked up the skills to solve it during pretraining. Second, we showed that the effect of the prefix-induced bias is more complicated and powerful when combined with downstream non-linear operations. However, it appears to be still strictly less expressive than full fine-tuning.

REPRODUCIBILITY STATEMENT

In order to facilitate the reproduction of our empirical results, validating our theoretical results, and further studying the properties of context-based fine-tuning, we [release all our code and resources](#) used in this work. Furthermore, in Appendix A we offer explicit constructions of transformers with the properties discussed in Section 3. We also provide Python implementations of these constructions that validate their correctness.

ACKNOWLEDGEMENTS

This work is supported by a UKRI grant Turing AI Fellowship (EP/W002981/1) and the EPSRC Centre for Doctoral Training in Autonomous Intelligent Machines and Systems (EP/S024050/1). AB has received funding from the Amazon Research Awards. We also thank the Royal Academy of Engineering and FiveAI.

REFERENCES

- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. 2022. [What learning algorithm is in-context learning? Investigations with linear models](#). In *International Conference on Learning Representations*.
- Jiaqi Bai, Zhao Yan, Jian Yang, Xinnian Liang, Hongcheng Guo, and Zhoujun Li. 2023. [Knowprefix-tuning: A two-stage prefix-tuning framework for knowledge-grounded dialogue generation](#). *arXiv preprint arXiv:2306.15430*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. [Language models are few-shot learners](#). *Advances in Neural Information Processing Systems*.
- YunSeok Choi and Jee-Hyong Lee. 2023. [CodePrompt: Task-agnostic prefix tuning for program and language generation](#). In *Findings of the Association for Computational Linguistics: ACL 2023*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. [A mathematical framework for transformer circuits](#). *Transformer Circuits Thread*.
- Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. 2021. [WARP: Word-level Adversarial ReProgramming](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*.
- Mohamad H Hassoun. 1995. *Fundamentals of artificial neural networks*. MIT press.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021a. [Towards a unified view of parameter-efficient transfer learning](#). In *International Conference on Learning Representations*.
- Tianxing He, Jun Liu, Kyunghyun Cho, Myle Ott, Bing Liu, James Glass, and Fuchun Peng. 2021b. [Analyzing the forgetting problem in pretrain-finetuning of open-domain dialogue response models](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for NLP](#). In *International Conference on Machine Learning*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*.
- Zhiqiang Hu, Yihuai Lan, Lei Wang, Wanyu Xu, Ee-Peng Lim, Roy Ka-Wei Lee, Lidong Bing, and Soujanya Poria. 2023. [LLM-Adapters: An adapter family for parameter-efficient fine-tuning of large language models](#). *arXiv preprint arXiv:2304.01933*.
- Sarthak Jain and Byron C. Wallace. 2019. [Attention is not explanation](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. [Scaling laws for neural language models](#). *arXiv preprint arXiv:2001.08361*.

-
- Andrej Karpathy. 2020. [minGPT GitHub Repository](#).
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners](#). *Advances in Neural Information Processing Systems*.
- Jannik Kossen, Tom Rainforth, and Yarin Gal. 2023. [In-context learning in large language models learns label relationships but is not conventional learning](#). *arXiv preprint arXiv:2307.12375*.
- Emanuele La Malfa, Aleksandar Petrov, Christoph Weinhuber, Simon Frieder, Ryan Burnell, Anthony G. Cohn, Nigel Shadbolt, and Michael Wooldridge. 2023. [The ARRT of Language-Models-as-a-Service: Overview of a new paradigm and its challenges](#). *arXiv preprint arXiv:2309.16573*.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-Tuning: Optimizing continuous prompts for generation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*.
- Vladislav Lialin, Vijeta Deshpande, and Anna Rumshisky. 2023. [Scaling down to scale up: A guide to parameter-efficient fine-tuning](#). *arXiv preprint arXiv:2303.15647*.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing](#). *ACM Computing Surveys*.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. [P-Tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*.
- Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. 2023. [An empirical study of catastrophic forgetting in large language models during continual fine-tuning](#). *arXiv preprint arXiv:2308.08747*.
- Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. 2021. [Recent advances in natural language processing via large pre-trained language models: A survey](#). *ACM Computing Surveys*.
- Jishnu Mukhoti, Yarin Gal, Philip H. S. Torr, and Puneet K. Dokania. 2023. [Fine-tuning can cripple your foundation model; preserving features may be the solution](#). *arXiv preprint arXiv:2308.13320*.
- Linyong Nan, Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xiangru Tang, Aadit Vyas, Neha Verma, Pranav Krishna, Yangxiaokang Liu, Nadia Irwanto, Jessica Pan, Faiaz Rahman, Ahmad Zaidi, Mutethia Mutuma, Yasin Tarabar, Ankit Gupta, Tao Yu, Yi Chern Tan, Xi Victoria Lin, Caiming Xiong, Richard Socher, and Nazneen Fatema Rajani. 2021. [DART: Open-domain structured data record to text generation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Yawen Ouyang, Yongchang Cao, Yuan Gao, Zhen Wu, Jianbing Zhang, and Xinyu Dai. 2023. [On prefix-tuning for lightweight out-of-distribution detection](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Abhishek Panigrahi, Nikunj Saunshi, Haoyu Zhao, and Sanjeev Arora. 2023. [Task-specific skill localization in fine-tuned language models](#). In *International Conference on Machine Learning*.
- Guanghui Qin and Jason Eisner. 2021. [Learning how to ask: Querying LMs with mixtures of soft prompts](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.

