# PROBABILITY SIGNATURE: BRIDGING DATA SEMANTICS AND EMBEDDING STRUCTURE IN LANGUAGE MODELS

Junjie Yao<sup>1,2</sup> and Zhi-Qin John Xu<sup>1,2,\*</sup>

<sup>1</sup>Institute of Natural Sciences, MOE-LSC, Shanghai Jiao Tong University <sup>2</sup>School of Mathematical Sciences, Shanghai Jiao Tong University \*Corresponding author

#### **ABSTRACT**

The embedding space of language models is widely believed to capture the semantic relationships; for instance, embeddings of digits often exhibit an ordered structure that corresponds to their natural sequence. However, the mechanisms driving the formation of such structures remain poorly understood. In this work, we interpret the embedding structures via the data distribution. We propose a set of probability signatures that reflect the semantic relationships among tokens. Through experiments on the composite addition tasks using the linear model and feedforward network, combined with theoretical analysis of gradient flow dynamics, we reveal that these probability signatures significantly influence the embedding structures. We further generalize our analysis to large language models (LLMs) by training the Qwen2.5 architecture on the subsets of the Pile corpus. Our results show that the probability signatures are faithfully aligned with the embedding structures, particularly in capturing strong pairwise similarities among embeddings. Our work uncovers the mechanism of how data distribution guides the formation of embedding structures, establishing a novel understanding of the relationship between embedding organization and semantic patterns.

# 1 Introduction

In recent years, deep neural network-based large language models (LLMs) have demonstrated remarkable performance (Comanici et al., 2025; OpenAI et al., 2024; DeepSeek-AI et al., 2025). The development of these models has largely followed what Richard Sutton termed "the bitter lesson"—that the most effective approach to improving AI performance has historically been to leverage greater computational resources, larger models, and more data, rather than incorporating human knowledge or specialized architectures (Sutton, 2019). This trend has been formalized through scaling laws, which quantify the relationship between model performance and factors such as model size, dataset size, and computational budget through power law relationships (Kaplan et al., 2020). While these scaling laws provide valuable quantitative predictions for model performance, they also reveal a concerning limitation: the power law relationship suggests that achieving further significant improvements may require prohibitively large increases in model and data size, making continued scaling increasingly impractical and resource-intensive.

One promising approach to address these limitations is to develop a deeper understanding of the underlying mechanisms that drive transformer models' success in natural language processing (NLP). The No Free Lunch theorem establishes that no single algorithm can perform optimally across all problem domains, highlighting the fundamental importance of understanding both the characteristics of the data and the properties of the algorithms that process it (Wolpert & Macready, 1997). Recent research has made significant progress in uncovering key properties of deep learning models, including the edge of stability phenomenon (Wu et al., 2018; Cohen et al., 2021), frequency principle (Xu et al., 2020, 2025a), attention patterns (Elhage et al., 2021; Olsson et al., 2022; Bhojanapalli et al., 2020), parameter distribution properties (Kovaleva et al., 2021; Dar et al., 2023), condensation phenomenon (Luo et al., 2021; Xu et al., 2025b), and embedding structure (Cai et al., 2021). There has also been some investigation into data characteristics—such as the power-law decay of correlations between elements (like pixels) as a function of their distance (Ruderman, 1994). However, there is a significant gap in understanding how these two fundamental aspects—algorithmic properties and data characteristics—interact to produce the remarkable performance.

The embedding space, which acts as the encoder of the language tokens, therefore provides an ideal entry point for investigating how algorithmic properties and data characteristics interact. Ideally, the structure of embeddings should reflect the semantic relationships among tokens. A concrete example involves digits such as  $0, 1, 2, \ldots$ , which possess a natural ordering. Their embedding vectors accordingly display an ordered structure consistent with this numerical sequence, reflecting basic reasoning capabilities in mathematical tasks (Mikolov et al., 2013b; Ethayarajh et al., 2019; Zhang et al., 2024; Yao et al., 2025). However, the cause of the consistency between embedding structures and semantic structures is still an open question, and the driving factors of the embedding structure are still not well characterized.

In this work, we identify a set of probability signatures that encode the semantic information in the data into the structure of the embedding space of language models. Such probability signatures are constructed based on label distribution, input distribution, input/output co-occurrence distribution, etc, that systematically capture inherent token-level relationships and reflect semantic structures. This result is achieved via utilizing an embedding-based model with gradient flow analysis of embedding vectors and unembedding vectors for well-designed variable-controlled experiments. We further extend our findings to LLMs with realistic corpora, such as the Qwen2.5 architecture (Team, 2024) and subsets of the Pile dataset (Gao et al., 2020; Biderman et al., 2022). The analysis approach with the controlled experiments offers a promising methodology to uncover more and more probability signatures that can bridge the data semantics and embedding structure in language models.

#### 2 Related Work

Parameter analysis in LLMs Investigating the underlying parameter properties in LLMs is crucial for understanding the foundation of models. Some works focus on the specific modules in models. Elhage et al. (2021); Olsson et al. (2022) uncover mechanisms such as induction heads from the attention module. Bhojanapalli et al. (2020) reveals the rank-collapse phenomenon of the attention matrix. Geva et al. (2021, 2022); Dai et al. (2022) investigates the characteristics and functions of the FFN in LLMs. Additionally, analysis of a single neuron has been widely employed in mechanism interpretation, particularly in circuits analysis Hanna et al. (2023); Wang et al. (2023); Hanna et al. (2024); Wang et al. (2025), sparse autoencoders (SAE) Huben et al. (2024); Bricken et al. (2023), transcoders Dunefsky et al. (2024), and cross-layer transcoders (CLT) Ameisen et al. (2025). There are also some studies investigating the global properties of all parameters. Dar et al. (2023); Katz et al. (2024) introduce a framework for interpreting all parameters of Transformer models by projecting them into the embedding space. Kovaleva et al. (2021); Yu et al. (2025) provide an analysis of the parameter distribution, demonstrating the significance of these outliers. In this work, we will focus on the embedding space, explaining the formation of its structure from both experimental and theoretical perspectives.

Embedding structure and representation learning Since the introduction of static word embeddings by Mikolov et al. (2013a); Pennington et al. (2014) and the adoption of contextualized embeddings (Devlin et al., 2019; Peters et al., 2018), significant attention has been devoted to analyzing embedding properties. Gao et al. (2019); Ethayarajh (2019); Timkey & van Schijndel (2021) explore the anisotropy of embedding space, while Cai et al. (2021) show that embeddings exhibit isotropy within clusters. Liu et al. (2022) offers insights into grokking by emphasizing the role of well-organized embedding structures. Zhang et al. (2024) establishes a connection between embedding structure and model generalization, and Yao et al. (2025) provides an analysis of this relationship. In contrast to prior work, our study focuses on the connection between embedding structure and data properties, offering a novel insight for understanding how embeddings are organized.

# 3 Embedding-based Model

We first explain the basic notation of embedding space. Given a vocabulary  $\mathcal{V} \subset \mathbb{N}^+$  with size  $d_{\mathrm{vob}}$ . We denote a trainable matrix  $\boldsymbol{W}^E \in \mathbb{R}^{d \times d_{\mathrm{vob}}}$  as the embedding matrix for  $\mathcal{V}$ , where d is the hidden dimension. For any  $x \in \mathcal{V}$ , we denote  $\boldsymbol{e}_x \in \mathbb{R}^{d_{\mathrm{vob}}}$  as its one-hot vector and  $\boldsymbol{W}_x^E := \boldsymbol{W}^E \boldsymbol{e}_x$  as the embedding vector of x, which is intuitively the x-th column of  $\boldsymbol{W}^E$ . For a sequence  $\boldsymbol{X} = [x_1, x_2, \cdots, x_L] \subset \mathcal{V}$ , we define its one-hot representation as  $\boldsymbol{e}_{\boldsymbol{X}} := [\boldsymbol{e}_{x_1}, \boldsymbol{e}_{x_2}, \cdots, \boldsymbol{e}_{x_L}] \in \mathbb{R}^{d_{\mathrm{vob}} \times L}$  and  $\boldsymbol{W}_{\boldsymbol{X}}^E := [\boldsymbol{W}_{x_1}^E, \boldsymbol{W}_{x_2}^E, \cdots, \boldsymbol{W}_{x_L}^E] \in \mathbb{R}^{d \times L}$  as the embedding sequence of  $\boldsymbol{X}$ . Similarly,  $\boldsymbol{W}^U \in \mathbb{R}^{d_{\mathrm{vob}} \times d}$  represents the unembedding matrix and  $\boldsymbol{W}_{\nu}^U := \boldsymbol{W}^{U,T} \boldsymbol{e}_{\nu}$  (the  $\nu$ -th row of  $\boldsymbol{W}^U$ ) means the unembedding vector for any  $\nu \in \mathcal{V}$ .

We denote the models functioning on the embedding of the input sequence as embedding-based models. We provide the following formulation:

**Definition 1.** Given a vocabulary  $V \subset \mathbb{N}^+$  with size  $d_{\text{vob}}$  and a sequence  $X \in V^L$  with length L. An embedding-based model F taking X as input could be formulated as

$$F\left(\boldsymbol{X}\right) = \boldsymbol{W}^{U}G\left(\boldsymbol{W}_{\boldsymbol{X}}^{E}\right),$$

where G means the mapping in the hidden space.

Embedding-based models have been widely applied in various domains, particularly in NLP. Investigating the characteristics of the embedding space via the gradient flow dynamics is essential for understanding the embedding structures. Given a dataset  $\{(\boldsymbol{X}^i, y^i)\}_{i=1}^N$ , we utilize cross-entropy as the loss function:

$$\ell^{i} = -\log \operatorname{Softmax}\left(F\left(\boldsymbol{X}^{i}\right)\right)_{y^{i}} = -\log \frac{\exp F\left(\boldsymbol{X}^{i}\right)_{y^{i}}}{\sum_{j=1}^{d_{\operatorname{vob}}} \exp F\left(\boldsymbol{X}^{i}\right)_{j}}.$$

Denote that  $p^i = \operatorname{Softmax}(F(X^i))$ . Let  $\odot$  represent the Hadamard (element-wise) product and T mean the transpose, we obtain the following results.

