



# Provably Transformers Harness Multi-Concept Word Semantics for Efficient In-Context Learning

Dake Bu, Wei Huang\*, Andi Han, Atsushi Nitanda, Taiji Suzuki, Qingfu Zhang, Hau-San Wong\*

# Contents

- Motivation & Background
  - *Observed linear latent geometry of LLM*
  - *Technical limitation of current work*
- Introduction
- Problem & Model Formulation
  - *Polysemous Word Model & Concept-specific Prompt Distribution*
  - *Transformer & SGD setup*
- Main Result
  - *Exponential Convergence*
  - *OOD results*
- Experiments
- Conclusion

# Motivation & Background

Japan: Sakura; China:



Japan: Sakura; China: Peony.

The peony is often considered a national flower of China and symbolizes prosperity, honor, and beauty. It holds cultural and historical significance, much like the cherry blossom does in Japan.

Japan: Sakura; Australia: Kangaroo; China:



Japan: Sakura; Australia: Kangaroo; China: Panda.

The giant panda is widely recognized as a symbol of China, representing peace, friendship, and the country's commitment to wildlife conservation.



# Motivation & Background

Japan: Sakura; China:



Japan: Sakura; China: Peony.

The peony is often considered a national flower of China and symbolizes prosperity, honor, and beauty. It holds cultural and historical significance, much like the cherry blossom does in Japan.

Japan: Sakura; Australia: Kangaroo; China:



Japan: Sakura; Australia: Kangaroo; China: Panda.

The giant panda is widely recognized as a symbol of China, representing peace, friendship, and the country's commitment to wildlife conservation.



➤ **Observation:** Different *task concepts* are identified from the same words in different prompt.

# Motivation & Background

Japan: Sakura; China:



Japan: Sakura; China: Peony.

The peony is often considered a national flower of China and symbolizes prosperity, honor, and beauty. It holds cultural and historical significance, much like the cherry blossom does in Japan.

Japan: Sakura; Australia: Kangaroo; China:



Japan: Sakura; Australia: Kangaroo; China: Panda.

The giant panda is widely recognized as a symbol of China, representing peace, friendship, and the country's commitment to wildlife conservation.



- **Observation:** Different *task concepts* are identified from the same words in different prompt.
- **Question:** Why can an additional demo-pair influence the outcome of ICL greatly?

# Motivation & Background

## ➤ Observed Multi-Concept Latent Geometric Linearity of LLM.

Existing studies [1-4] suggest the multi-concepts are encoded linearly in the latent representation of LLM.

- Representations *within-concepts (topics)* have positive inner products
- Representations *cross-concepts (topics)* exhibit near-orthogonal relationships
- ICA is more suitable than PCA when extracting meaningful concepts

[1] Yamagiwa et al. Discovering universal geometry in embeddings with ICA. EMNLP 2023

[2] Li et al. How Do Transformers Learn Topic Structure: Towards a Mechanistic Understanding. ICML 2023

[3] Park et al. 2023: The linear representation hypothesis and the geometry of large language models. ICML 2024

[4] Jiang et. al. On the origins of linear representations in large language models. ICML 2024

[5] Reizinger et al. Position: Understanding LLMs Requires More Than Statistical Generalization. ICML 2024

# Motivation & Background

## ➤ Observed Multi-Concept Latent Geometric Linearity of LLM.

Existing studies [1-4] suggest the multi-concepts are encoded linearly in the latent representation of LLM.

- Representations *within-concepts (topics)* have positive inner products
- Representations *cross-concepts (topics)* exhibit near-orthogonal relationships
- ICA is more suitable than PCA when extracting meaningful concepts

### Essential Question

Whether and how do the observed latent geometry facilitate transformer in ICL, especially in OOD scenario?

Remark: This question is also raised as a research question **Question 5.1.4** in [5], available after our submission.