-
- Yujia Qin, Xiaozhi Wang, Yusheng Su, Yankai Lin, Ning Ding, Jing Yi, Weize Chen, Zhiyuan Liu, Juanzi Li, Lei Hou, et al. 2021. [Exploring universal intrinsic task subspace via prompt tuning](#). *arXiv preprint arXiv:2110.07867*.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#).
- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, et al. 2021. [Scaling language models: Methods, analysis & insights from training Gopher](#).
- Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. 2017. [Learning multiple visual domains with residual adapters](#). In *Advances in Neural Information Processing Systems*.
- Laria Reynolds and Kyle McDonell. 2021. [Prompt programming for large language models: Beyond the few-shot paradigm](#). In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*.
- Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang, Junlin Wu, and Yi-Shin Chen. 2018. [CARER: Contextualized affect representations for emotion recognition](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. [Auto-Prompt: Eliciting knowledge from language models with automatically generated prompts](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Mirko Solazzi and Aurelio Uncini. 2004. [Regularising neural networks using flexible multivariate activation function](#). *Neural Networks*, 17(2):247–260.
- Yusheng Su, Xiaozhi Wang, Yujia Qin, Chi-Min Chan, Yankai Lin, Huadong Wang, Kaiyue Wen, Zhiyuan Liu, Peng Li, Juanzi Li, Lei Hou, Maosong Sun, and Jie Zhou. 2022. [On transferability of prompt tuning for natural language processing](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. 2022. [Black-box tuning for language-model-as-a-service](#). In *International Conference on Machine Learning*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. [LLaMA: Open and efficient foundation language models](#). *arXiv preprint arXiv:2302.13971*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*.
- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, Joao Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. 2023. [Transformers learn in-context by gradient descent](#). In *International Conference on Machine Learning*.
- Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou’, and Daniel Cer. 2022. [SPoT: Better frozen model adaptation through soft prompt transfer](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022a. [Finding skill neurons in pre-trained transformer-based language models](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*.
- Zhen Wang, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, Huan Sun, and Yoon Kim. 2022b. [Multitask prompt tuning enables parameter-efficient transfer learning](#). In *International Conference on Learning Representations*.

-
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. [Finetuned language models are zero-shot learners](#). In *International Conference on Learning Representations*.
- Gail Weiss, Yoav Goldberg, and Eran Yahav. 2021. [Thinking like transformers](#). In *International Conference on Machine Learning*.
- Sarah Wiegrefe and Yuval Pinter. 2019. [Attention is not not explanation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*.
- Hui Wu and Xiaodong Shi. 2022. [Adversarial soft prompt tuning for cross-domain sentiment analysis](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Yuanhang Zheng, Zhixing Tan, Peng Li, and Yang Liu. 2023. [Black-box prompt tuning with subspace learning](#). *arXiv preprint arXiv:2305.03518*.
- Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and Dacheng Tao. 2022. [PANDA: Prompt transfer meets knowledge distillation for efficient model adaptation](#). *arXiv preprint arXiv:2208.10160*.
- Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. 2023. [Universal and transferable adversarial attacks on aligned language models](#). *arXiv preprint arXiv:2307.15043*.

A CONSTRUCTING TRANSFORMERS THAT UTILIZE THE CAPACITY OF THE EMBEDDING SPACE

A.1 UNCONDITIONAL GENERATION FOR A SINGLE VIRTUAL TOKEN

This section provides an explicit construction of a transformer with the properties described in Theorem 1. The goal is to construct a transformer that, by varying the choice of the virtual token, can generate any sequence of N tokens.

First, we need to specify how we encode the target sequence (Y_1, \dots, Y_N) into the virtual token \mathbf{s}_1 . We chose the size of the embedding (and hence of \mathbf{s}_1) to be N . This way, each element of \mathbf{s}_1 can represent one position of the target sequence. We then represent the token value by discretizing each element of \mathbf{s}_1 into V levels:

$$\mathbf{s}_1 = ((Y_1-1)/V, \dots, (Y_N-1)/V).$$

Note that this means that each element of \mathbf{s}_1 is in $[0, 1)$.

When predicting the token for the $i + 1$ position, the transformer needs to pick the i -th element of \mathbf{s}_1 , and then decode the corresponding value as a one-hot encoding representing the Y_i -th token.