**Proposition 1.** Given an embedding-based model F with an embedding matrix  $\mathbf{W}^E$ . For any token  $x \in \mathcal{V}$ , the gradient flow of  $\mathbf{W}_x^E$  (the x-th column of  $\mathbf{W}^E$ ) can be formulated as:

$$\frac{d\boldsymbol{W}_{x}^{E}}{dt} = \sum_{\nu \in \mathcal{V}} \frac{r_{x,\nu}}{N_{x,\nu}} \left( \boldsymbol{W}^{U,T} \boldsymbol{e}_{\nu} \right) \odot \sum_{i=1}^{N_{x,\nu}} G^{(1)} \left( \boldsymbol{W}_{\boldsymbol{X}_{(x,\nu)}^{i}}^{E} \right) - \frac{r_{x}^{\text{in}}}{N_{x}^{\text{in}}} \sum_{i=1}^{N_{x}^{\text{in}}} G^{(1)} \left( \boldsymbol{W}_{\boldsymbol{X}_{x}^{i}}^{E} \right) \odot \left( \boldsymbol{W}^{U,T} \boldsymbol{p}_{x}^{i} \right),$$

where  $N_x^{\text{in}}$ ,  $N_{x,\nu}$  denotes the count of sequences containing x and the count of sequences containing x with label  $\nu$ ,  $r_x^{\text{in}} = \frac{N_x^{\text{in}}}{N}$ ,  $r_{x,\nu} = \frac{N_{s,\nu}}{N}$ ,  $G^{(1)}$  represents the derivative of G with respect to  $\mathbf{W}_x^E$ .  $\mathbf{X}_x^i$  is the i-th sequence containing x and  $\mathbf{X}_{(x,\nu)}^i$  denotes the i-th sequence containing x with label  $\nu$ .  $\mathbf{p}_x^i$  means the output probability of the model over  $\mathbf{X}_x^i$ .

**Proposition 2.** Given an embedding-based model F with an unembedding matrix  $\mathbf{W}^U$ . For any token  $\nu \in \mathcal{V}$ , the gradient flow of  $\mathbf{W}^U_{\nu}$  (the  $\nu$ -th row of  $\mathbf{W}^U$ ) can be written as

$$\frac{d\boldsymbol{W}_{\nu}^{U}}{dt} = \frac{r_{\nu}^{\text{out}}}{N_{\nu}^{\text{out}}} \sum_{i=1}^{N_{\nu}^{\text{out}}} \left[ G\left(\boldsymbol{W}_{\boldsymbol{X}_{(\cdot,\nu)}^{E}}^{E}\right) \right]^{T} - \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{p}^{i,\nu} \left[ G\left(\boldsymbol{W}_{\boldsymbol{X}^{i}}^{E}\right) \right]^{T},$$

where  $N_{\nu}^{\text{out}}$  denotes the count of sequences with label  $\nu$  and  $r_{\nu_j}^{\text{out}} = \frac{N_{\nu}^{\text{out}}}{N}$ .  $\boldsymbol{X}_{(\cdot,\nu)}^i$  means the i-th sample which takes  $\nu$  as the label.  $\boldsymbol{p}^{i,\nu}$  means the  $\nu$ -th element of  $\boldsymbol{p}^i$ .

In this work, we employ three embedding-based architectures:

- Linear model.  $F_{\text{lin}}(\boldsymbol{X}) = \boldsymbol{W}^U \sum_{x \in \boldsymbol{X}} \boldsymbol{W}_x^E$ .
- Feedforward network.  $F_{\text{ffn}}(\boldsymbol{X}) = \boldsymbol{W}^U \sigma\left(\sum_{x \in \boldsymbol{X}} \boldsymbol{W}_x^E\right)$ , where  $\sigma$  denotes the element-wise nonlinear activation.
- Transformer-based architecture. We employ the Qwen2.5 architecture in Section 6 and the Llama 2 architecture (Touvron et al., 2023) in Appendix C.4.

#### 4 Probability Signature

In the field of deep learning, the data characteristics play a critically important role in both the training dynamics and the final performance of the model. Given a training dataset  $\{(\boldsymbol{X}^i, y^i)\}_{i=1}^N$  sampled from distribution  $\pi$ , we define the following *probability signatures*:

$$\begin{split} \phi_{x}^{y} &= \sum_{\nu \in \mathcal{V}} \mathbb{P}_{\pi} \left( y = \nu \mid x \in \boldsymbol{X} \right) \boldsymbol{e}_{\nu}, \quad \phi_{x}^{\boldsymbol{X}} = \sum_{x' \in \mathcal{V}} \mathbb{P}_{\pi} \left( x' \in \boldsymbol{X} \mid x \in \boldsymbol{X} \right) \boldsymbol{e}_{x'}, \\ \phi_{x}^{\boldsymbol{X} \mid y} &= \sum_{\nu \in \mathcal{V}} \boldsymbol{e}_{\nu} \left( \sum_{x' \in \mathcal{V}} \mathbb{P} \left( x' \in \boldsymbol{X} \mid x \in \boldsymbol{X}, y = \nu \right) \boldsymbol{e}_{x'} \right)^{T}, \quad \boldsymbol{\varphi}_{\nu}^{\boldsymbol{X}} = \sum_{x \in \mathcal{V}} \mathbb{P}_{\pi} \left( x \in \boldsymbol{X} \mid y = \nu \right) \boldsymbol{e}_{x}. \end{split}$$

 $\phi_x^y$  denotes the distribution of the label given that the input sequence X contains element x. It captures the association between a specific input element x and its label.  $\phi_x^X$  represents the probability that, given X contains x, it also contains another input element x'. It characterizes the co-occurrence relationship between different elements.  $\phi_x^{X|y}$  indicates the probability that, when X contains x and the label is fixed, the sequence additionally includes x'. It further delineates the relationship among different elements in the sequence under a specified label condition.  $\varphi_x^X$  denotes the probability distribution of the element x contained in x conditional with the label being x. It reflects the dependency between input elements and the label from a reverse perspective.

# 5 Addition Task

The addition task has become an important benchmark for studying the characteristics of language models. Many studies have found that the digits' embeddings exhibit an ordered structure consistent with the natural sequence of numbers (Zhang et al., 2024; Yao et al., 2025). To demonstrate the significant influence of probability signatures on the model's embedding space, we design three types of composite addition tasks to perform *variable-controlled experiments*. Assuming all tokens belong to positive integers, and we denote an anchor set by  $\mathcal{A}$ , whose elements represent different addition operations, i.e., anchor  $\alpha_1$  means addition with  $\alpha_1$ . Given a input sequence  $\mathbf{X} = [z, \alpha_1, \alpha_2]$ , which means the composite function  $(\alpha_1, \alpha_2)$  on the key z, we define the following tasks:

- Addition task.  $f_{\mathrm{add}}(X) = z + \alpha_1 + \alpha_2$ ,  $\alpha_1, \alpha_2 \in \mathcal{A}$ . For each anchor pair  $(\alpha_1, \alpha_2)$ , z is sampled from the same set  $\mathcal{Z}$  with  $\mathcal{Z} \cap \mathcal{A} = \emptyset$ . In  $f_{\mathrm{add}}, \phi_{\alpha}^y$  are distinct with varying anchor  $\alpha$  while  $\phi_{\alpha}^X$  are identical across anchors.
- Addition task with the same value domain.  $\tilde{f}_{add}(X) = z + \alpha_1 + \alpha_2$ ,  $\alpha_1, \alpha_2 \in \mathcal{A}$ . For anchor pair  $(\alpha_1, \alpha_2)$ ,  $z \in \mathcal{Z}_{(\alpha_1, \alpha_2)} = \mathcal{Y} \alpha_1 \alpha_2$  where  $\mathcal{Y}$  denotes the label domain, which is identical for all anchor pairs. In  $\tilde{f}_{add}$ ,  $\phi_{\alpha}^{\mathbf{X}}$  are distinct with varying anchor  $\alpha$  while  $\phi_{\alpha}^{y}$  are identical for all  $\alpha \in \mathcal{A}$ .
- Module addition.  $f_{\mathrm{mod}}(X) = \min \mathcal{Z} + (z + \alpha_1 + \alpha_2 \mod |\mathcal{Z}|)$ ,  $\alpha_1, \alpha_2 \in \mathcal{A}$  and  $z \in \mathcal{Z}$ . Both  $\phi_{\alpha}^X$  and  $\phi_{\alpha}^y$  are identical with different anchors, while  $\phi_{\alpha}^{X|y}$  are distinct.

In this work, we set  $\mathcal{A}=\{11,12,\cdots,20\}$  and  $\mathcal{Y}=\mathcal{Z}=\{101,102,\cdots,140\}$ . Figure 1 A displays the probability signature, which is distinct across the  $\alpha$  in each task, revealing that the difference among  $\alpha$  lies in the global horizontal shift. The detailed formulations are provided in Appendix B.1.

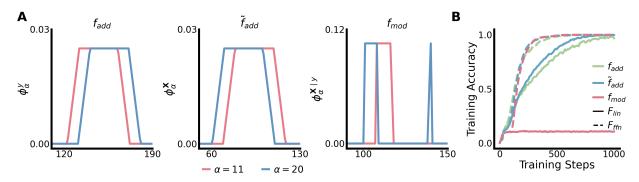


Figure 1: A: Probability signature which is distinct with varying  $\alpha$  in each task ( $f_{\rm add} \to \phi_{\alpha}^y$ ,  $\tilde{f}_{\rm add} \to \phi_{\alpha}^X$ ,  $f_{\rm mod} \to \phi_{\alpha}^{X|y}$ ),  $\alpha=11$  (red) and 20 (blue). In  $\phi_{\alpha}^{X|y}$ , it's displayed with y=160. B: Training accuracy of the  $F_{\rm lin}$  and  $F_{\rm ffn}$  on the three addition tasks.

For notation convenience, we denote that  $\boldsymbol{W}_{\mathcal{A}}^{E} = \left[\boldsymbol{W}_{\alpha}^{E}\right]_{\alpha\in\mathcal{A}}, \phi_{\mathcal{A}} = \left[\phi_{\alpha}\right]_{\alpha\in\mathcal{A}}$  for  $\phi_{\alpha} = \phi_{\alpha}^{y}, \phi_{\alpha}^{\boldsymbol{X}}, \phi_{\alpha}^{\boldsymbol{X}|y}$ , and  $\cos\left(\boldsymbol{W}_{\mathcal{A}}^{E}\right) := \left[\cos\left(\boldsymbol{W}_{\alpha}^{E}, \boldsymbol{W}_{\alpha'}^{E}\right)\right]_{\alpha,\alpha'\in\mathcal{A}}, \cos\left(\phi_{\mathcal{A}}\right) := \left[\cos\left(\phi_{\alpha},\phi_{\alpha'}\right)\right]_{\alpha,\alpha'\in\mathcal{A}}.$  Similarly,  $\cos\left(\boldsymbol{W}_{\mathcal{V}}^{U}\right) := \left[\cos\left(\boldsymbol{W}_{\mathcal{V}}^{U}, \boldsymbol{W}_{\nu'}^{U}\right)\right]_{\nu,\nu'\in\mathcal{V}}$  and  $\cos\left(\boldsymbol{\varphi}_{\mathcal{V}}^{\boldsymbol{X}}\right) := \left[\cos\left(\varphi_{\nu},\varphi_{\nu'}\right)\right]_{\nu,\nu'\in\mathcal{V}}.$ 

#### 5.1 Results

We train these addition tasks using the  $F_{\rm lin}$  and  $F_{\rm ffn}$  with d=200. Inspired by the work of Luo et al. (2021); Xu et al. (2025b), we initialize the model parameters by  $W_{i,j} \sim \mathcal{N}\left(0,d^{-0.8}\right)$ , indicating a small initialization scale. The

complete training configurations are provided in Appendix A. Figure 1 B shows the training accuracy of  $F_{\rm lin}$  and  $F_{\rm ffn}$  on the three addition tasks. The results reveal that both  $f_{\rm add}$  and  $\tilde{f}_{\rm add}$  are learned well by the linear model, whereas  $f_{\rm mod}$  requires the nonlinear model to achieve an effective fit.

