# **Duality Between Prompting and Weight Update in Transformers**

**Presenter:** Tianyang Hu (Huawei Noah's Ark Lab)

**Collaborators:** Brian Chen (NUS), Hui Jin (Huawei), Hwee Kuan Lee (A\*Star), Kenji Kawaguchi (NUS)

# **Duality Between Prompt and Δw --- Motivation**

Motivation:

Prompting (ICL) and Finetuning can both significantly alter the behavior of LLMs.

- Which one is preferred under what scenarios?
- Can the two be converted back-and-forth?





## Theoretical:

- Understanding how model weights are learned from data is an ultimate goal of deep learning theorists.
- Generally demystifying the duality between training data and model weights might be too challenging in the era of LLMs.
- Duality between prompt (extra data) and Δw (extra weight) may be more manageable!

# **Duality Between Prompt and Δw --- Motivation**

Motivation:

Prompting (ICL) and Finetuning can both significantly alter the behavior of LLMs.

- Which one is preferred under what scenarios?
- Can the two be converted back-and-forth?



#### Empirically From $P \rightarrow \Delta w$ :

• Model Editing: Localized and more precise SFT



• Speed up inference:



# **Duality Between Prompt and Δw --- Motivation**

Motivation:

Prompting (ICL) and Finetuning can both significantly alter the behavior of LLMs.

- Which one is preferred under what scenarios?
- Can the two be converted back-and-forth?



## Empirically From $\Delta w \rightarrow P$ :

• New tool for understanding weights of LLM



Transferability to downstream tasks:



# **Duality Between Prompt and Δw --- Results Overview**

Motivation:

Prompting (ICL) and Finetuning can both significantly alter the behavior of LLMs.

- Which one is preferred under what scenarios?
- Can the two be converted back-and-forth?



## From $P \rightarrow \Delta w$ :

- Exact conversion **impossible** without architectural modification
- We introduce an **extra bias term** in the **KV circuit** to incorporate prompts and only the attention layers need to be modified
- Exact conversion for linearized attention LLMs
- Approximate conversion for SoftMax attention

## From $\Delta w \rightarrow P$ :

• Capacity: Prompt tuning may have **universal approximation** capability, as long as the context length can be arbitrarily large.

## **Duality Between Prompt and Δw --- Notations**

**Definition 3.1** (Single-headed Attention Layer). Given input data  $X \in \mathbb{R}^{N \times d_{in}}$ , a single attention layer is characterized by trainable matrices  $W_Q, W_K \in \mathbb{R}^{d_{in} \times d_K}, W_V \in \mathbb{R}^{d_{in} \times d_V}$ . The Single-headed attention layer takes the form:

$$O(Q,K,V) = \operatorname{Attn}(Q,K,V) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V, \quad (1)$$

where O is the output,  $Q = XW_Q$ ,  $K = XW_K$  and  $V = XW_V$ .

For any given matrix Z, let  $Z_i$  return the *i*-th row as a vector and denote  $sim(\cdot, \cdot)$  as a similarity score (Note vectors are all vertical). Drawing upon the work of (Katharopoulos et al., 2020), the single-headed attention layer can be generalized as follows:

$$O_{i}^{T} = \frac{\sum_{j=1}^{N} \sin(Q_{i}, K_{j}) V_{j}^{T}}{\sum_{j=1}^{N} \sin(Q_{i}, K_{j})}.$$
(2)

**Definition 3.2** (linearized attention Layer). Given input data  $X \in \mathbb{R}^{N \times d_{in}}$ , a single attention layer is characterized by trainable matrices  $W_Q, W_K \in \mathbb{R}^{d_{in} \times d_K}, W_V \in \mathbb{R}^{d_{in} \times d_V}$ . The linearized attention layer takes the form of:

$$O_{i}^{T} = \frac{\sum_{j=1}^{N} \phi(Q_{i})^{T} \phi(K_{j}) V_{j}^{T}}{D_{1}(Q_{i})^{T} D_{2}(X)}$$

$$= \frac{\phi(Q_{i})^{T} \sum_{j=1}^{N} \phi(K_{j}) V_{j}^{T}}{D_{1}(Q_{i})^{T} D_{2}(X)},$$
(3)

where  $\phi(\cdot)$  is the associated feature map and  $D_1(Q_i)^T D_2(X)$ is a normalizing term depending on the network architecture and input X. Typically  $D_2$  is a function of  $K = XW_K$ , but we keep it as a function of X itself to capture a more general class of models. For instance, when implementing kernel softmax we would take  $D_1(Q_i) = \phi(Q_i), D_2(X) =$  $\sum_{j=1}^N \phi(K_j) = \sum_{j=1}^N \phi(X_j W_K)$  and when implementing linear attention we would take  $D_1(Q_i) = 1, D_2(X) = 1$ .