[1] Yamagiwa et al. Discovering universal geometry in embeddings with ICA. EMNLP 2023

[2] Li et al. How Do Transformers Learn Topic Structure: Towards a Mechanistic Understanding. ICML 2023

[3] Park et al. 2023: The linear representation hypothesis and the geometry of large language models. ICML 2024

[4] Jiang et. al. On the origins of linear representations in large language models. ICML 2024

[5] Reizinger et al. Position: Understanding LLMs Requires More Than Statistical Generalization. ICML 2024

# Motivation & Background

## ➤ **Observed Multi-Concept Latent Geometric Linearity of LLM.**

Existing studies [1-4] suggest the multi-concepts are encoded linearly in the latent representation of LLM.

- Representations *within-concepts (topics)* have positive inner products
- Representations *cross-concepts (topics)* exhibit near-orthogonal relationships
- ICA is more suitable than PCA when extracting meaningful concepts

## ➤ **Existing transformer theories suffer from unrealistic settings.**

- Prior theories are conducted on unrealistic settings such as linear or ReLU transformers, MLP-free attention-only models, QK-combined softmax attention and impractical loss functions like square / hinge loss.
- Due to their technical limitation, they only obtain linear or sub-linear convergence rates.

[1] Yamagiwa et al. Discovering universal geometry in embeddings with ICA. EMNLP 2023

[2] Li et al. How Do Transformers Learn Topic Structure: Towards a Mechanistic Understanding. ICML 2023

[3] Park et al. 2023: The linear representation hypothesis and the geometry of large language models. ICML 2024

[4] Jiang et. al. On the origins of linear representations in large language models. ICML 2024

[5] Reizinger et al. Position: Understanding LLMs Requires More Than Statistical Generalization. ICML 2024



# Contents

- Motivation & Background
  - *Observed linear latent geometry of LLM*
  - *Technical limitation of current work*
- **Introduction**
- Problem & Model Formulation
  - *Polysemous Word Model & Concept-specific Prompt Distribution*
  - *Transformer & SGD setup*
- Main Result
  - *Exponential Convergence*
  - *OOD results*
- Experiments
- Conclusion

# Introduction

Grounded in the studies of the LLM linear concept representation, we conduct theoretical analysis on a concept-specific sparse coding prompt distribution for ICL bi-classification tasks. Our main contributions are highlighted as below.

- We are the first to analyze the realistic setting: *softmax* attention + *ReLU* MLP transformer, which is trained using the *cross-entropy loss* via stochastic gradient descent

# Introduction

Grounded in the studies of the LLM linear concept representation, we conduct theoretical analysis on a concept-specific sparse coding prompt distribution for ICL bi-classification tasks. Our main contributions are highlighted as below.

- We are the first to analyze the realistic setting: *softmax* attention + *ReLU* MLP transformer, which is trained using the *cross-entropy loss* via stochastic gradient descent
- We are the first to showcase the *exponential 0-1 loss convergence* over the highly non-convex training dynamics in ICL theory

# Introduction

Grounded in the studies of the LLM linear concept representation, we conduct theoretical analysis on a concept-specific sparse coding prompt distribution for ICL bi-classification tasks. Our main contributions are highlighted as below.

- We are the first to analyze the realistic setting: *softmax* attention + *ReLU* MLP transformer, which is trained using the *cross-entropy loss* via stochastic gradient descent
- We are the first to showcase the *exponential 0-1 loss convergence* over the highly non-convex training dynamics in ICL theory
- We provably show that transformers can perform *certain OOD ICL tasks* by leveraging the multi-concept semantic linearity after training, highlighting their *innovative potential* for large models.