## 5.2 Embedding Matrix

In the addition tasks, the anchors exhibit a strict ordering due to the numerical sequence. This provides an ideal setting for the embedding space to develop a corresponding ordered relationship. To formally quantify the formation of the ordered structure, we define the following metric:

$$R_{\text{order}}\left(\boldsymbol{W}_{\mathcal{A}}^{E}\right) = \text{Corr}\left(\cos\left(\boldsymbol{W}_{\mathcal{A}}^{E}\right), \left\{\mid \alpha - \alpha'\mid\right\}_{\alpha, \alpha' \in \mathcal{A}}\right).$$

 $R_{\mathrm{order}}\left(oldsymbol{W}_{\mathcal{A}}^{E}\right)$  reflects the relationship between embedding similarity and anchor difference. A strong negative  $R_{\mathrm{order}}\left(oldsymbol{W}_{\mathcal{A}}^{E}\right)$  (approximately -1) indicates that the similarity decreases systematically with increasing anchor difference, confirming the presence of a hierarchical organization in the anchor embeddings. Figure 2 A represents the distribution of  $\cos\left(oldsymbol{W}_{\mathcal{A}}^{E}\right)$  for the three tasks with  $F_{\mathrm{lin}}$  and  $F_{\mathrm{ffn}}$ , respectively, and Figure 2 B depicts the corresponding evolution of  $R_{\mathrm{order}}\left(oldsymbol{W}_{\mathcal{A}}^{E}\right)$ . In the case of  $f_{\mathrm{add}}$ , anchor embeddings quickly form an ordered structure, where the cosine similarity gets smaller as the anchor distance gets larger. For the task  $\tilde{f}_{\mathrm{add}}$ , the anchor embeddings also develop a similar hierarchical structure. However, its construction requires more steps, indicating that the driving factors of the structure in  $f_{\mathrm{add}}$  and  $\tilde{f}_{\mathrm{add}}$  are different. In  $f_{\mathrm{mod}}$ , although the linear model fails to learn it effectively, the anchor embeddings still undergo noticeable changes from the initial stage. Specifically, all embedding vectors become aligned in nearly the same direction. Furthermore, the anchor embeddings of  $f_{\mathrm{mod}}$  in  $F_{\mathrm{ffn}}$  construct an ordered structure with more steps, suggesting that the activation provides another factor in deriving such an embedding structure.

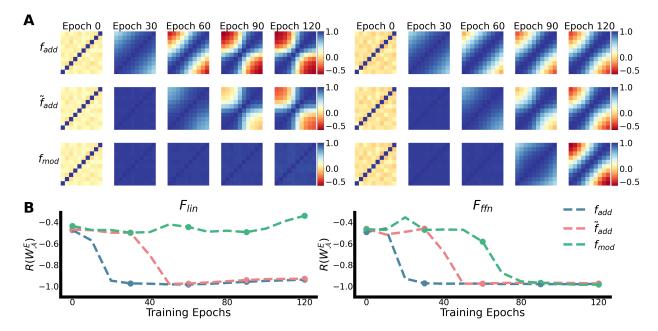


Figure 2: A: The heatmap of  $\cos\left(\boldsymbol{W}_{\mathcal{A}}^{E}\right)$  in  $F_{\mathrm{lin}}$  (left) and  $F_{\mathrm{ffn}}$  (right) during the training process. Each row corresponds to  $f_{\mathrm{add}}$ ,  $\tilde{f}_{\mathrm{add}}$ , and  $f_{\mathrm{mod}}$ , respectively. B: Dynamics of  $R_{\mathrm{order}}\left(\boldsymbol{W}_{\mathcal{A}}^{E}\right)$  in  $F_{\mathrm{lin}}$  (left) and  $F_{\mathrm{ffn}}$  (right). Line colors represent task types.

We derive the reasons for the embedding structure in each task from the gradient flow. With Proposition 1, we obtain the following approximation:

**Corollary 1** (Embedding of Linear Model). Let  $N \to \infty$ ,  $\pi$  denotes the data distribution over the training set. The gradient flow of  $W_{\alpha}^{E}$  in  $F_{\text{lin}}$  can be approximated by

$$\frac{d\mathbf{W}_{\alpha}^{E}}{dt} = \mathbf{W}^{U,T} r_{x}^{\text{in}} \left( \boldsymbol{\phi}_{\alpha}^{y} - \frac{1}{d_{\text{vob}}} \mathbf{W}^{U} \mathbf{W}^{E} \boldsymbol{\phi}_{\alpha}^{X} + \boldsymbol{\eta} \right), \tag{1}$$

where  $\eta$  denotes the data-independent and higher-order terms.

**Remark 1.** Deep learning methods fundamentally comprise three essential components: the model, the data, and the optimization algorithm. Corollary 1 clearly illustrates how these three elements are coupled (model:  $F_{\rm lin}$ ; data: probability signatures; optimization algorithm: the gradient flow) and influence the embedding space, providing crucial insights for understanding and investigating their joint impact on the performance of deep learning approaches.

Corollary 1 indicates that the dynamics of  $W_{\alpha}^{E}$  in  $F_{\text{lin}}$  are primarily impacted by the probability signatures  $\phi_{\alpha}^{y}$  and  $\phi_{\alpha}^{X}$ , demonstrating the connection between data distribution and the embedding space. As we mentioned, the  $\phi_{\alpha}^{y}$  is distinct for different  $\alpha$ , while the  $\phi_{\alpha}^{X}$  is identical for all  $\alpha$  in  $f_{\text{add}}$ ; the opposite holds for  $\tilde{f}_{\text{add}}$ . Figure 3 A depicts  $\cos\left(\phi_{\mathcal{A}}^{y}\right)$  (left) and  $\cos\left(\phi_{\mathcal{A}}^{X}\right)$  (middle column) in  $f_{\text{add}}$  (top) and  $\tilde{f}_{\text{add}}$  (middle row), revealing that  $\phi_{\alpha}^{y}$  and  $\phi_{\alpha}^{X}$  are significant for the formation of the hierarchy embedding structure in  $f_{\text{add}}$  and  $\tilde{f}_{\text{add}}$ , respectively. Furthermore, (1) indicates that  $\phi_{\alpha}^{y}$  acts as a leading term and the effect of  $\phi_{\alpha}^{X}$  is weaker in the early training process since it times a small magnitude term  $\frac{1}{d_{\text{vob}}}W^{U}W^{E}$ . This results in the formation speed of the structure in  $\tilde{f}_{\text{add}}$  being slower than  $f_{\text{add}}$ , which is consistent with the phenomenon in Figure 2.

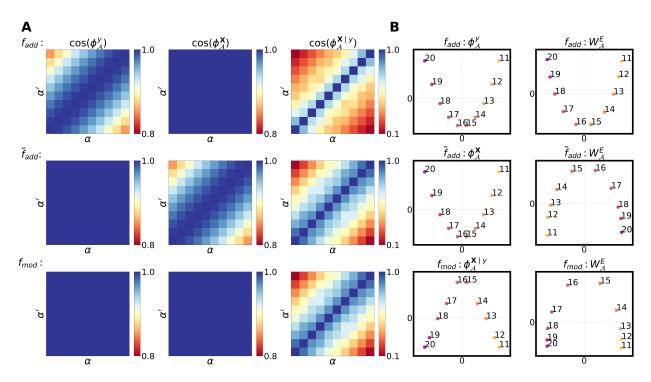


Figure 3: A: Cosine similarity among different anchor  $\alpha$  of  $\phi_{\alpha}^{y}$ ,  $\phi_{\alpha}^{X}$ ,  $\phi_{\alpha}^{X|y}$  (see (4)) in each task. B: The PCA projection of the key factors  $(f_{\text{add}} \to \phi_{\alpha}^{y}, \tilde{f}_{\text{add}} \to \phi_{\alpha}^{X}, f_{\text{mod}} \to \phi_{\alpha}^{X|y})$  and the embedding vectors in different tasks  $(F_{\text{fin}}, 120 \text{ epoch})$ .

In task  $f_{\rm mod}$ ,  $\phi_{\alpha}^y$  and  $\phi_{\alpha}^X$  are both identical across different anchors  $\alpha$ . Figure 3 A (bottom) indicates that  $\cos{(\phi_{\mathcal{A}}^y)}$  and  $\cos{(\phi_{\mathcal{A}}^X)}$  are 1 for all anchor pairs, which leads these embedding vectors to converge to almost the same direction, consistent with the observation in  $F_{\rm lin}$ . To identify the key factors that contribute to the formation of the ordered embedding structure for  $f_{\rm mod}$  in  $F_{\rm fin}$ , we perform a similar analysis of its gradient flow and obtain the following result.

**Corollary 2** (Embedding of FFN). Let  $N \to \infty$ ,  $\pi$  denotes the data distribution over the training set. The gradient flow of  $W_{\alpha}^{E}$  in  $F_{\mathrm{ffn}}$  could be approximated by

$$\frac{d\mathbf{W}_{\alpha}^{E}}{dt} = \mathbb{T} \cdot \left(\boldsymbol{\phi}_{\alpha}^{\mathbf{X}|y}\right)^{T} + \boldsymbol{\eta}_{\boldsymbol{\phi}_{\alpha}^{y}} + \frac{1}{d_{\text{vob}}} \boldsymbol{\eta}_{\boldsymbol{\phi}_{\alpha}^{\mathbf{X}}} + \tilde{\boldsymbol{\eta}}, \tag{2}$$

where  $\mathbb{T} \in \mathbb{R}^{d \times d_{\text{vob}} \times d_{\text{vob}}}$ ,  $\mathbb{T}_{:,:,\nu} = r_{\alpha,\nu} \text{diag}\left(\boldsymbol{W}_{\nu}^{U}\right) \boldsymbol{W}^{E}$  for  $\nu \in \mathcal{V}$  and 0 otherwise.  $\boldsymbol{\eta}_{\boldsymbol{\phi}_{\alpha}^{u}}$  and  $\boldsymbol{\eta}_{\boldsymbol{\phi}_{\alpha}^{x}}$  denotes the term related with  $\boldsymbol{\phi}_{\alpha}^{y}$  and  $\boldsymbol{\phi}_{\alpha}^{x}$ , respectively.  $\tilde{\boldsymbol{\eta}}$  represents the higher-order term.

Corollary 2 indicates that the ordered embedding structure of  $f_{\rm mod}$  primarily relies on probability signature  $\phi_{\alpha}^{X|y}$ . Figure 3 A (right) depicts the  $\cos\left(\phi_{\mathcal{A}}^{X|y}\right)$  in tasks, which reveals that  $\phi_{\alpha}^{X|y}$  in  $f_{\rm mod}$  constructs an ordered structure, resulting in the ordered embedding structure in  $F_{\rm ffn}$ . Furthermore, Figure 3 C shows the PCA projection on probability signatures and the embedding space in  $F_{\rm ffn}$ , revealing a high consistency. This comparison demonstrates the impact of the probability signatures in shaping the embedding space.