# **Naive Conversion of Prompt to Weight**

This **naïve** method has been widely explored

• Fix Prompt, generate data with various Q's;



Use the generated (Q, A) pairs for finetuning;

## Pro:

- Straightforward to implement
- Empirically works fine

## Con:

- Cumbersome: FT needs large amount of data and lots of computation resource
- Cannot achieve exact conversion
- No guarantee to generalize, especially out-ofdistribution

This **ideal** method should be: Exact, Generalizable, Fast

However, this is **impossible** without **architectural** changes

Consider the simplest **linear attention** transformer. Let X denote the input and X` denote the prompt before X. Then:

We want to find corresponding weights  $W_1, W_2$  such that the following equivalence holds true for all  $X \in \mathbb{R}^{N \times d_{in}}$ ,

 $W_1^T X^T X W_2 = W_K^T X^T X W_V + W_K^T X'^T X' W_V.$ (7)

Since X = 0 is a potential prompt, which is not true in general.



Which modifications should we make?

# **Key-Value Matrix and Addition of Extra Parameters**

**Definition 3.2** (linearized attention Layer). Given input data  $X \in \mathbb{R}^{N \times d_{in}}$ , a single attention layer is characterized by trainable matrices  $W_Q, W_K \in \mathbb{R}^{d_{in} \times d_K}, W_V \in \mathbb{R}^{d_{in} \times d_V}$ . The linearized attention layer takes the form of:

$$O_{i}^{T} = \frac{\sum_{j=1}^{N} \phi(Q_{i})^{T} \phi(K_{j}) V_{j}^{T}}{D_{1}(Q_{i})^{T} D_{2}(X)} = \frac{\phi(Q_{i})^{T} \sum_{j=1}^{N} \phi(K_{j}) V_{j}^{T}}{D_{1}(Q_{i})^{T} D_{2}(X)},$$
(3)

**Definition 3.4** (Key-Value Matrix). Given the definition of linearized attention found in equation 3, we define the Key-Value matrix *A* as:

$$A = \sum_{j=1}^{N} \phi(K_j) V_j^T.$$
(6)

We can see that the addition of extra tokens to the input will not change the  $\phi(Q_i)$  term but only the Key-Value matrix. Furthermore, the Key-Value matrix does not depend on the value of i. Hence, conversion within the Key-Value matrix A allows for conversion within the linearized attention Layer. Due to the **linearity** of the attention module, **bias** arises as a natural option.

Consider an extra bias term in the KV matrix.

L

$$\operatorname{inAttn}(Q, K, V, b_{KV}) = Q(K^T V + b_{KV}).$$

This can be seen as an architecture modification



$$LinAttn(Q, K, V, b_{KV}) = Q(K^T V + b_{KV})$$
  
= LinAttn(Q, K, V) + Qb<sub>KV</sub>

# **Key-Value Matrix and Addition of Extra Parameters**



**Theorem 4.1** (Bias conversion for Linear Attention). For weights  $W_Q, W_K, W_V$ , input prompt X, ICL prompt X' and initial bias  $b_{VK}$ , we have that:

 $\operatorname{LinAttn}([X';X]W_Q,[X';X]W_K,[X':X]W_V,b_{V,K}) = \operatorname{LinAttn}(XW_Q,XW_K,XW_V,b'_{V,K}).$ 

holds for all X if we set

 $b'_{V,K} = b_{V,K} + W_K^T X'^T X' W_V.$ 

Algorithm:

- All the attention layers needs to add this bias term
- Only attention layers needs to be modified

#### Will the extra bias term affect LLM's abilities?



#### How will numerical errors accumulate?

Table 1.	Relative	error	of	logits
				0

model size	205K	1.99M	19.8M	198M	1.98B
relative err.	2.87e-7	4.38e-7	8.27e-7	1.73e-6	4.28e-6

# **Necessary Conditions for Exact Conversion**

### Autoregressive:

Autoregressive prevents the ICL tokens (X`) from being affected by the input tokens (X).