We extract the i -th element of \mathbf{s}_1 using one attention block of two heads. The `fst` head always looks at the first position which is our virtual token \mathbf{s}_1 . For that purpose we create an attention head that always has $\mathbf{A}_{ij}^{\text{fst}} = 1$ if $j = 1$ and $\mathbf{A}_{ij}^{\text{fst}} = 0$ otherwise together with a value matrix $\mathbf{W}_V^{\text{fst}}$ that extracts the embedding. This is achieved with

$$\mathbf{W}_Q^{\text{fst}} = [\mathbf{0}_N, \mathbf{1}_N], \quad \mathbf{W}_K^{\text{fst}} = [\mathbf{0}_N, \mathbf{1}, \mathbf{1}_{N-1}], \quad \mathbf{W}_V^{\text{fst}} = [\mathbf{I}_N, \mathbf{0}_{N \times N}], \quad (9)$$

and a sufficiently high inverse temperature parameter T .

The `pos` head instead extracts the one-hot encoding of the current position. This can be done with an attention head that always attends only to the current position and a value matrix $\mathbf{W}_V^{\text{pos}}$ that extracts the position embedding as a one-hot vector:

$$\mathbf{W}_Q^{\text{pos}} = [\mathbf{0}_{N \times N}, \mathbf{I}_N], \quad \mathbf{W}_K^{\text{pos}} = [\mathbf{0}_{N \times N}, \mathbf{I}_N], \quad \mathbf{W}_V^{\text{pos}} = [\mathbf{0}_{N \times N}, \mathbf{I}_N]. \quad (10)$$

When the outputs of these two attention heads are summed, then only the element of \mathbf{s}_1 that corresponds to the current position will be larger than 1. From Equation (2) the output at the i -th position of the attention block is:

$$\mathbf{t}_i = \sum_{j=1}^p \mathbf{A}_{ij}^{\text{fst}} \mathbf{x}_j + \sum_{j=1}^p \mathbf{A}_{ij}^{\text{pos}} \mathbf{e}_N(j) = \mathbf{s}_1 + \mathbf{e}_N(i),$$

where $\mathbf{x}_1 = \mathbf{s}_1$ and $\mathbf{x}_j = E_{:, Y_{j-1}}$ for $j > 1$.

We can extract the value of \mathbf{s}_1 corresponding to the current position by subtracting 1 from the hidden state and apply ReLU: $\hat{\mathcal{L}}_{\text{ex}} = \hat{\mathcal{L}}[\mathbf{I}_N, -\mathbf{1}_N]$. Now, we are left with only one non-zero entry and that's the one corresponding to the next token. We can retain only the non-zero entry if we just sum all the entries of the hidden state with $\hat{\mathcal{L}}_{\text{sum}} = \hat{\mathcal{L}}[\mathbf{1}_N^\top, \mathbf{0}]$.

The final step is to map this scalar to a V -dimensional vector which has its maximum value at index Y_i . This task is equivalent to designing V linear functions, each attaining its maximum at one of $0, 1/V, \dots, (V-1)/V$. To construct this, we use the property of convex functions that their tangent is always under the plot of the function. Therefore, given a convex function $\gamma(x)$, we construct the i -th linear function to be simply the tangent of γ at i^{-1}/V . If we take $\gamma(x) = (x - 1/2)^2$, this results in the following linear layer:

$$\mathcal{L}_{\text{proj}} = \mathcal{L} \left[\left[\frac{2(1-1)}{V} - 1, \dots, \frac{2(V-1)}{V} - 1 \right]^T, \left[\frac{1}{4} - \frac{(1-1)^2}{V^2}, \dots, \frac{1}{4} - \frac{(V-1)^2}{V^2} \right]^T \right]. \quad (11)$$

Figure 5 shows the predictors for each individual token id.

With just two attention heads and three linear layers, the transformer $\mathcal{A}[(\mathbf{W}_Q^{\text{fst}}, \mathbf{W}_Q^{\text{pos}}), (\mathbf{W}_K^{\text{fst}}, \mathbf{W}_K^{\text{pos}}), (\mathbf{W}_V^{\text{fst}}, \mathbf{W}_V^{\text{pos}})] \circ \hat{\mathcal{L}}_{\text{ex}} \circ \hat{\mathcal{L}}_{\text{sum}} \circ \mathcal{L}_{\text{proj}} \circ \text{softmax}$ achieves the upper bound of V^N unique outputs by controlling a single virtual token at its input. Note that for this

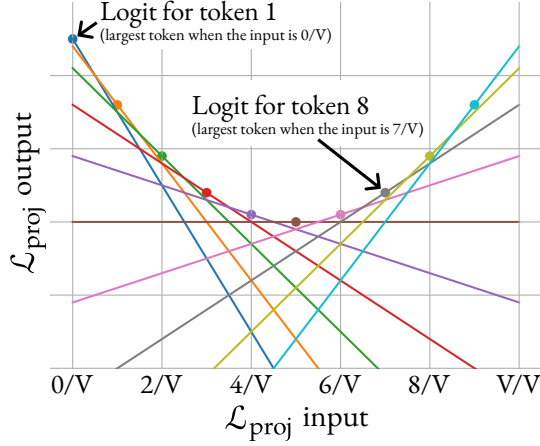


Figure 5: Illustration of the predictors for each token in the $\mathcal{L}_{\text{proj}}$ linear layer for $V = 10$. The layer is constructed in such a way that the i -th token has the highest confidence when the input is $i-1/V$.