#### 5.3 Unembedding Matrix

The *i*-th row of the unembedding matrix can also be viewed as the feature for the *i*-th token. As shown in Figure 4 A, a similar ordered structure emerges among the unembedding vectors with the label index in  $F_{\rm lin}$  across all tasks. Specifically,  $\mathbf{W}_{\nu}^{U}$  in  $f_{\rm mod}$  constructs a ring where the similarity between small  $\nu$  and large  $\nu'$  is also large since  $\mathcal{Z}_{\rm max}+1=\mathcal{Z}_{\rm min}$  in  $f_{\rm mod}$ .

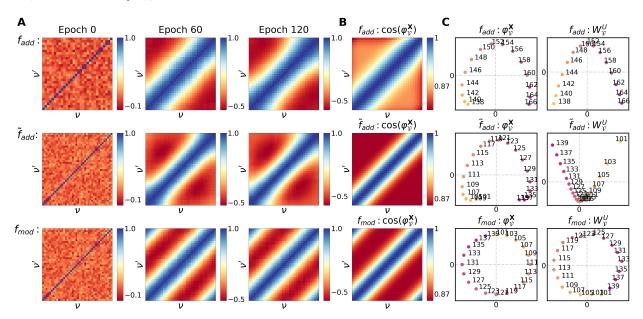


Figure 4: A: The heatmap of the  $\cos\left(\boldsymbol{W}_{\mathcal{V}}^{U}\right)$  with label index in  $F_{\text{lin}}$  during the training process. B: The heatmap of  $\cos\left(\boldsymbol{\varphi}_{\mathcal{V}}^{\boldsymbol{X}}\right)$  across different tasks. C: PCA projection of  $\boldsymbol{\varphi}_{\mathcal{V}}^{\boldsymbol{X}}$  and  $\boldsymbol{W}_{\mathcal{V}}^{U}$  (epoch 120).

Similarly, we identify the driving factors of this specific structure by examining the gradient flow of  $W^U$ . Since this phenomenon occurs in both  $F_{\text{lin}}$  and  $F_{\text{ffn}}$ , it suffices to analyze the linear model. Based on Proposition 2, we derive the following result:

**Corollary 3** (Unembedding of Linear Model). Let  $N \to \infty$ ,  $\pi$  denotes the data distribution over the training set. The gradient flow of  $W^U_{\nu}$  in  $F_{\rm lin}$  could be approximated by

$$\frac{d\mathbf{W}_{\nu}^{U}}{dt} = Lr_{\nu}^{\text{out}} \left( \mathbf{W}^{E} \boldsymbol{\varphi}_{\nu}^{\mathbf{X}} \right)^{T} + \boldsymbol{\eta}, \tag{3}$$

where  $\eta$  denotes the output term.

Corollary 3 illustrates that  $\varphi^X_{\nu}$  plays a significant role in shaping the unembedding matrix. Figure 4 B depicts the distribution of  $\cos\left(\varphi^X_{\mathcal{V}}\right)$ , which is aligned with the distribution of the  $\cos\left(W^U_{\mathcal{V}}\right)$  in Figure 4 A. Furthermore, Figure 4 C compares the PCA projection of  $\varphi^X_{\mathcal{V}}$  and  $W^U_{\mathcal{V}}$  in all tasks, revealing a high consistency and validating our analysis.

## 6 Language Model

We have demonstrated the influence of data distribution on the embedding space in the addition tasks. In this section, we explore how to extend this analysis to real-world language models. Most contemporary language models are built upon the Transformer decoder architecture. Assuming the input sequence is denoted as X with length L, we define a

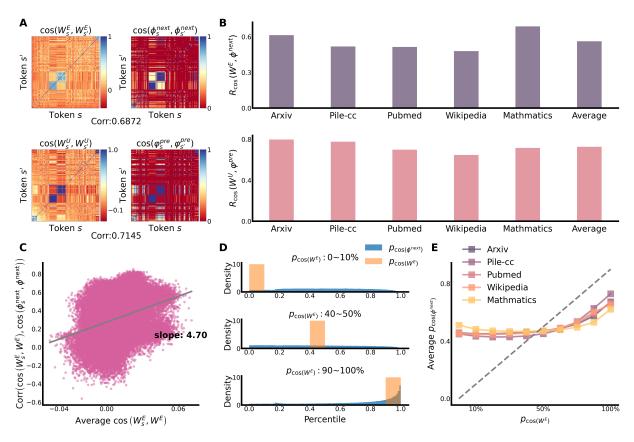


Figure 5: A: Heatmap of the cosine similarity of  $\boldsymbol{W}^E, \boldsymbol{W}^U, \phi^{\text{next}}$  and  $\varphi^{\text{pre}}$ . B:  $R_{\cos}\left(\boldsymbol{W}^E, \phi^{\text{next}}\right)$  (top) and  $R_{\cos}\left(\boldsymbol{W}^U, \varphi^{\text{pre}}\right)$  (bottom) with different datasets. C: Relation between  $\operatorname{Corr}\left(\cos\left(\boldsymbol{W}_s^E, \boldsymbol{W}^E\right), \cos\left(\phi_s^{\text{next}}, \phi^{\text{next}}\right)\right)$  and the average value of  $\cos\left(\boldsymbol{W}_s^E, \boldsymbol{W}^E\right)$ . Each point denotes a token s. D: Distribution of  $p_{\cos(\phi^{\text{next}})}$ , conditioned on intervals  $0 \sim 10\%, 40 \sim 50\%$  and  $90 \sim 100\%$  of the  $p_{\cos(\boldsymbol{W}^E)}$ . E: Average value of  $p_{\cos(\phi^{\text{next}})}$  within each interval of  $p_{\cos(\boldsymbol{W}^E)}$ .

language model  $F_{lan}$  as:

$$F_{\mathrm{lan}}\left(oldsymbol{X}
ight) = oldsymbol{W}^{U}\left(oldsymbol{W}_{oldsymbol{X}}^{E} + \tilde{F}\left(oldsymbol{X}
ight)
ight)$$

Given the training corpus  $\{X^i\}_{i=1}^N$ , we define the following probability signatures for any  $s \in \mathcal{V}$ :

$$\phi_{s}^{\text{next}} = \sum_{s' \in \mathcal{V}} \mathbb{P}_{\pi} \left( \bigcup_{t=1}^{L-1} \left\{ X_{t+1} = s' \mid X_{t} = s \right\} \right) e_{s'},$$

$$\varphi_{s}^{\text{pre}} = \sum_{s' \in \mathcal{V}} \mathbb{P}_{\pi} \left( \bigcup_{t=1}^{L-1} \left\{ X_{t} = s' \mid X_{t+1} = s \right\} \right) e_{s'},$$
(4)

and  $\phi^{\mathrm{next}} = [\phi^{\mathrm{next}}_s]_{s \in \mathcal{V}}$ ,  $\varphi^{\mathrm{pre}} = [\phi^{\mathrm{pre}}_s]_{s \in \mathcal{V}}$ . We derive the following result:

**Corollary 4.** Let  $N \to \infty$ ,  $\pi$  denotes the token distribution in the training dataset. The gradient flow of the embedding vector  $W_s^E$  of token s could be fomulated as

$$\frac{d\mathbf{W}_s^E}{dt} = r_s^{\text{in}} \mathbf{W}^{U,T} \boldsymbol{\phi}_s^{\text{next}} + \boldsymbol{\eta}^E.$$

Furthermore, the gradient flow of the unembedding vector  $\mathbf{W}_s^U$  could be approximated as

$$\frac{d\mathbf{W}_{s}^{U}}{dt} = r_{s}^{\text{out}} \left( \mathbf{W}^{E} \boldsymbol{\varphi}_{s}^{\text{pre}} \right)^{T} + \boldsymbol{\eta}^{U}.$$

The  $\eta^E$  and  $\eta^U$  denote the output probability and the higher-order term.

Corollary 4 suggests that given any token s, the distributions of its next token and previous token significantly impact its embedding. We trained a group of Qwen2.5 models on different subsets of the Pile. Figure 5 A shows these similarity matrices for the dataset Pile-dm-mathematics, where the tokens displayed are those that occur most frequently in the corpus. We define the following correlation coefficient  $R_{\cos}\left(\boldsymbol{W}^{E},\phi^{\mathrm{next}}\right):=\mathrm{Corr}\left(\cos\left(\boldsymbol{W}^{E}\right),\cos\left(\phi^{\mathrm{next}}\right)\right)$ , and similarly  $R_{\cos}\left(\boldsymbol{W}^{U},\varphi^{\mathrm{pre}}\right)$ . Figure 5 B depicts both metrics across all subsets, suggesting that probability signatures significantly impact the structure of the embedding space and reflect the relationships among embeddings. Furthermore, we find that the probability signatures reflect the strong connections of embeddings more faithfully. As shown in Figure 5 C, the correlation between  $\mathrm{Corr}\left(\cos\left(\boldsymbol{W}_{s}^{E},\boldsymbol{W}^{E}\right),\cos\left(\phi_{s}^{\mathrm{next}},\phi^{\mathrm{next}}\right)\right)$  and  $\cos\left(\boldsymbol{W}_{s}^{E},\boldsymbol{W}^{E}\right)$  is plotted against for all tokens s, demonstrating stronger consistency in high-similarity regions. We define  $p_{\cos(\boldsymbol{W}^{E})}$  and  $p_{\cos(\boldsymbol{\psi}^{\mathrm{next}})}$  as the percentile matrix of each elements in  $\cos\left(\boldsymbol{W}^{E}\right)$  and  $\cos\left(\phi^{\mathrm{next}}\right)$ , respectively. Figure 5 D displays the distribution of  $p_{\cos(\boldsymbol{\psi}^{\mathrm{next}})}$ , conditioned on different intervals of the  $p_{\cos(\boldsymbol{W}^{E})}$ , and Figure 5 E shows the average value of  $p_{\cos(\boldsymbol{\psi}^{\mathrm{next}})}$  within each interval of  $p_{\cos(\boldsymbol{W}^{E})}$ . It can be observed that the alignment is significantly stronger in the regions with large embedding similarity. In Appendix C, we provide a detailed method explanation, a specific case of the token group with large similarity, and an analysis with the Llama-2 architecture to validate the generalization of our analysis.

Since general-purpose pretrained base models are trained on broad corpora, we attempt to directly estimate their embedding structure with a subset of general text. We combine all datasets employed in Figure 5 and define  $\tilde{\phi} = \phi^{\text{next}} + \varphi^{\text{pre}}$  (Since the tied embedding, the detail is provided in Appendix C.1). We compare the  $\cos\left(\tilde{\phi}\right)$  with  $\cos\left(W^E\right)$  of Qwen2.5-3B-base. As shown in Figure 6 A, the structure of  $\tilde{\phi}$  could capture the main properties of the pre-trained model's embedding structure, particularly the presence of sub-blocks with high similarity. Furthermore, we examine the instance for the digits ranging from 0 to 9. Figure 6 B illustrates the  $\cos\left(W^E\right)$  and  $\cos\left(\tilde{\phi}\right)$  of such digits, both revealing an ordered organization that aligns with their numerical sequence. It should be noted that this estimation may not generalize across all open-source base models, as it is sensitive to both the initialization of the pre-trained model and the true training dataset.