## **Relative Positional Encoding:**

- Not a must but significantly simplifies the process.
- Using absolute positional encodings needs to include an extra transformation to correct for the positional shift.

For simplicity, we consider **RoPE** in this work.

$$O_i^T(X) = (R_{\Theta,i}^{d_K} \phi(Q_i))^T \left\{ \frac{1}{D_1(Q_i)^T D_2(X)} \sum_{j=1}^i R_{\Theta,j}^{d_K} \phi(K_j) V_j^T \right\}$$

 $R_{\Theta,i}^{d_K}$  is the rotary matrix that applies a shift by i positions.

With the inclusion of bias terms, linearized transformers with RoPE become:

$$O_{i}^{T}(X,b_{KV},b_{D}) = (R_{\Theta,i}^{d_{K}}\phi(Q_{i}))^{T} \left\{ \frac{1}{D_{1}(Q_{i})^{T}(D_{2}(X)+b_{D})} \left[\sum_{j=1}^{i} R_{\Theta,j}^{d_{K}}\phi(K_{j})V_{j}^{T}+b_{KV}\right] \right\}.$$

## Separation property of normalizing factor:

Assume that for all token sequences A and B, for the given normalizing factor  $D_2$  in equation(3) the following is true:  $D_2([A;B]) - D_2([B]) = D_2^*([A])$ , i.e., it only depends on A.

## This is not a strong condition:

- 1. Kernelized attention (Katharopoulos et al., 2020): Set  $\phi(\cdot)$  as the representation function of the kernel, set  $D_1(Q_i) = \phi(Q_i)$  and  $D_2(X) = \sum_{j=1}^M \phi(K'_i)$ .
- 2. RetNet (Sun et al., 2023): Set  $\phi(\cdot) = I(\cdot)$  be the identity function, set  $D_1(Q_i) = 1$  and  $D_2(X) = 1$ .

**Theorem 4.5.** Given model weights  $W_V, W_K, W_Q$ , input prompt X and ICL prompt X', by taking:

$$b'_{KV} = R^{d_K}_{\Theta, -M} b_{KV} + \sum_{j=1}^M R^{d_K}_{\Theta, j-M} \phi(K'_j) V'^T_j, \quad (14)$$

$$b_D' = b_D + D_2^*(X'), \tag{15}$$

then for all X:

$$O_i([X';X], b_{KV}, b_D) = O_i(X, b'_{KV}, b'_D).$$
(16)

# **Extending to SoftMax Attention**

For regular attention model, consider using dot product to approximate the SoftMax function.

Only the prompt (P) needs to be approximated



 $\begin{aligned} \operatorname{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}, \boldsymbol{b}_{KV}) &= \operatorname{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) + \boldsymbol{\phi}(\mathbf{Q})\boldsymbol{b}_{KV} \\ &= \operatorname{Softmax}(\mathbf{Q}\mathbf{K}^T)\mathbf{V} + \boldsymbol{\phi}(\mathbf{Q})\boldsymbol{b}_{KV} \\ &\approx \boldsymbol{\phi}(\mathbf{Q})\boldsymbol{\phi}(\mathbf{K}^T)\mathbf{V} + \boldsymbol{\phi}(\mathbf{Q})\boldsymbol{b}_{KV} \\ &= \boldsymbol{\phi}(\mathbf{Q})(\boldsymbol{\phi}(\mathbf{K}^T)\mathbf{V} + \boldsymbol{b}_{KV}) \end{aligned}$ 

With RoPE, the resulting output takes the form:

$$O_{i}^{T} = \frac{\left[\sum_{j=1}^{N} \sin(Q_{i}, K_{j}) V_{j}^{T}\right] + \phi(Q_{i})^{T} b_{KV}}{\left[\sum_{j=1}^{N} \sin(Q_{i}, K_{j})\right] + \phi(Q_{i})^{T} b_{D}}.$$

$$b_{KV} = \sum_{j=1}^{M} R_{\Theta,j-M}^{d_K} \phi(K'_j) V'^T_j, \quad b_D = \sum_{j=1}^{M} \phi(K'_j).$$

The accuracy of this approximate conversion depends directly on the accuracy of the linearization of SoftMax.

• Various linearization methods can all be utilized, e.g., Performer [Choromanski et al. 2020], Hedgehog [Zhang et al. 2024], etc.

$$\phi_{\mathrm{mlp}}(\boldsymbol{x}) = \left[\exp(\boldsymbol{w}_1^{\top}\boldsymbol{x} + \boldsymbol{b}), \dots, \exp(\boldsymbol{w}_d^{\top}\boldsymbol{x} + \boldsymbol{b})
ight]$$

• **New setting** for linearized attention study.