construction, the choice of embedding matrix $E \in \mathbb{R}^{N \times V}$ does not matter. The same transformer architecture can generate only V unique outputs if we only control the first token instead. Therefore, it is indeed the case that the embedding space has exponentially more capacity for control than the token space. You can see this transformer implemented and running in practice in Section 2 of the `constructions` Jupyter notebook.

A.2 CONDITIONAL GENERATION FOR A SINGLE VIRTUAL TOKEN ($n_X = n_Y = 1$)

This section provides an explicit construction of a transformer with the properties described in Theorem 2. The goal is to construct a transformer that, by varying the choice of the virtual token, can cause the model to act as any map $m : [1, \dots, V] \rightarrow [1, \dots, V]$. In other words, by selecting the virtual token, we can fully control how the model will respond to any token the user may provide.

First, we need to specify how the map m will be encoded in the virtual token s_1 . We choose the embedding size d_e to be V . Now, we can use the same encoding scheme as before, but now each element in s_1 corresponds to a different user token, rather than to a position in the generated sequence:

$$s_1 = (m(1)/V, \dots, m(V)/V).$$

Therefore, the first element of s_1 designates the response if the user provides token 1, the second element is the response to the token 2, and so on.

Extracting the Y_i -th value from s_1 and decoding it can be done in a very similar way as for the unconditional case. The only difference is that instead of looking at the user input position, we look at its value. Take $E = I_V$ and $N = 2$.

Hence we have the following `val` head (only differing in the W_V matrix from Equation (10)):

$$W_Q^{\text{val}} = [0_{2 \times V}, I_2], \quad W_K^{\text{val}} = [0_{2 \times V}, I_2], \quad W_V^{\text{val}} = [I_V, 0_{V \times 2}].$$

We also need embedding of the first token, so we have a modified version of Equation (9):

$$W_Q^{\text{fst}} = [0_V, 1, 1], \quad W_K^{\text{fst}} = [0_V, 1, 0], \quad W_V^{\text{fst}} = [I_V, 0_{V \times 2}].$$

And hence the output of this attention block at the second position would be:

$$t_2 = \sum_{j=1}^2 A_{ij}^{\text{fst}} x_j + \sum_{j=1}^2 A_{ij}^{\text{val}} W_V^{\text{fst}} x_j = s_1 + e_V(Y_1).$$

Similarly to the unconditional case, only the entry of t_2 corresponding to the user token will have a value above 1 and that value would be $1 + m(x_1)/V$.

We can now extract the one-hot representation of the target token using the same approach as before, just adjusting for the different hidden state size: $\hat{\mathcal{L}}_{\text{ex}} = \hat{\mathcal{L}}[\mathbf{I}_V, -\mathbf{1}_V]$, $\hat{\mathcal{L}}_{\text{sum}} = \hat{\mathcal{L}}[\mathbf{1}_V^\top, 0]$, and the same projection had as before (Equation (11)). The final transformer is then: $\mathcal{A}[(\mathbf{W}_Q^{\text{fst}}, \mathbf{W}_Q^{\text{val}}), (\mathbf{W}_K^{\text{fst}}, \mathbf{W}_K^{\text{val}}), (\mathbf{W}_V^{\text{fst}}, \mathbf{W}_V^{\text{val}})] \circ \hat{\mathcal{L}}_{\text{ex}} \circ \hat{\mathcal{L}}_{\text{sum}} \circ \mathcal{L}_{\text{proj}} \circ \text{softmax}$. You can see this transformer implemented and running in practice in Section 3 of the `constructions` Jupyter notebook.

A.3 CONDITIONAL GENERATION FOR LONGER RESPONSES ($n_X = 1, n_Y > 1$)

We can obtain longer responses via a simple extension. If the response length is N_0 , then we can encode the map $m : [1, \dots, V] \rightarrow [1, \dots, V]^{N_0}$ in N_0 virtual tokens, each corresponding to one of the target positions:

$$\mathbf{s}_i = (m^{(1)}_i/V, \dots, m^{(V)}_i/V) \text{ for } i = 1, \dots, N_0.$$

For this model we would then have $N = 2N_0$ and $d_e = V$.