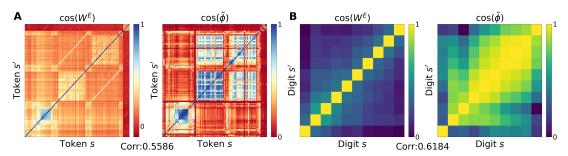


Figure 6: Cosine similarity of  $W^E$  of the Qwen2.5-3B-base and  $\tilde{\phi}$ , respectively, with the frequently-appearing tokens (A) and the digits from 0 to 9 (B).

## 7 Conclusion

In this work, we investigate the formation of embedding structures in language models. By interpreting the relationship between embedding organization and semantic structure through the lens of data distribution, we propose the probability signatures and design the addition tasks to conduct variable-controlled experiments. Our findings demonstrate that probability signatures play a crucial role in shaping the embedding structure and reflecting underlying semantic relationships. An extended analysis of LLMs further confirms our analysis. This study establishes a bridge between data semantics and embedding space, offering new insights into the understanding of the joint impact of model, data and optimization method. For future work, we plan to extend our theoretical analysis into a comprehensive framework. Besides, we aim to incorporate the self-attention mechanism into our analysis of the LLMs, which is essential for capturing more subtle and complex relationships among embeddings that remain beyond the reach of our current methods.

## References

Emmanuel Ameisen, Jack Lindsey, Adam Pearce, Wes Gurnee, Nicholas L. Turner, Brian Chen, Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar, Adly Templeton, Trenton Bricken, Cal-

- lum McDougall, Hoagy Cunningham, Thomas Henighan, Adam Jermyn, Andy Jones, Andrew Persic, Zhenyi Qi, T. Ben Thompson, Sam Zimmerman, Kelley Rivoire, Thomas Conerly, Chris Olah, and Joshua Batson. Circuit tracing: Revealing computational graphs in language models. *Transformer Circuits Thread*, 2025. URL https://transformer-circuits.pub/2025/attribution-graphs/methods.html.
- Srinadh Bhojanapalli, Chulhee Yun, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Low-rank bottleneck in multi-head attention models. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 864–873. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/bhojanapalli20a.html.
- Stella Biderman, Kieran Bicheno, and Leo Gao. Datasheet for the pile. arXiv preprint arXiv:2201.07311, 2022.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. https://transformer-circuits.pub/2023/monosemantic-features/index.html.
- Xingyu Cai, Jiaji Huang, Yuchen Bian, and Kenneth Church. Isotropy in the contextual embedding space: Clusters and manifolds. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=xYGN0860WDH.
- Jeremy M Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. *arXiv preprint arXiv:2103.00065*, 2021.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, and Inderjit Dhillon et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons in pretrained transformers. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8493–8502, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.581. URL https://aclanthology.org/2022.acl-long.581/.
- Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. Analyzing transformers in embedding space. In *Proceedings* of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2023.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, and Runxin Xu et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. URL https://arxiv.org/abs/1810.04805.
- Jacob Dunefsky, Philippe Chlenski, and Neel Nanda. Transcoders find interpretable llm feature circuits. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural Information Processing Systems, volume 37, pp. 24375–24410. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper\_files/paper/2024/file/2b8f4db0464cc5b6e9d5e6bea4b9f308-Paper-Conference.pdf.
- 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. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. https://transformer-circuits.pub/2021/framework/index.html.
- Kawin Ethayarajh. How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 55–65, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1006. URL https://aclanthology.org/D19-1006/.
- Kawin Ethayarajh, David Duvenaud, and Graeme Hirst. Towards understanding linear word analogies. In *Proceedings* of the 57th Annual Meeting of the Association for Computational Linguistics, 2019.
- Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tieyan Liu. Representation degeneration problem in training natural language generation models. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=SkEYojRqtm.

- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv* preprint arXiv:2101.00027, 2020.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.446. URL https://aclanthology.org/2021.emnlp-main.446/.
- Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 30–45, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.3. URL https://aclanthology.org/2022.emnlp-main.3/.
- Michael Hanna, Ollie Liu, and Alexandre Variengien. How does GPT-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=p4PckNQR8k.
- Michael Hanna, Sandro Pezzelle, and Yonatan Belinkov. Have faith in faithfulness: Going beyond circuit overlap when finding model mechanisms. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=TZOCCGDcuT.
- Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=F76bwRSLeK.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. arXiv preprint arXiv:2001.08361, 2020.
- Shahar Katz, Yonatan Belinkov, Mor Geva, and Lior Wolf. Backward lens: Projecting language model gradients into the vocabulary space. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 2390–2422, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.142. URL https://aclanthology.org/2024.emnlp-main.142/.
- Olga Kovaleva, Saurabh Kulshreshtha, Anna Rogers, and Anna Rumshisky. BERT busters: Outlier dimensions that disrupt transformers. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 3392–3405, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.300. URL https://aclanthology.org/2021.findings-acl.300/.
- Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams. Towards understanding grokking: An effective theory of representation learning. In *Advances in Neural Information Processing Systems*, 2022.
- Tao Luo, Zhi-Qin John Xu, Zheng Ma, and Yaoyu Zhang. Phase diagram for two-layer relu neural networks at infinite-width limit. *Journal of Machine Learning Research*, 22(71):1–47, 2021.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013a. URL https://arxiv.org/abs/1301.3781.
- Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2013b.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022. https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, and Florencia Leoni Aleman et al. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans (eds.), *Proceedings of the 2014 Conference on Empirical*

- Methods in Natural Language Processing (EMNLP), pp. 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1162. URL https://aclanthology.org/D14-1162/.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 2227–2237, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1202. URL https://aclanthology.org/N18-1202/.
- Daniel L Ruderman. The statistics of natural images. *Network: computation in neural systems*, 5(4):517, 1994. Richard Sutton. The bitter lesson. *Incomplete Ideas (blog)*, 13(1):38, 2019.
- Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.github.io/blog/gwen2.5/.
- William Timkey and Marten van Schijndel. All bark and no bite: Rogue dimensions in transformer language models obscure representational quality. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 4527–4546, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.372. URL https://aclanthology.org/2021.emnlp-main.372/.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. URL https://arxiv.org/abs/2307.09288.
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=NpsVSN6o4ul.
- Xu Wang, Yan Hu, Wenyu Du, Reynold Cheng, Benyou Wang, and Difan Zou. Towards understanding fine-tuning mechanisms of LLMs via circuit analysis. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=45EIiFd60a.
- David H Wolpert and William G Macready. No free lunch theorems for optimization. *IEEE transactions on evolutionary computation*, 1(1):67–82, 1997.
- Lei Wu, Chao Ma, et al. How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective. *Advances in Neural Information Processing Systems*, 31, 2018.
- Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle: Fourier analysis sheds light on deep neural networks. *Communications in Computational Physics*, 28(5):1746–1767, 2020.
- Zhi-Qin John Xu, Yaoyu Zhang, and Tao Luo. Overview frequency principle/spectral bias in deep learning. *Communications on Applied Mathematics and Computation*, 7(3):827–864, 2025a.
- Zhi-Qin John Xu, Yaoyu Zhang, and Zhangchen Zhou. An overview of condensation phenomenon in deep learning. *arXiv preprint arXiv:2504.09484*, 2025b.
- Junjie Yao, Zhongwang Zhang, and Zhi-Qin John Xu. An analysis for reasoning bias of language models with small initialization. In *Forty-second International Conference on Machine Learning*, 2025.
- Mengxia Yu, De Wang, Qi Shan, Colorado Reed, and Alvin Wan. The super weight in large language models, 2025. URL https://arxiv.org/abs/2411.07191.
- Zhongwang Zhang, Pengxiao Lin, Zhiwei Wang, Yaoyu Zhang, and Zhi-Qin John Xu. Initialization is critical to whether transformers fit composite functions by reasoning or memorizing. In *Advances in Neural Information Processing Systems*, 2024.

# A Experimental Setups

**Addition tasks** For each type of addition task, we trained a linear model  $F_{\rm lin}$  and a Feedforward network  $F_{\rm ffn}$ . The hidden size d=200, and we employed the ReLU as the activation function. Each dataset contains 50000 data pairs. The training is conducted for 1000 epochs with a batch size of 100. The AdamW optimizer is employed with an initial learning rate of  $10^{-5}$ .

**Language Models** In the analysis of the LLMs, we employ the Qwen2.5 architecture with 12 layers and 12 attention heads in each layer. We set up that the hidden size is 512, and the intermediate size in FFN is 1024. The dimension of the key vectors and value vectors in each head is 64. Similarly, we initialize the parameter by  $W_{i,j} \sim \mathcal{N}\left(0, d_{\text{in}}^{-1}\right)$  where  $d_{\text{in}}$  means the input dimension of W. We select five subsets of Pile, including Pile-arxiv, Pile-dm-mathematics, Pile-cc, Pile-pubmed-central, and Pile-wikipedia-en. The length of each sequence is 2048. The training is conducted for 1 epoch in each experiment, with the AdamW optimizer and a cosine learning rate schedule utilized. The initial learning rate is  $10^{-4}$ .

#### **B** Addition Task

#### **B.1** Probability Signatures in Addition Tasks

We provide a formulation of the following probability in the three addition tasks. We denote  $U(\mathcal{A})$  and  $U(\mathcal{Z})$  as the discrete uniform distribution over  $\mathcal{A}$  and  $\mathcal{Z}$ , respectively. A and Z are the random variables following  $U(\mathcal{A})$  and  $U(\mathcal{Z})$ . For the task  $f_{\mathrm{add}}$ , we have that

$$\mathbb{P}_{\pi} (y = \nu \mid \alpha \in \mathbf{X}) = \mathbb{P}_{\pi} (A + Z = \nu - \alpha), \quad \mathbb{P}_{\pi} (z \in \mathcal{X} \mid \alpha \in \mathbf{X}) = \frac{1}{|\mathcal{Z}|},$$

$$\mathbb{P}_{\pi} (z \in \mathbf{X} \mid \alpha \in \mathbf{X}, y = \nu) = \mathbb{P}_{\pi} (A = \nu - \alpha - z) = \frac{1}{|\mathcal{A}|} \delta_{\nu - \alpha - z \in \mathcal{A}},$$

$$\mathbb{P}_{\pi} (\alpha' \in \mathbf{X} \mid \alpha \in \mathbf{X}, y = \nu) = \mathbb{P}_{\pi} (Z = \nu - \alpha - \alpha') = \frac{1}{|\mathcal{Z}|} \delta_{\nu - \alpha - \alpha' \in \mathcal{Z}},$$

$$\mathbb{P}_{\pi} (z \in \mathbf{X} \mid y = \nu) = \mathbb{P}_{\pi} (A + A = \nu - z), \quad \mathbb{P}_{\pi} (\alpha \in \mathbf{X} \mid y = \nu) = \mathbb{P}_{\pi} (A + Z = \nu - \alpha),$$