# **Experiments**

**Toy Case with linear attention + RoPE:** Inspired by [Bietti et al. 2023], we consider the following toy case with **Trigger tokens** where the sequences follow a

mixed Markov distribution.

Tokens after the trigger tokens {a, b, c, d, e} will be the same in each sequence (the first appearance will be random). Other tokens are completely random, e.g., "**a**k...**a**k...**b**H...**b**H...", "**a**M...**c**k...**a**M..."

This case is used to study induction head and ICL. The prompt to be converted can be the trigger setup: [X`, X] = ["**a**M", "...**c**k...**c**k...**a**"] and we want to converted model to output "M" after "...ck...a".

### Table 2. In context accuracy of three models.

model	In context accuracy
$M_{ m old,ICL}$	99.95%
$M_{ m old,no\ ICL}$	2.23%
$M_{ m new,no\ ICL}$	99.95%

#### **Real Data with Pretrained GPT-2:**

Table 3. Comparison of generated texts.

ICL prompt: Bill Gates was born at Seattle, Washington.

#### Model before conversion, with ICL prompt

**Input Prompt:** Bill Gates was born at Seattle, Washington. He is the CEO of

**Generated Text:** Bill Gates was born at Seattle, Washington. He is the CEO of Microsoft, the world's largest software company. He is also the co-founder and CEO of Microsoft.com, a leading online marketplace for developers and publishers. He is also the co-founder and CEO of Microsoft's Bing search engine.

#### Model after conversion, without ICL prompt

Input Prompt: He is the CEO of

**Generated Text:** He is the CEO of Microsoft, the world's largest company. He is a member of the Microsoft Foundation. He is a member of the Council of Economic Advisers.

#### Model before conversion, without ICL prompt

#### Input Prompt: He is the CEO of

**Generated Text:** He is the CEO of European Cryptocurrency Exchange, where he rapidly explores the various urban applications for cryptocurrencies and cryptocurrencies in development. The Crypto Central Service crew is looking for cash and other cryptocurrencies based on ethics or finance.

# **Duality Between Prompt and Δw --- Results Overview**

Motivation:

Prompting (ICL) and Finetuning can both significantly alter the behavior of LLMs.

- Which one is preferred under what scenarios?
- Can the two be converted back-and-forth?



## From $P \rightarrow \Delta w$ :

- Exact conversion **impossible** without architectural modification
- We introduce an **extra bias term** in the **KV circuit** to incorporate prompts and only the attention layers need to be modified
- Exact conversion for linearized attention LLMs
- Approximate conversion for SoftMax attention

## From $\Delta w \rightarrow P$ :

• Capacity: Prompt tuning may have **universal approximation** capability, as long as the context length can be arbitrarily large.

# Duality Between Prompt and Δw --- Prompt to Δw

With the previous formulation of  $\Delta w \rightarrow P$ , we can **rethink** the working mechanism of Prompt Tuning.

Prompt tuning can be seen as modifying the bias term  $\boldsymbol{b}_{KV}$  in every attention layer:

$$O_i^T = \frac{\left[\sum_{j=1}^N \sin(Q_i, K_j) V_j^T\right] + \phi(Q_i)^T b_{KV}}{\left[\sum_{j=1}^N \sin(Q_i, K_j)\right] + \phi(Q_i)^T b_D}.$$

$$b_{KV} = \sum_{j=1}^{M} R^{d_K}_{\Theta,j-M} \phi(K'_j) V'^T_j, \quad b_D = \sum_{j=1}^{M} \phi(K'_j).$$

How far can be go by modifying the bias term  $m{b}_{KV}$  ?

- $\phi(Q)b_{KV}$  can be regarded as an **MLP** with nonlinear activation  $\phi$  (with hidden dimension M)
- With M large enough, MLP can be universal approximators

#### Hypothesized arguments:

- The bias term  $m{b}_{KV}$  can dominate the attention output
- For any attention weight modification Δw and ε > 0, there exists b<sub>KV</sub>(Δw, ε) such that the attention module can be approximated within error ε.
- If any weight updates in the MLP layer can be arbitrarily approximated by weight changes in the attention layer, prompt tuning can have universal approximation capability.

# Thank you.