First, we need a head that always looks at the token provided by the user, which will be at position $N_0 + 1$:

$$\mathbf{W}_Q^{\text{user}} = [\mathbf{0}_V, \mathbf{1}_N], \quad \mathbf{W}_K^{\text{user}} = [\mathbf{0}_{(V+N_0)}, 1, \mathbf{0}_{(N_0-1)}], \quad \mathbf{W}_V^{\text{user}} = [\mathbf{I}_V, \mathbf{0}_{V \times N}].$$

In order to consume the map at the right location, we need to also look at the embedding of the token N_0 positions before the one we are trying to generate:

$$\mathbf{W}_Q^{\text{back}} = \begin{bmatrix} \mathbf{0}_{N_0 \times (N_0+V)} & \mathbf{I}_{N_0} \\ \mathbf{0}_{N_0 \times (N_0+V)} & \mathbf{0}_{N_0 \times N_0} \end{bmatrix}, \quad \mathbf{W}_K^{\text{back}} = [\mathbf{0}_{N \times V}, \mathbf{I}_N], \quad \mathbf{W}_V^{\text{back}} = [\mathbf{I}_V, \mathbf{0}_{V \times N}].$$

From here on, the decoding is exactly the same as in the $n_X = n_Y = 1$ case. The final transformer is then: $\mathcal{A}[(\mathbf{W}_Q^{\text{user}}, \mathbf{W}_Q^{\text{back}}), (\mathbf{W}_K^{\text{user}}, \mathbf{W}_K^{\text{back}}), (\mathbf{W}_V^{\text{user}}, \mathbf{W}_V^{\text{back}})] \circ \hat{\mathcal{L}}_{\text{ex}} \circ \hat{\mathcal{L}}_{\text{sum}} \circ \mathcal{L}_{\text{proj}} \circ \text{softmax}$. You can see this transformer implemented and running in practice in Section 4 of the `constructions` Jupyter notebook.

A.4 CONDITIONAL GENERATION FOR LONGER USER INPUTS ($n_X > 1, n_Y = 1$)

Finally, we consider the case when the user input X is longer. This is a bit more complicated because we need to search through a domain of size V^V . We will only consider the case with $n_X = 2$ where we would need two attention layers. A similar approach can be used to construct deeper models for $n_X > 2$. Finally, combining the strategy in the previous section for longer responses with the strategy in this section for longer user inputs allows us to construct transformers that map from arbitrary length user strings to arbitrary length responses.

In order to encode a map $m : [1, \dots, V]^2 \rightarrow [1, \dots, V]$ into a single virtual token we would need a more involved construction than before. Similarly to how we discretized each element of the virtual token \mathbf{s}_1 in V levels before, we are going to now discretize it into V^V levels. Each one of these levels would be one of the V^V possible maps from the *second* user token to the response. The first user token would be used to select the corresponding element of \mathbf{s}_1 . Then this scalar will be “unpacked” into a new vector of V elements using the first attention block. Then, the second user token will select an element from this unpacked vector, which will correspond to the target token.

We construct the virtual token as follows:

$$\mathbf{s}_1 = \left[\sum_{i=1}^V m_1(i) \times \frac{V^{i-1}}{V^V}, \dots, \sum_{i=1}^V m_V(i) \times \frac{V^{i-1}}{V^V} \right],$$

where $m_f(x) = m(f, x)$ is a map from the second user token to the response when the first token is fixed to be f .

An additional change from the previous constructions is that we are going to divide the residual stream into two sections. This is in line with the theory that different parts of the residual stream specialize for different communications needs by different attention heads (Elhage et al., 2021). We will use the first half of the residual stream to extract and “unpack” the correct mapping from second

token to target token, while the second half of the residual stream will be used to copy the second token value so that the second attention layer can use it to extract the target. As usual, the embedding matrix will be the identity matrix: $\mathbf{E} = \mathbf{I}_V$. Finally, for convenience, we will also use a dummy zero virtual token that we will attend to when we want to not attend to anything. This results in context size $N = 4$ with the input being

$$\left(\begin{bmatrix} 0_V \\ \mathbf{e}_N(1) \end{bmatrix}, \begin{bmatrix} \mathbf{s}_1 \\ \mathbf{e}_N(2) \end{bmatrix}, \begin{bmatrix} \mathbf{E}_{:,X_1} \\ \mathbf{e}_N(3) \end{bmatrix}, \begin{bmatrix} \mathbf{E}_{:,X_2} \\ \mathbf{e}_N(4) \end{bmatrix} \right) = \left(\begin{bmatrix} \mathbf{0}_V \\ \mathbf{e}_N(1) \end{bmatrix}, \begin{bmatrix} \mathbf{s}_1 \\ \mathbf{e}_N(2) \end{bmatrix}, \begin{bmatrix} \mathbf{e}_V(X_1) \\ \mathbf{e}_N(3) \end{bmatrix}, \begin{bmatrix} \mathbf{e}_V(X_2) \\ \mathbf{e}_N(4) \end{bmatrix} \right).$$

We want the output at the last position to be the target $m(X_1, X_2)$, that is:

$$\arg \max_{u \in \{1, \dots, V\}} \mathbf{y}_{4,u} = m(X_1, X_2) \text{ for any } m, X_1, X_2.$$

The first attention block will have three attention heads.