# Contents

- Motivation & Background
  - *Observed linear latent geometry of LLM*
  - *Technical limitation of current work*
- Introduction
- **Problem & Model Formulation**
  - *Polysemous Word Model & Concept-specific Prompt Distribution*
  - *Transformer & SGD setup*
- **Main Result**
  - *Exponential Convergence*
  - *OOD results*
- Experiments
- Conclusion

## Problem & Model Formulation

### ➤ Polysemous Word Model. $(\mathcal{D}_x, \mathcal{D}_y, \mathcal{D}_z, \mathcal{D}_{\xi_x}, \mathcal{D}_{\xi_y})$

Define the feature and label dictionaries:

$$\mathbf{M} = [\boldsymbol{\mu}_1^+, \boldsymbol{\mu}_1^-, \dots, \boldsymbol{\mu}_{K_1}^+, \boldsymbol{\mu}_{K_1}^-, \boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_{K_2}]$$

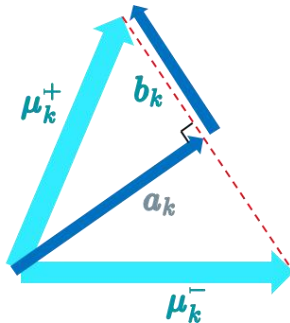
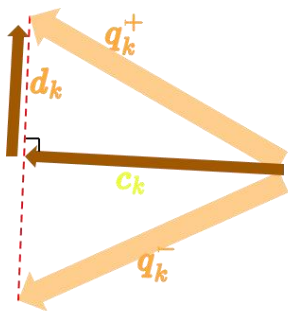
$$\mathbf{Q} = [\mathbf{q}_1^+, \mathbf{q}_1^-, \dots, \mathbf{q}_{K_1}^+, \mathbf{q}_{K_1}^-, 0, \dots, 0]$$

satisfying *within-concepts positive inner product* and *cross-concepts orthogonal* relationships.

There exists  $0 < \kappa_x, \kappa_y < 1$  such that

$$0 < \cos\langle \boldsymbol{\mu}_{k_1}^+, \boldsymbol{\mu}_{k_1}^- \rangle \leq \kappa_x, \quad 0 < \cos\langle \mathbf{q}_{k_1}^+, \mathbf{q}_{k_1}^- \rangle \leq \kappa_y$$

We can *naturally* define the high-level concept features  $\mathbf{a}_k := (\boldsymbol{\mu}_k^+ + \boldsymbol{\mu}_k^-)/2$  and the low-level semantic label features  $\mathbf{b}_k := (\boldsymbol{\mu}_k^+ - \boldsymbol{\mu}_k^-)/2$ . Also we define  $\mathbf{c}_k := (\mathbf{q}_k^+ + \mathbf{q}_k^-)/2$ ,  $\mathbf{d}_k := (\mathbf{q}_k^+ - \mathbf{q}_k^-)/2$ .



## Problem & Model Formulation

➤ **Polysemous Word Model.**  $(\mathcal{D}_x, \mathcal{D}_y, \mathcal{D}_z, \mathcal{D}_{\xi_x}, \mathcal{D}_{\xi_y})$

Then, the  $z, \xi_x, \xi_y$  are generated from  $\mathcal{D}_z$ , and Gaussian distributions  $\mathcal{D}_{\xi_x}, \mathcal{D}_{\xi_y}$  independently.

By reparameterization we define

$$\mathbf{x} := \mathbf{M}z + \xi_x \sim \mathcal{D}_x, \quad \mathbf{y} := \mathbf{Q}z + \xi_y \sim \mathcal{D}_y$$

➤ **Concept-specific Prompt Distribution.**  $(\mathcal{D}_S)$

$$\mathcal{D}_S = \sum_{k=1}^{K_1} \left( \pi_k^+ \mathcal{P}_{k,L+1}^+ + \pi_k^- \mathcal{P}_{k,L+1}^- \right)$$