where  $\alpha, \alpha' \in \mathcal{A}, z \in \mathcal{Z}$ . It's noted that besides the co-occurrence probability  $\mathbb{P}_{\pi}$   $(z \in \mathcal{X} \mid \alpha \in \mathbf{X})$ , the value of other ones is dependent on  $\alpha$  or  $\nu$ . Figure 7 (left) displays the distribution of these probabilities, which intuitively reveals the cause of the hierarchy structure in the similarity matrix. Similarly, for  $\tilde{f}_{\text{add}}$ , denote  $Y \sim U(\mathcal{Y})$  and we have

$$\mathbb{P}_{\pi} (y = \nu \mid \alpha \in \mathbf{X}) = \frac{1}{|\mathcal{Y}|}, \quad \mathbb{P}_{\pi} (z \in \mathcal{X} \mid \alpha \in \mathbf{X}) = \mathbb{P}_{\pi} (Y - A = z + \alpha),$$

$$\mathbb{P}_{\pi} (z \in \mathbf{X} \mid \alpha \in \mathbf{X}, y = \nu) = \mathbb{P}_{\pi} (A = \nu - \alpha - z) = \frac{1}{|\mathcal{A}|} \delta_{\nu - \alpha - z \in \mathcal{A}},$$

$$\mathbb{P}_{\pi} (\alpha' \in \mathbf{X} \mid \alpha \in \mathbf{X}, y = \nu) = \frac{1}{|\mathcal{Z}|},$$

$$\mathbb{P}_{\pi} (z \in \mathbf{X} \mid y = \nu) = \mathbb{P}_{\pi} (A + A = \nu - z), \quad \mathbb{P}_{\pi} (\alpha \in \mathbf{X} \mid y = \nu) = \mathbb{P}_{\pi} (A + Z = \nu - \alpha).$$

For  $f_{\text{mod}}$ , we have

$$\begin{split} & \mathbb{P}_{\pi}\left(y=\nu \mid \alpha \in \boldsymbol{X}\right) = \frac{1}{|\mathcal{Z}|}, \quad \mathbb{P}_{\pi}\left(z \in \mathcal{X} \mid \alpha \in \boldsymbol{X}\right) = \frac{1}{|\mathcal{Z}|}, \\ & \mathbb{P}_{\pi}\left(z \in \boldsymbol{X} \mid \alpha \in \boldsymbol{X}, y=\nu\right) = \frac{1}{|\mathcal{A}|} \delta_{\nu-\min \mathcal{Z}-(\alpha-z \bmod{|\mathcal{Z}|}) \in (A \bmod{|\mathcal{Z}|})}, \\ & \mathbb{P}_{\pi}\left(\alpha' \in \boldsymbol{X} \mid \alpha \in \boldsymbol{X}, y=\nu\right) = \frac{1}{|\mathcal{Z}|}, \\ & \mathbb{P}_{\pi}\left(z \in \boldsymbol{X} \mid y=\nu\right) = \mathbb{P}_{\pi}\left((A+A \bmod{|\mathcal{Z}|}) = \nu - \min \mathcal{Z} - (z \bmod{|\mathcal{Z}|})\right), \\ & \mathbb{P}_{\pi}\left(\alpha \in \boldsymbol{X} \mid y=\nu\right) = \mathbb{P}_{\pi}\left((A+Z \bmod{|\mathcal{Z}|}) = \nu - \min \mathcal{Z} - (\alpha \bmod{|\mathcal{Z}|})\right). \end{split}$$

Figure 7 depicts all these probability distributions.

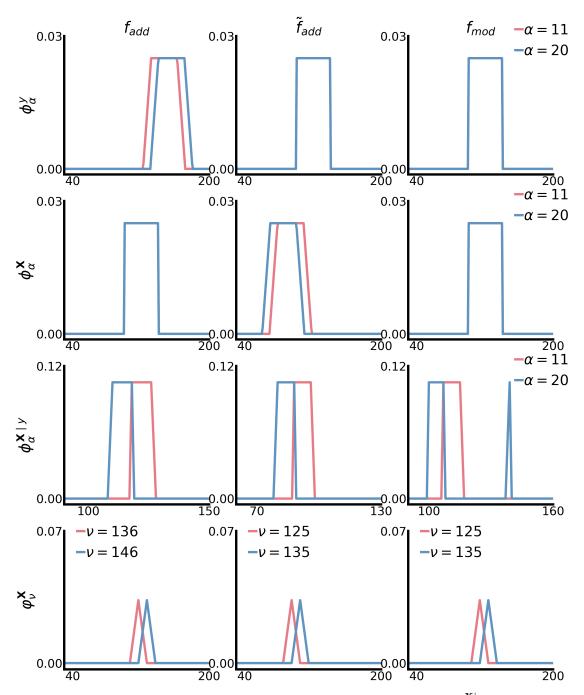


Figure 7: Probability signatures in each task under distinct  $\alpha$  and  $\nu$ . In the distribution of  $\phi_{\alpha}^{X|y}$ , y=150 is displayed in  $f_{\rm add}$  and y=120 in  $\tilde{f}_{\rm add}$  and  $f_{\rm mod}$ , since 150 and 120 are the average label value in each task.

## **B.2** Embedding matrix in Linear Model

Figure 8 depicts the PCA projection of the anchor embeddings in  $F_{\rm lin}$ , revealing that  $f_{\rm add}$  and  $\tilde{f}_{\rm add}$  both establish an ordered structure while the anchor embeddings in  $f_{\rm mod}$  are chaotic.

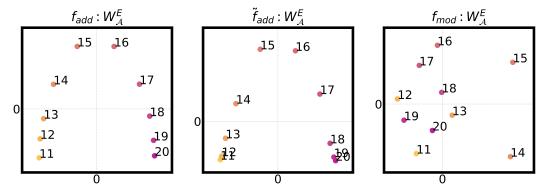


Figure 8: PCA projection of  $W_A^E$  in  $F_{\text{lin}}$  (epoch 120).

# B.3 Umembedding matrix in Feedforward Network

Figure 9 displays the structure of the unembedding matrix in  $F_{\rm ffn}$  with the three types of addition tasks. The distribution of  $\cos\left(\mathbf{W}_{\nu}^{U}\right)$  (A) and the PCA projection (B) jointly reveal that the unembedding vectors of those label tokens establish a hierarchy structure, which is consistent with their natural sequence.

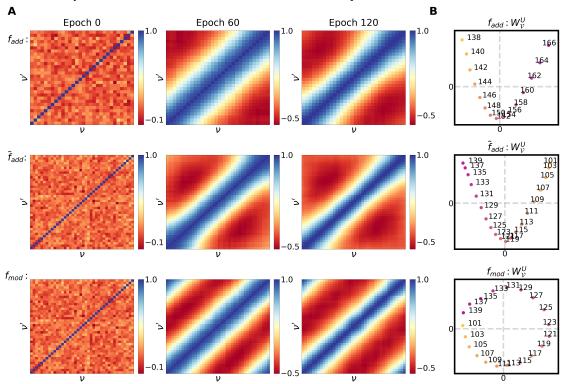


Figure 9: A: The heatmap of the  $\cos (\mathbf{W}_{\mathcal{V}}^U)$  with label index in  $F_{\text{ffn}}$  during the training process. B: PCA projection of  $\mathbf{W}_{\mathcal{V}}^U$  in  $F_{\text{ffn}}$  (epoch 120).

# C Language Models

#### C.1 Technical Details

**Remark about Figure 5 C** In each subset  $D_i$ ,  $i = 1, 2, \dots M$ , we define the set  $S_i = \left\{s_j^i\right\}_{j=1}^{C_i}$  as the set of the  $C_i$  tokens which appear most frequently in  $D_i$ . Based on the dataset  $D_i$ , and denote  $W^{E_i}$  as the embedding matrix of the model corresponding to dataset  $D_i$ , we compute that

$$\cos_{D_i}\left(\boldsymbol{W}_{s_j^i}^{E}, \boldsymbol{W}^{E}\right) = \left[\cos\left(\boldsymbol{W}_{s_j^i}^{E_i}, \boldsymbol{W}_{s'}^{E_i}\right)\right]_{s' \in \mathcal{S}_i} \in \mathbb{R}^{C_i},$$

and

$$\cos_{D_i}\left(oldsymbol{\phi}_{s^i_j}^{ ext{next}}, oldsymbol{\phi}^{ ext{next}}
ight) = \left[\cos\left(oldsymbol{\phi}_{s^i_j}^{ ext{next}}, oldsymbol{\phi}_{s'}^{ ext{next}}
ight)
ight]_{s' \in \mathcal{S}_i} \in \mathbb{R}^{C_i}.$$

for any token  $s_i^i \in \mathcal{S}_i$ . Then we define the correlation coefficient

$$R_{D_{i}}\left(s_{j}^{i}\right) = \operatorname{Corr}\left(\cos_{D_{i}}\left(\boldsymbol{W}_{s_{j}^{i}}^{E}, \boldsymbol{W}^{E}\right), \cos_{D_{i}}\left(\boldsymbol{\phi}_{s_{j}^{i}}^{\operatorname{next}}, \boldsymbol{\phi}^{\operatorname{next}}\right)\right)$$

and the average embedding similarity as

$$\operatorname{Mean}_{\boldsymbol{W}^{E},D_{i}}\left(s_{j}^{i}\right) = \frac{1}{C_{i}} \cos_{D_{i}}\left(\boldsymbol{W}_{s_{j}^{i}}^{E}, \boldsymbol{W}^{E}\right) \cdot \mathbf{1}.$$

Then we concatenate the metrics with all token  $s_j^i \in \mathcal{S}_i, j=1,2,\cdots,C_i$  and all datasets  $\mathcal{S}_i, i=1,2,\cdots,M$ , i.e.

$$\operatorname{Corr}\left(\cos\left(\boldsymbol{W}_{s}^{E},\boldsymbol{W}^{E}\right),\cos\left(\boldsymbol{\phi}_{s}^{\operatorname{next}},\boldsymbol{\phi}^{\operatorname{next}}\right)\right)=\left[R_{D_{i}}\left(s_{j}^{i}\right)\right]_{j=1,2,\cdots,C_{i}}^{i=1,2,\cdots,M}\in\mathbb{R}^{\sum_{i=1}^{M}C_{i}},$$

$$\operatorname{Mean}\left(\cos\left(\boldsymbol{W}_{s}^{E},\boldsymbol{W}^{E}\right)\right)=\left[\operatorname{Mean}_{\boldsymbol{W}^{E},D_{i}}\left(s_{j}^{i}\right)\right]_{j=1,2,\cdots,C_{i}}^{i=1,2,\cdots,M}\in\mathbb{R}^{\sum_{i=1}^{M}C_{i}}.$$

Figure 5 displays the relation between  $\operatorname{Corr}\left(\cos\left(\boldsymbol{W}_{s}^{E},\boldsymbol{W}^{E}\right),\cos\left(\boldsymbol{\phi}_{s}^{\operatorname{next}},\boldsymbol{\phi}^{\operatorname{next}}\right)\right)$  and  $\operatorname{Mean}\left(\cos\left(\boldsymbol{W}_{s}^{E},\boldsymbol{W}^{E}\right)\right)$ , revealing a positive correlation. In our work, M=5, and we set up  $C_{i}=10000$  for each dataset.