As before, we want to extract the value of \mathbf{s}_1 that corresponds to the first token the user provided (X_1) and place it in the first half of the residual stream. We want only the third position to do that, while the rest of the positions keep the first half of their residual stream with zeros. Hence we have the following `fst` head:

$$\mathbf{W}_Q^{\text{fst}} = \left[\begin{array}{c|cccc} \mathbf{0}_{2 \times V} & 1 & 1 & 0 & 1 \\ \hline & 0 & 0 & 1 & 0 \end{array} \right], \quad \mathbf{W}_K^{\text{fst}} = \left[\begin{array}{c|cccc} \mathbf{0}_{2 \times V} & 1 & 0 & 0 & 0 \\ \hline & 0 & 1 & 0 & 0 \end{array} \right], \quad \mathbf{W}_V^{\text{fst}} = \begin{bmatrix} \mathbf{I}_V & \mathbf{0}_{V \times N} \\ \mathbf{0}_{V \times V} & \mathbf{0}_{V \times N} \end{bmatrix}.$$

The `user1` head extracts the value of the first user-provided token (X_1) and also places it in the first half of the residual stream:

$$\mathbf{W}_Q^{\text{user1}} = \left[\begin{array}{c|cccc} \mathbf{0}_{2 \times V} & 1 & 1 & 0 & 1 \\ \hline & 0 & 0 & 1 & 0 \end{array} \right], \quad \mathbf{W}_K^{\text{user1}} = \left[\begin{array}{c|cccc} \mathbf{0}_{2 \times V} & 1 & 0 & 0 & 0 \\ \hline & 0 & 0 & 1 & 0 \end{array} \right], \quad \mathbf{W}_V^{\text{user1}} = \begin{bmatrix} \mathbf{I}_V & \mathbf{0}_{V \times N} \\ \mathbf{0}_{V \times V} & \mathbf{0}_{V \times N} \end{bmatrix}.$$

And the `user2` head does the same for the value of the second user-provided token (X_2), placing it in the second half of the residual stream:

$$\mathbf{W}_Q^{\text{user2}} = \left[\begin{array}{c|cccc} \mathbf{0}_{2 \times V} & 1 & 1 & 1 & 0 \\ \hline & 0 & 0 & 0 & 1 \end{array} \right], \quad \mathbf{W}_K^{\text{user2}} = \left[\begin{array}{c|cccc} \mathbf{0}_{2 \times V} & 1 & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 & 1 \end{array} \right], \quad \mathbf{W}_V^{\text{user2}} = \begin{bmatrix} \mathbf{0}_{V \times V} & \mathbf{0}_{V \times N} \\ 2\mathbf{I}_V & \mathbf{0}_{V \times N} \end{bmatrix},$$

where the factor 2 is there because, as usual, the first linear layer will subtract 1 from everything in order to extract the value selected by the first token.

This linear layer looks as usual: $\hat{\mathcal{L}}_{\text{ex2}} = \hat{\mathcal{L}}[\mathbf{I}_{2V}, -\mathbf{1}_{2V}]$. The result is that the first V elements will be 0 except one which designates which map from second user token to output we should use, and the second V elements have a one hot-encoding of the second user token. Constructing an MLP that unpacks the mapping can become quite involved so we do not provide an explicit form for it. But from the universal approximation theorems and the finiteness of the domain and range, we know that such an MLP should exist. We thus designate by `unpack` the MLP that decodes the first half of the residual stream to:

$$\left(\frac{m_{X_1}(1)}{V}, \dots, \frac{m_{X_1}(V)}{V} \right)$$

and keeps the second half unchanged.

And now, by using two attention heads, the second attention block extracts the value of the above vector at the position designated by the second token, in a fashion not dissimilar to all the previous cases:

$$\begin{aligned} \mathbf{W}_Q^{\text{emb}} &= [\mathbf{0}_V^\top, \mathbf{1}_V^\top], & \mathbf{W}_K^{\text{emb}} &= [\mathbf{1}_V^\top, \mathbf{0}_V^\top], & \mathbf{W}_V^{\text{emb}} &= [\mathbf{I}_V & \mathbf{0}_{V \times V}], \\ \mathbf{W}_Q^{\text{user2}'} &= [\mathbf{0}_V^\top, \mathbf{1}_V^\top], & \mathbf{W}_K^{\text{user2}'} &= [\mathbf{0}_V^\top, \mathbf{1}_V^\top], & \mathbf{W}_V^{\text{user2}'} &= [\mathbf{0}_{V \times V} & \mathbf{I}_V], \end{aligned}$$