$\mathbf{x}_1$	$\cdots$	$\mathbf{x}_L$	$\mathbf{x}_{L+1}$
$\mathbf{y}_1$	$\cdots$	$\mathbf{y}_L$	$\mathbf{0}$

$$= \begin{pmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_L & \mathbf{x}_{\text{query}} \\ \mathbf{y}_1 & \mathbf{y}_2 & \cdots & \mathbf{y}_L & \mathbf{0} \end{pmatrix}$$

where  $\pi_k^\pm = (2K_1)^{-1}$  denotes the equal chance of  $\mathcal{P}_{k,L+1}^\pm$ ;  $\mathcal{P}_{k,L+1}^{y_{S_n}}$  represents the  $k$ -th concept-specific prompt distribution;  $y_{S_n} \in [\pm 1]$  is the true label of a prompt. Each demo-pair  $(\mathbf{x}_l^n, \mathbf{y}_l^n)$  in  $\mathcal{P}_{k,L+1}^e$  includes either  $(\boldsymbol{\mu}_k^+, \mathbf{q}_k^+)$  or  $(\boldsymbol{\mu}_k^-, \mathbf{q}_k^-)$  with equal chance. Furthermore,  $\forall l \in [L+1]$ ,  $\mathbb{P}(z_{l, -(2k-1 \vee 2k)}^n = 1) = K^{-1}$ , indicating an equal chance of diverse task-irrelevant feature presence.

## Problem & Model Formulation

➤ **Transformer Model.**  $\Psi' := \{\mathbf{W}_Q^x, \mathbf{W}_K^x, \mathbf{W}_O^y\}$

$$\mathbf{H} = \mathbf{E}(S) = \begin{pmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_L & \mathbf{x}_{\text{query}} \\ \mathbf{y}_1 & \mathbf{y}_2 & \cdots & \mathbf{y}_L & \mathbf{0} \end{pmatrix} := (\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_{\text{query}}) \in \mathbb{R}^{(d_x + d_y) \times (L+1)}$$

$$f(\mathbf{H}; \Psi) = \mathbf{r}^\top \sigma_R(\mathbf{W}_O \text{attn}(\mathbf{H}; \Psi)),$$

$$\text{attn}(\mathbf{H}; \Psi) = \sum_{l=1}^L \mathbf{W}_V \mathbf{h}_l \sigma_S \left( (\mathbf{W}_K \mathbf{h}_l)^\top \mathbf{W}_Q \mathbf{h}_{\text{query}} \right),$$

$$\mathbf{W}_Q = \begin{pmatrix} \mathbf{W}_Q^x & * \\ * & * \end{pmatrix}, \quad \mathbf{W}_K = \begin{pmatrix} \mathbf{W}_K^x & * \\ * & * \end{pmatrix}, \quad \mathbf{W}_V = \begin{pmatrix} * & * \\ * & \mathbf{W}_V^y \end{pmatrix}, \quad \mathbf{W}_O = (* \quad \mathbf{W}_O^y),$$

where  $\mathbf{W}_Q^x, \mathbf{W}_K^x \in \mathbb{R}^{d_x \times d_x}$ ,  $\mathbf{W}_V^y \in \mathbb{R}^{(m_v - d_x) \times d_y}$ ,  $\mathbf{W}_O^y \in \mathbb{R}^{m \times d_y}$ . Here, we set the elements other than  $\mathbf{W}_Q^x, \mathbf{W}_K^x, \mathbf{W}_V^y$  and  $\mathbf{W}_O^y$  to be zero. Besides, we fix  $\mathbf{W}_V^y$  to be  $\mathbf{I}_{(m_v - d_x) \times d_y}$ . We sample  $\mathbf{r}_i$  from a uniform distribution  $\text{Unif}\{-1, 1\}$  and fixed during the training process. Based on this setting, the trainable part we need to consider is actually  $\Psi' := \{\mathbf{W}_Q^x, \mathbf{W}_K^x, \mathbf{W}_O^y\}$ . This problem remains highly non-convex and challenging.