**Remark about Figure 5 D & E** In each subset  $D_i$ ,  $i = 1, 2, \dots M$ , we define the set  $S_i = \left\{s_j^i\right\}_{j=1}^{C_i}$  as the set of the  $C_i$  tokens which appear most frequently in  $D_i$ . We compute that

$$\cos_{D_i}\left(\mathbf{W}^E\right) = \left[\cos\left(\mathbf{W}_s^{E_i}, \mathbf{W}_{s'}^{E_i}\right)\right]_{s,s' \in \mathcal{S}_i} \in \mathbb{R}^{C_i \times C_i}$$

and

$$\cos_{D_i}\left(\boldsymbol{\phi}^{ ext{next}}\right) = \left[\cos\left(\boldsymbol{\phi}^{ ext{next}}_s, \boldsymbol{\phi}^{ ext{next}}_{s'}\right)\right]_{s,s' \in S_i} \in \mathbb{R}^{C_i \times C_i}.$$

Then translate the similarity matrix into a percentile formulation, i.e.

$$p_{\cos_{D_i}(\mathbf{W}^E)} = \text{Percentile}\left(\cos_{D_i}\left(\mathbf{W}^E\right)\right), \quad p_{\cos_{D_i}(\boldsymbol{\phi}^{\text{next}})} = \text{Percentile}\left(\cos_{D_i}\left(\boldsymbol{\phi}^{\text{next}}\right)\right)$$

 $\text{and } p_{\cos(\boldsymbol{W}^E)} = \left[p_{\cos_{D_i}(\boldsymbol{W}^E)}\right]_{i=1,2,\cdots,M}, \quad p_{\cos(\boldsymbol{\phi}^{\text{next}})} = \left[p_{\cos_{D_i}(\boldsymbol{\phi}^{\text{next}})}\right]_{i=1,2,\cdots,M}. \text{ Figure 5 D and E reveal the distribution and average value of } p_{\cos(\boldsymbol{\phi}^{\text{next}})}, \text{ where } k\times 10\% \leq p_{\cos(\boldsymbol{W}^E)} < (k+1)\times 10\%, k=0,1,2,\cdots,9.$ 

**Tied Embedding** In the Qwen2.5-3B-base model, the embedding matrix and unembedding matrix are the same one, which aims for computational source saving. Under this condition, we have that

$$\begin{split} \frac{d\boldsymbol{W}_{s}^{E}}{dt} &= r_{s}^{\mathrm{in}} \boldsymbol{W}^{U,T} \boldsymbol{\phi}_{s}^{\mathrm{next}} + r_{s}^{\mathrm{out}} \boldsymbol{W}^{E} \boldsymbol{\varphi}_{s}^{\mathrm{pre}} + \boldsymbol{\eta} \\ &= \boldsymbol{W}^{E} \left( r_{s}^{\mathrm{in}} \boldsymbol{\phi}_{s}^{\mathrm{next}} + r_{s}^{\mathrm{out}} \boldsymbol{\varphi}_{s}^{\mathrm{pre}} \right) + \boldsymbol{\eta}. \end{split}$$

Since the next-token-prediction, each token will be an input and an output, except the last token in a sequence, resulting in  $r_s^{\rm in} \approx r_s^{\rm out}$ . Denote  $r_s = r_s^{\rm in}$  and  $\tilde{\phi}_s = \phi_s^{\rm next} + \varphi_s^{\rm pre}$ , then we have

$$\frac{d\mathbf{W}_s^E}{dt} = r_s \mathbf{W}^E \tilde{\boldsymbol{\phi}}_s + \boldsymbol{\eta}.$$

## **C.2** Complete results

Figure 10 represents the cosine similarity distribution of  $W^E$ ,  $\phi^{\text{next}}$ ,  $W^U$  and  $\varphi^{\text{pre}}$  in the other 4 subsets of Pile we selected, exhibiting an analogous phenomenon with the observation in Figure 5. The distribution representations  $\phi^{\text{next}}$  and  $\varphi^{\text{pre}}$  could effectively capture the high similarity among embedding vectors and unembedding vectors. Figure 11 displays the completed result of Figure 5 D.

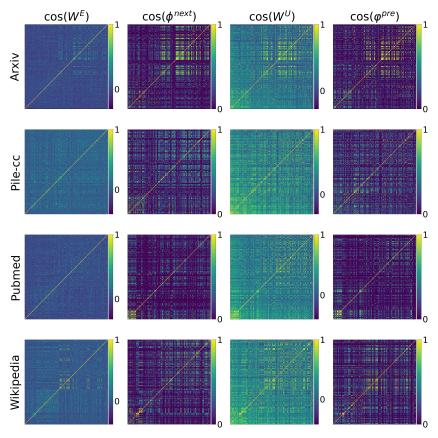


Figure 10: Cosine similarity distribution of  $\mathbf{W}^E, \phi^{\text{next}}, \mathbf{W}^U, \varphi^{\text{pre}}$  in each experiment with distinct dataset. The tokens displayed are those with the most appearances in the dataset.

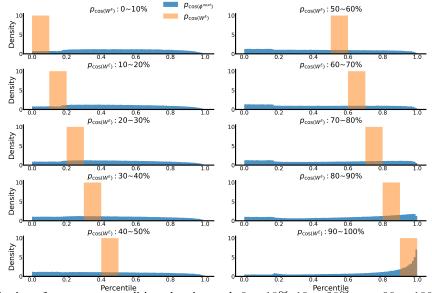
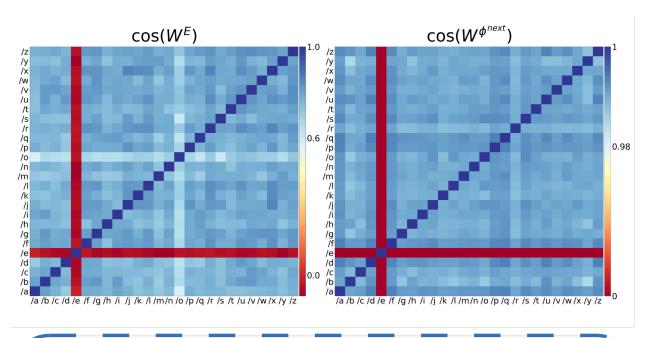


Figure 11: Distribution of  $p_{\cos(\phi^{\text{next}})}$ , conditioned on intervals  $0 \sim 10\%$ ,  $10 \sim 20\%$ ,  $\cdots$ ,  $90 \sim 100\%$  of the  $p_{\cos(\boldsymbol{W}^E)}$ .

## C.3 Case Analysis

We provide a detailed case to explain the group of tokens exhibiting high embedding similarities. In experiments on the Pile-dm-mathematics dataset, tokens such as "/a", "/b", "/c", and "/d" often serve as denominators in mathematical expressions. Figure 12 shows the cosine similarities of both their embedding vectors and distribution representations, which are notably high for all tokens except "/e", which does not appear in the dataset. These tokens share highly similar semantics and also exhibit very similar next-token distributions, most frequently followed by "\*" or ")". This similarity in next-token distribution leads to strong similarities in their embedding vectors. This example vividly illustrates how data distribution shapes semantic structure within the embedding space, particularly in the case of tokens with high semantic affinity.



```
Simplify ((k**(-9)*k*k/((k**0/k*k)/k)*k)/(k**(1/2))**(-26))/(((k*k**(2/3)*k)/k**0)/(k**(-1)*k)**14) assuming k is positive. k**(16/3)

Simplify (((a**31/a*a)/a)/a)/(a/(a/(a/(a/a**(6/5)))/a)))*(a*a**(-1/67))/((a*(a*a/(a*(a**14*a)/a))/a)/a) assuming a is positive. a**
(14668/335)

Simplify ((j*(j**(-10/11)/j)/j)/j**(3/4)*(j/j**4)**7)**(-22/7) assuming j is positive.

Simplify (u*u**(-19))/u**(-1/3)*u/(u**(-5)*u)*u**0 assuming u is positive. u**(-38/3)

Simplify h*h/(h/(h**14/h)*h*h)*h/(h/(h/h**12))*h**(-11)/(((h**(-4)*h)/h)/h) assuming h is positive. h**(-5)

Simplify ((r/r**(-2/3))**4*r*r*r*r**(-2/9)*r*r/((r**0*r)/r))**33 assuming r is positive
```

Figure 12: A case analysis of the token group "/a", "/b", "/c", etc. The first row depicts the cosine similarity of their embeddings (left) and distribution representations (right). The second row exhibits the contexts containing these tokens, which are highlighted by different colors.

#### C.4 Results of LLama 2

To assess the generalizability of our analysis in Section 6 across different model architectures and tokenizers, we replicate the experiment using the Llama 2 architecture. We employ the same dataset from Pile, and the training configurations are the same as the experiments of Qwen2.5. As shown in Figure 13, the probability signatures effectively capture structural relationships in the embedding space, especially in regions exhibiting high embedding similarity. These results align closely with those in Figure 5, indicating that our analytical approach is robust to variations in model architecture.

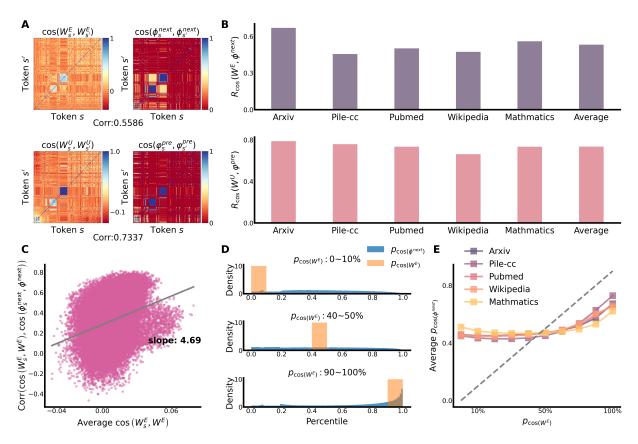


Figure 13: Results with Llama-2 architecture. A: Heatmap of the cosine similarity of  $\boldsymbol{W}^E, \boldsymbol{W}^U, \boldsymbol{\phi}^{\text{next}}$  and  $\boldsymbol{\varphi}^{\text{pre}}$ . B:  $R_{\cos}\left(\boldsymbol{W}^E, \boldsymbol{\phi}^{\text{next}}\right)$  (top) and  $R_{\cos}\left(\boldsymbol{W}^U, \boldsymbol{\varphi}^{\text{pre}}\right)$  (bottom) with different datasets. C: Relation between Corr  $\left(\cos\left(\boldsymbol{W}_s^E, \boldsymbol{W}^E\right), \cos\left(\boldsymbol{\phi}_s^{\text{next}}, \boldsymbol{\phi}^{\text{next}}\right)\right)$  and the average value of  $\cos\left(\boldsymbol{W}_s^E, \boldsymbol{W}^E\right)$ . Each point denotes a token s. D: Distribution of  $p_{\cos(\boldsymbol{\phi}^{\text{next}})}$ , conditioned on intervals  $0 \sim 10\%, 40 \sim 50\%$  and  $90 \sim 100\%$  of the  $p_{\cos(\boldsymbol{W}^E)}$ . E: Average value of  $p_{\cos(\boldsymbol{\phi}^{\text{next}})}$  within each interval of  $p_{\cos(\boldsymbol{W}^E)}$ .