And finally, with $\hat{\mathcal{L}}_{\text{ex}} = \hat{\mathcal{L}}[\mathbf{I}_V, -\mathbf{1}_V]$, $\hat{\mathcal{L}}_{\text{sum}} = \hat{\mathcal{L}}[\mathbf{1}_V^\top, 0]$, and the same projection had as before (Equation (11)), we get the target token. The final transformer is then: $\mathcal{A}[(\mathbf{W}_Q^{\text{fst}}, \mathbf{W}_Q^{\text{user1}}, \mathbf{W}_Q^{\text{user2}}), (\mathbf{W}_K^{\text{fst}}, \mathbf{W}_K^{\text{user1}}, \mathbf{W}_K^{\text{user2}}), (\mathbf{W}_V^{\text{fst}}, \mathbf{W}_V^{\text{user1}}, \mathbf{W}_V^{\text{user2}})] \circ \hat{\mathcal{L}}_{\text{ex2}} \circ \text{unpack} \circ \mathcal{A}[(\mathbf{W}_Q^{\text{emb}}, \mathbf{W}_Q^{\text{user2}'})', (\mathbf{W}_K^{\text{emb}}, \mathbf{W}_K^{\text{user2}'})', (\mathbf{W}_V^{\text{emb}}, \mathbf{W}_V^{\text{user2}'})'] \circ \hat{\mathcal{L}}_{\text{ex}} \circ \hat{\mathcal{L}}_{\text{sum}} \circ \mathcal{L}_{\text{proj}} \circ \text{softmax}$. You can see this transformer implemented and running in practice in Section 5 of the `constructions` Jupyter notebook.

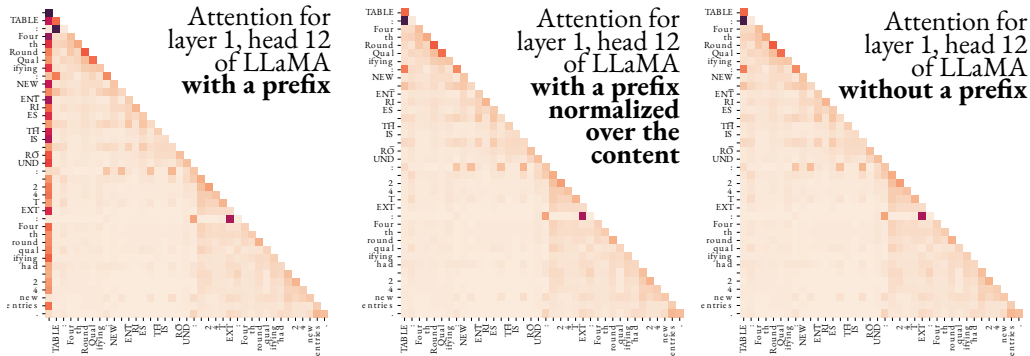


Figure 6: The attention of the twelfth head of the first layer of LLaMA (Touvron et al., 2023). The left plot shows the attention with a prefix of length one. The second plot shows the same attention but normalized such that the attention over the non-prefix positions sums to 1. The right plot shows the attention of the pre-trained model (without prefix). The center and the right plots are the same, illustrating that the presence of the prefix indeed only scales down the attention over the content (non-prefix positions) but does not change its relative distribution, providing empirical validation of Equation (6). The test sequence is TABLE: Fourth Round Qualifying : NEW_ENTRIES_THIS_ROUND : 24 TEXT: Fourth round qualifying had 24 new entries. from the DART table-to-test dataset (Nan et al., 2021).

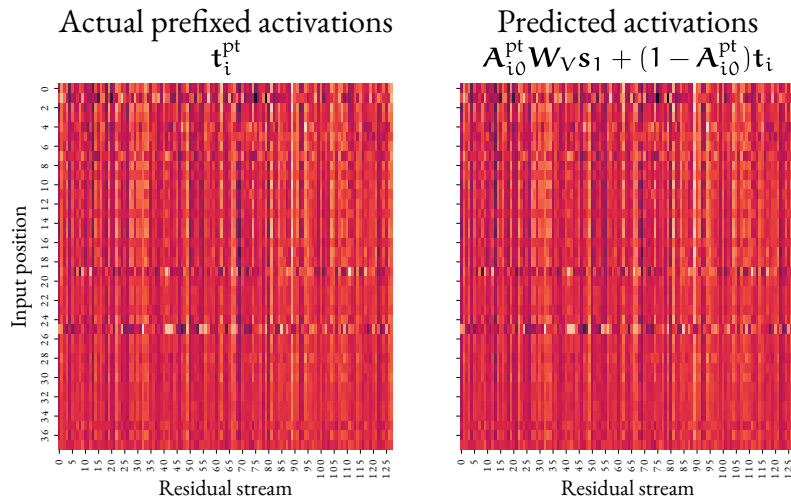


Figure 7: The activations of the twelfth head of the first layer of LLaMA (Touvron et al., 2023). The left plot shows the activations in the presence of the prefix. The right plot shows the activations t_i of the pretrained model, scaled by one minus the attention that the prefix would take and then biased in the direction $W_V s_1$. The two plots are the same, illustrating that our theory, Equation (7) in particular, also holds for real-world large transformer models. The test sequence is the same as in Figure 6.