# Problem & Model Formulation

## ➤ Stochastic Gradient Descent.

$$L_{\mathcal{B}_t}(\Psi) = L_{\mathcal{B}_t}(\Psi') := \frac{1}{B} \sum_{n \in \mathcal{B}_t} \ell(y_{S_n} \cdot f(\mathbf{H}; \Psi)) + \frac{\lambda}{2} \|\Psi'\|_F^2,$$

where  $\ell(z) = \log(1 + \exp(-z))$   $\|\Psi'\|_F^2$  represents  $\|\mathbf{W}_Q^x\|_F^2 + \|\mathbf{W}_K^x\|_F^2 + \|\mathbf{W}_O^y\|_F^2$   $\eta_t = \frac{2}{\lambda(\gamma+t)}$

**Initialization Setting.** All initial values of  $\mathbf{W}_O^y$  are sampled from a i.i.d. Gaussian distributions with mean 0 and variance  $\sigma_1^2$ . The initialization of  $\mathbf{W}_Q^x$  and  $\mathbf{W}_K^x$  are diagonal matrices  $\sigma_0 \mathbb{I}$

---

### Algorithm 1 Training algorithm

---

**Input:** Training distribution  $\mathcal{D}_S$ , Test distribution  $\mathcal{D}^*$ , Batch size  $B$ , step size  $\eta_t = \frac{2}{\lambda(\gamma+t)}$ , stopping criterion  $\varepsilon$  and total epochs  $T$ .

Initialize model parameters  $\Psi'^{(0)}$ .

**for**  $t = 0, 1, \dots, T - 1$  **do**

    If  $L_{\mathcal{D}^*}^{0-1}(\Psi^{(t)}) \leq \varepsilon$  stop else continue.

    Randomly sample mini batches  $\mathcal{B}_t$  of size  $B$  from  $\mathcal{D}_S$ .

    Update model parameters:  $\Psi'^{(t+1)} = \Psi'^{(t)} - \eta_t \nabla_{\Psi'} L_{\mathcal{B}_t}(\Psi'^{(t)})$ .

**end for**

---

# Contents

- Motivation & Background
  - *Observed linear latent geometry of LLM*
  - *Technical limitation of current work*
- Introduction
- Problem & Model Formulation
  - *Polysemous Word Model & Concept-specific Prompt Distribution*
  - *Transformer & SGD setup*
- **Main Result**
  - *Exponential Convergence*
  - *OOD results*
- Experiments
- Conclusion

## Main Result

### ➤ Exponential Convergence of 0-1 loss under low-noise condition

Theorem 1. Under Condition 1, for  $\forall \varepsilon > 0$ ,  $\exists C_1, C_2 > 0$ , with probability no less than  $1 - \delta$ , for  $T \geq \hat{T}$ , we have

$$L_{\mathcal{D}^*}^{0-1}(\Psi^{(T)}) \leq \exp\left(-\frac{C_2 \nu^2 m \lambda^2 (\gamma + T)}{K_1 \|\mathbf{q}\|^2 ((L-1)\|\mathbf{u}\|^2 + 1)}\right).$$

➤ Thus after  $T_\varepsilon = \frac{K_1 \|\mathbf{q}\|^2 ((L-1)\|\mathbf{u}\|^2 + 1)}{C_2 \nu^2 m \lambda^2} \log\left(\frac{1}{\varepsilon}\right)$  iterations, we have

$$L_{\mathcal{D}^*}^{0-1}(\Psi^{(T)}) \leq \varepsilon$$

➤ Importantly,  $\hat{T}$  is independent of  $\varepsilon$  and does not affect the convergence rate as  $\varepsilon \rightarrow 0$ .

## Main Result

### ➤ Out-of-Distribution-Generalization.

Proposition 1. Under Condition 1, for  $\forall \varepsilon > 0$ , The learned model satisfies  $L_{\mathcal{D}_S^*}^{0-1}(\Psi^{(T^*)}) \leq \varepsilon$  for  $T^* \geq T\varepsilon$ , where the  $\mathcal{D}_S^*$  can enjoy the following distribution shifts.

- The prompt length can be any positive integer.
- $\mathcal{D}_z^*$  can enjoy any shift, with each prompt sharing  $\geq 1$  co-concept, and equal chance to be  $\pm 1$ .
- $\mathcal{D}_x^* \times \mathcal{D}_y^*$  can enjoy great shift. The new  $M^*$  and  $Q^*$  satisfying that

$$\mu_k^{\pm*} = a_k^* \pm b_k^*, \quad q_k^{\pm*} = c_k^* \pm d_k^*, \quad \nu_{k_2} = \nu_{k_2}^*$$

The  $a_k^*$ ,  $b_k^*$ ,  $c_k^*$ ,  $d_k^*$ ,  $\nu_{k_2}$  are any vectors in the conic hulls of

$$\{a_k\}_{k=1}^{K_1}, \{b_k\}_{k=1}^{K_1}, \{c_k\}_{k=1}^{K_1}, \{d_k\}_{k=1}^{K_1}, \{\pm \nu_k\}_{k_2}^{K_2}$$

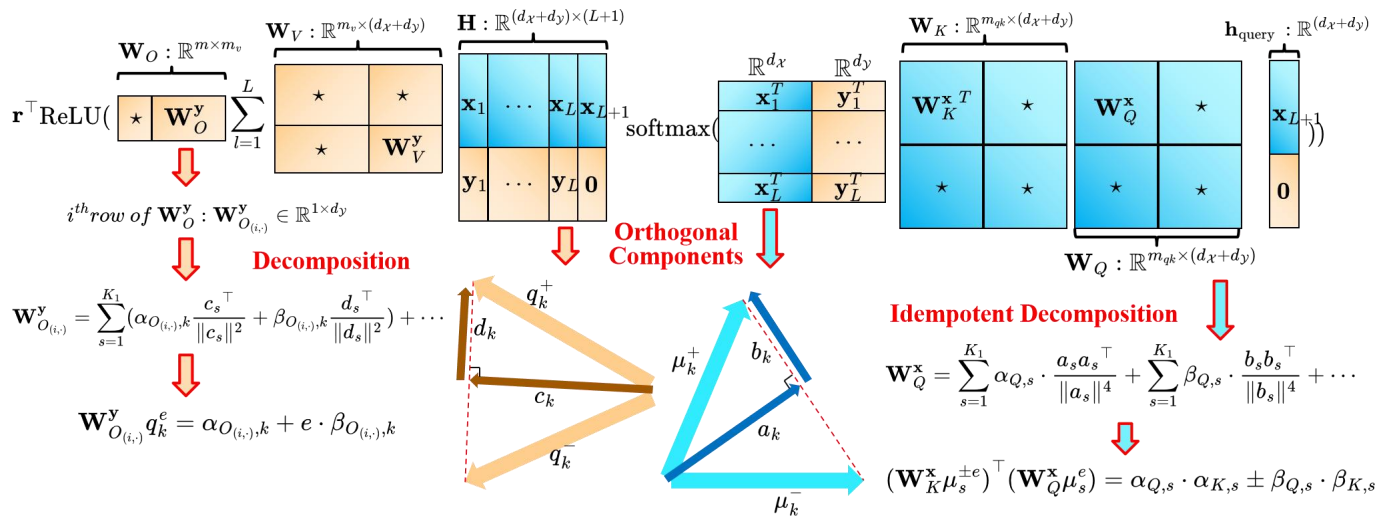
respectively.  $\|b_k^*\| \geq \|a_k^*\| = \Theta(\|\mathbf{u}\|)$ ,  $\|d_k^*\| \geq \|c_k^*\| = \Theta(\|\mathbf{q}\|)$  and  $\nu_{k_2}^* = \Theta(\|\mathbf{u}\|)$

# Main Result

## ➤ Proof Strategy: Convergence of Expectation - Exponential Variance Reduction [1]

In a big picture, we extend the standard techniques in SGD [1] to our model under **strong low-noise condition**

(i) The expected estimator would fastly converge; (ii) The variance can converge exponentially by the property of tails



With a good initialization and a symmetric low-noise prompt distribution, we can decompose the expected (over the stochastic batches) NN matrices along concept and semantic directions.