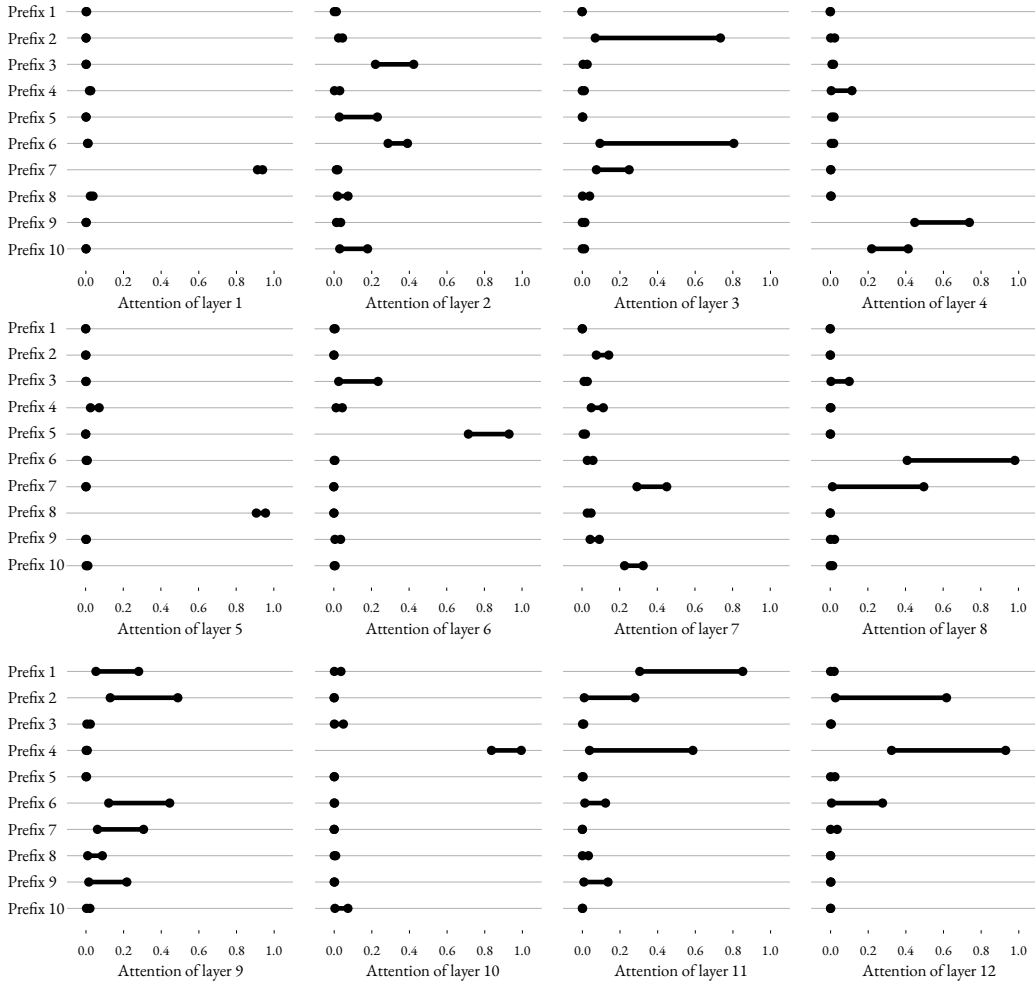


Figure 8: The range of attention (1st to 99th percentile) for a single GPT-2 (Radford et al., 2019) prefix trained on the Emotion dataset (Saravia et al., 2018). The prefix is of size 10 ($n_S = 10$). This is the attention of the last user input token (n_X) because this is the position at which the class prediction is done. For illustration purposes, we have normalized the attention so that the attention over the 10 prefix positions sums to 1. The range of attention over the 10 positions for each layer are shown.

B ATTENTION DISTRIBUTION OVER THE PREFIX

As discussed in Section 4, longer prefixes define a subspace from which the bias for the attention block is selected. For a prefix of size n_S , that means that this subspace is n_S -dimensional. Each of prefix position j defines a basis vector \mathbf{W}_{V,S_j} for this subspace, while the attention A_{i,S_j}^{pt} on this position determines how much of this basis component contributes to the bias.

In order to span the whole subspace and make full use of the capacity of the prefix, A_{i,S_j}^{pt} should vary between 0 and 1 for different inputs. However, we observe that this does not happen in practice. Figure 8 shows the ranges of attention the different prefix positions take for the GPT-2 model (Radford et al., 2019). For layer 1, for example, the attention each prefix positions gets is almost constant hence, the effective subspace is collapsed and there is a single bias vector that’s applied to the attention layer output, regardless of the user input X .

Some other layers show slightly higher variation. For example, layer 3 has three prefix positions with large variations. Therefore, the effective bias subspace is 3-dimensional and the user input X governs which bias vector from this subspace will be selected.