# Contents

- Motivation & Background
  - *Observed linear latent geometry of LLM*
  - *Technical limitation of current work*
- Introduction
- Problem & Model Formulation
  - *Polysemous Word Model & Concept-specific Prompt Distribution*
  - *Transformer & SGD setup*
- Main Result
  - *Exponential Convergence*
  - *OOD results*
- **Experiments**
- **Conclusion**

# Experiments

## In-Distribution Test Distribution.

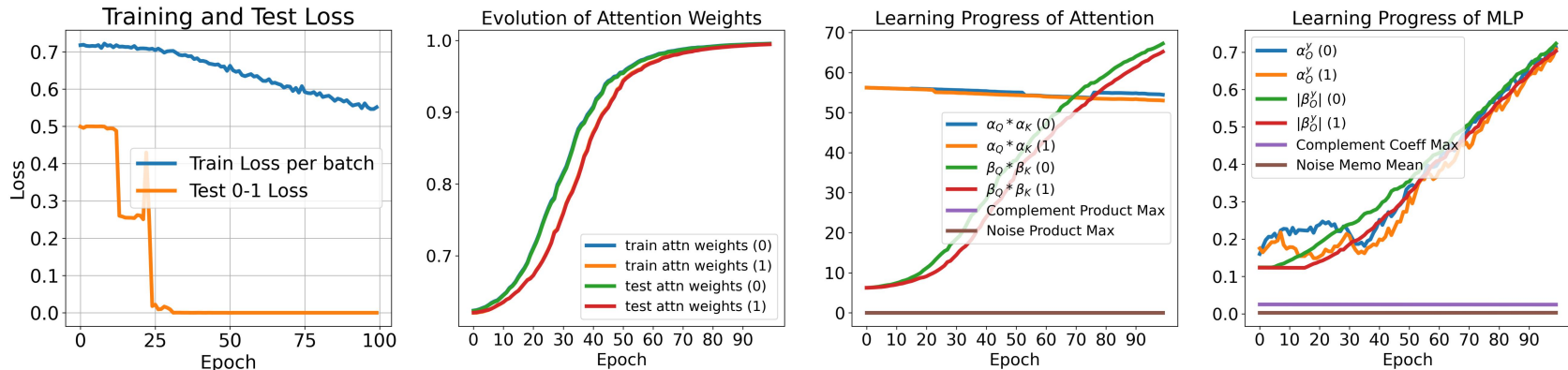
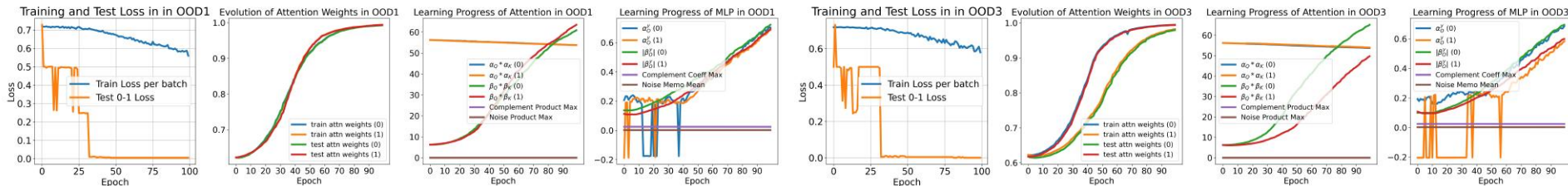


Figure 2: Learning dynamics: (i) training and test loss; (ii) correct attention weight; (iii) maximum values of  $\alpha_{Q,s} \cdot \alpha_{K,s}$ ,  $\beta_{Q,s} \cdot \beta_{K,s}$ , maximum values of the complement products  $\tau_{Q,r} \cdot \tau_{K,r}$  or  $\rho_{Q,2} \cdot \rho_{K,2}$ , and maximum values of product-with-noise  $(\mathbf{W}_K^x \xi_x)^\top \mathbf{W}_Q^x \xi_x$ ; (iv) maximum values of  $\alpha_{O(i,\cdot),k}$  and  $|\beta_{O(i,\cdot),k}|$ , maximum values of the complement coefficients  $\rho_{O(i,\cdot),w}$  and maximum values of product-with-noise  $\mathbf{W}_{O(i,\cdot)}^y \xi_y$ . The parameter settings are:  $L = 4$ ,  $K_1 = 2$ ,  $K = 104$ ,  $n_{\text{test}} = 5000$ ,  $d_{\mathcal{X}} = d_{\mathcal{Y}} = 1000$ ,  $m = 50$ ,  $\|\mathbf{u}\| = \|\mathbf{q}\| = 10$ ,  $\forall k \in [K_1]$ ,  $\langle \boldsymbol{\mu}_k^+, \boldsymbol{\mu}_k^- \rangle / \|\mathbf{u}\|^2 = \langle \mathbf{q}_k^+, \mathbf{q}_k^- \rangle / \|\mathbf{q}\|^2 = 0.5$ ,  $\sigma_0 = 0.1$ ,  $\sigma_1 = 0.01$ ,  $\sigma_\xi = 0.01$ ,  $\lambda = 0.002$ ,  $B = 16$ ,  $\gamma = 10000$ , and the total training epochs is 100.

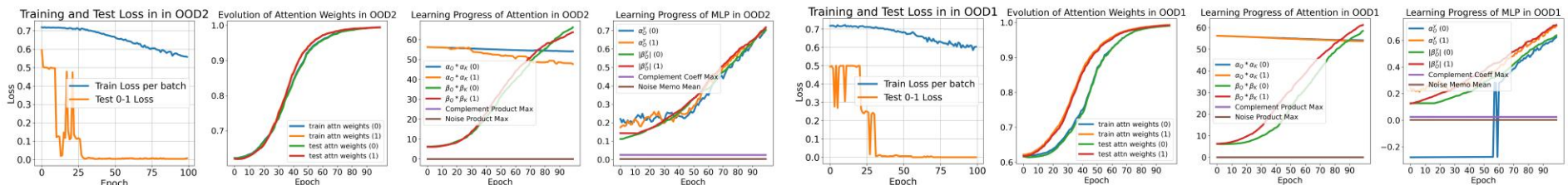
# Experiments

## OOD Test Distribution.



(a) OOD Scenario 1(i):  $L^* = 5$  during testing.

(b) OOD Scenario 1(ii):  $L^* = 2$  during testing.



(c) OOD Scenario 2: 0.8 fraction for the first and 0.2 fraction for the second concept during testing.

(d) OOD Scenario 3: Shift the data as  $\mu_1^{\pm} = \mathbf{a}_1 \pm \mathbf{b}_2$  and  $\mu_2^{\pm*} = \mathbf{a}_2 \pm \mathbf{b}_1$  during testing.

Figure 3: Learning dynamic in three OOD scenarios. The training settings and plotting methods are identical to those used in Figure 2. The consistency of the results validates Proposition 1.



# Contents

- Motivation & Background
  - *Observed linear latent geometry of LLM*
  - *Technical limitation of current work*
- Introduction
- Problem & Model Formulation
  - *Polysemous Word Model & Concept-specific Prompt Distribution*
  - *Transformer & SGD setup*
- Main Result
  - *Exponential Convergence*
  - *OOD results*
- Experiments
- **Conclusion**

## Conclusion

### ➤ **Advancing the Theory of Transformers and ICL.**

We provide a fine-grained analysis of the learning dynamics for a three-layer transformer model, comprising an **softmax** attention followed by a **ReLU**-activated feedforward network. We showcase the asymptotic properties governing the coupled learning of the attention and MLP layers.

### ➤ **Exponential Convergence of 0-1 Loss.**

Despite the highly non-convex nature of the problem, we are the first to prove an exponential convergence rate for the 0-1 loss utilizing techniques in stochastic optimization literature.

### ➤ **Connecting Multi-Concept Semantics to Efficient ICL.**

We provably show how the multi-concept encoded linear geometry of representations can enable transformer to conduct certain OOD ICL tasks.

Thanks for Listening