#### **D** Theoretical Details

## D.1 Proof of Proposition 1

**Lemma 1.** Given a model F and data pair  $(\mathbf{X},y) \in \mathbb{N}^{+,L} \times \mathbb{N}^{+}$ ,  $\ell = -\log Softmax (F(\mathbf{X}))_{y}$ , we have that

$$\frac{\partial \ell}{\partial F(\mathbf{X})} = \mathbf{p} - \mathbf{e}_y,\tag{5}$$

where p = softmax(X).

*Proof.* It's noted  $\ell = -F(\boldsymbol{X})_y + \log \sum_{j=1}^{d_{\text{vob}}} \exp F(\boldsymbol{X})_j$ , then we have

$$\frac{\partial \ell}{\partial F\left(X\right)_{i}} = -\delta_{i=y} + \frac{\exp F\left(\boldsymbol{X}\right)_{i}}{\sum_{j=1}^{d_{\text{vob}}} \exp F\left(\boldsymbol{X}\right)_{j}} = \boldsymbol{p}_{i} - \delta_{i=y},$$

where  $\delta_{i=y}=1$  if i=y else 0. This indicates that  $\frac{\partial \ell}{\partial F(\mathbf{X})}=\mathbf{p}-\mathbf{e}_y$ .

With Lemma 1, we could obtain the derivative of  $\ell$  with respect to  $\boldsymbol{W}_{x}^{E}$  for any  $x \in \mathcal{V}$  as follows:

$$\frac{\partial \ell^{i}}{\partial \boldsymbol{W}_{x}^{E}} = \frac{\partial F\left(\boldsymbol{X}^{i}\right)}{\partial \boldsymbol{W}_{x}^{E}} \frac{\partial \ell^{i}}{\partial F\left(\boldsymbol{X}^{i}\right)} \\
= \left(\boldsymbol{W}^{U,T}\left(\boldsymbol{p}^{i} - \boldsymbol{e}_{y^{i}}\right)\right) \odot G^{(1)}\left(\boldsymbol{W}_{\boldsymbol{X}^{i}}^{E}\right).$$

Then the gradient flow of  $vW_x^E$  could be obtained by

$$\frac{d\boldsymbol{W}_{x}^{E}}{dt} = -\frac{1}{N} \sum_{i=1}^{N} \frac{\partial \ell^{i}}{\partial \boldsymbol{W}_{x}^{E}} = \frac{1}{N} \sum_{i=1}^{N} \left( \boldsymbol{W}^{U,T} \left( \boldsymbol{p}^{i} - \boldsymbol{e}_{y^{i}} \right) \right) \odot G^{(1)} \left( \boldsymbol{W}_{\boldsymbol{X}^{i}}^{E} \right),$$

Since diag  $\left(G^{(1)}\left(\boldsymbol{W}_{\boldsymbol{X}^{i}}^{E}\right)\right)=0$  if  $x\notin\boldsymbol{X}^{i}$ , we have that

$$\begin{split} \frac{d\boldsymbol{W}_{x}^{E}}{dt} &= \frac{1}{N} \sum_{i=1}^{N_{x}^{\text{in}}} \left( \boldsymbol{W}^{U,T} \left( \boldsymbol{e}_{y_{x}^{i}} - \boldsymbol{p}_{x}^{i} \right) \right) \odot G^{(1)} \left( \boldsymbol{W}_{\boldsymbol{X}_{x}^{i}}^{E} \right) \\ &= \frac{r_{x}^{\text{in}}}{N_{x}^{\text{in}}} \sum_{i=1}^{N_{x}^{\text{in}}} \left( \boldsymbol{W}^{U,T} \left( \boldsymbol{e}_{y_{x}^{i}} - \boldsymbol{p}_{x}^{i} \right) \right) \odot G^{(1)} \left( \boldsymbol{W}_{\boldsymbol{X}_{x}^{i}}^{E} \right). \end{split}$$

Since that  $y_x^i$  takes value  $\nu \in \mathcal{V}$ , we can rewrite this formation as

$$\frac{d\boldsymbol{W}_{x}^{E}}{dt} = \sum_{\nu \in \mathcal{V}} \frac{r_{x,\nu}}{N_{x,\nu}} \left( \boldsymbol{W}^{U,T} \boldsymbol{e}_{\nu} \right) \odot \sum_{i=1}^{N_{x,\nu}} G^{(1)} \left( \boldsymbol{W}_{\boldsymbol{X}_{(x,\nu)}^{i}}^{E} \right) - \frac{r_{x}^{\text{in}}}{N_{x}^{\text{in}}} \sum_{i=1}^{N_{x}^{\text{in}}} G^{(1)} \left( \boldsymbol{W}_{\boldsymbol{X}_{x}^{i}}^{E} \right) \odot \left( \boldsymbol{W}^{U,T} \boldsymbol{p}_{x}^{i} \right).$$

# D.2 Proof of Proposition 2

Similar with the analysis of  ${m W}_x^E$ , we derive the gradient flow of  ${m W}_
u^U$  as follows:

$$\begin{split} \frac{d\boldsymbol{W}_{\nu}^{U}}{dt} &= -\frac{1}{N} \sum_{i=1}^{N} \frac{\partial \ell^{i}}{\partial \boldsymbol{W}_{\nu}^{U}} \\ &= & \frac{1}{N} \sum_{i=1}^{N} \left( \boldsymbol{e}_{\boldsymbol{y}^{i,\nu}} - \boldsymbol{p}^{i,\nu} \right) \left[ \boldsymbol{G} \left( \boldsymbol{W}_{\boldsymbol{X}^{i}}^{E} \right) \right]^{T}. \end{split}$$

Since  $e_{y^{i,\nu}}=1$  if  $y^i=\nu$  else 0, we have that

$$\frac{d\boldsymbol{W}_{\nu}^{U}}{dt} = \frac{r_{\nu}^{\text{out}}}{N_{\nu}^{\text{out}}} \sum_{i=1}^{N_{\nu}^{\text{out}}} \left[ G\left(\boldsymbol{W}_{\boldsymbol{X}_{(\cdot,\nu)}^{E}}^{E}\right) \right]^{T} - \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{p}^{i,\nu} \left[ G\left(\boldsymbol{W}_{\boldsymbol{X}^{i}}^{E}\right) \right]^{T}.$$

#### D.3 Proof of Corollary 1

With proposition 1, we have that

$$\frac{d\boldsymbol{W}_{\alpha}^{E}}{dt} = \boldsymbol{W}^{U,T} \left( \sum_{\nu \in \mathcal{V}} r_{\alpha,\nu} \boldsymbol{e}_{\nu} - \frac{r_{\alpha}^{\text{in}}}{N_{\alpha}^{\text{in}}} \sum_{i=1}^{N_{x}^{\text{in}}} \boldsymbol{p}_{\alpha}^{i} \right) \\
= \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \left( \sum_{\nu \in \mathcal{V}} \frac{r_{\alpha,\nu}}{r_{\alpha}^{\text{in}}} \boldsymbol{e}_{\nu} - \frac{1}{N_{\alpha}^{\text{in}}} \sum_{i=1}^{N_{x}^{\text{in}}} \boldsymbol{p}_{\alpha}^{i} \right).$$

Utilizing that softmax  $(f) = \frac{1}{d_{\text{vob}}} \mathbf{1} + \frac{1}{d_{\text{vob}}} f + \mathcal{O}\left(d_{\text{vob}}^{-2} f\right)$ , we obtain that

$$\frac{d\boldsymbol{W}_{\alpha}^{E}}{dt} = \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \left( \sum_{\nu \in \mathcal{V}} \frac{r_{\alpha,\nu}}{r_{\alpha}^{\text{in}}} \boldsymbol{e}_{\nu} - \frac{1}{N_{\alpha}^{\text{in}}} \sum_{i=1}^{N_{\alpha}^{\text{in}}} \left( \frac{1}{d_{\text{vob}}} \mathbf{1} - \frac{1}{d_{\text{vob}}} \boldsymbol{W}^{U} \left( \boldsymbol{W}_{z_{i}}^{E} + \boldsymbol{W}_{\alpha_{i}}^{E} + \boldsymbol{W}_{\alpha}^{E} \right) + \mathcal{O} \left( d_{\text{vob}}^{-2} \boldsymbol{W}^{U} \boldsymbol{W}_{\alpha}^{E} \right) \right) \right) \\
= \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \left( \sum_{\nu \in \mathcal{V}} \frac{r_{\alpha,\nu}}{r_{\alpha}^{\text{in}}} \boldsymbol{e}_{\nu} - \frac{1}{d_{\text{vob}}} \mathbf{1} + \frac{1}{d_{\text{vob}}} \boldsymbol{W}^{U} \left( \frac{1}{N_{\alpha}^{\text{in}}} \sum_{i=1}^{N_{\alpha}^{\text{in}}} \left( \boldsymbol{W}_{z_{i}}^{E} + \boldsymbol{W}_{\alpha_{i}}^{E} \right) + \boldsymbol{W}_{\alpha}^{E} \right) + \mathcal{O} \left( d_{\text{vob}}^{-2} \boldsymbol{W}^{U} \boldsymbol{W}_{\alpha}^{E} \right) \right) \\
= \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \left( \sum_{\nu \in \mathcal{V}} \frac{r_{\alpha,\nu}}{r_{\alpha}^{\text{in}}} \boldsymbol{e}_{\nu} - \frac{1}{d_{\text{vob}}} \mathbf{1} + \frac{1}{d_{\text{vob}}} \boldsymbol{W}^{U} \left( \sum_{x \in (\mathcal{Z} \cap \mathcal{A})} \frac{N_{\alpha,x}^{\text{in}}}{N_{\alpha}^{\text{in}}} \boldsymbol{W}_{x}^{E} + \boldsymbol{W}_{\alpha}^{E} \right) + \mathcal{O} \left( d_{\text{vob}}^{-2} \boldsymbol{W}^{U} \boldsymbol{W}_{\alpha}^{E} \right) \right).$$

Let  $N \to \infty$ , we have that

$$\begin{split} \frac{d\boldsymbol{W}_{\alpha}^{E}}{dt} = & \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \left( \sum_{\nu \in \mathcal{V}} \mathbb{P}_{\pi} \left( y = \nu \mid \alpha \in \boldsymbol{X} \right) \boldsymbol{e}_{\nu} - \frac{1}{d_{\text{vob}}} \boldsymbol{1} \right. \\ & \left. + \frac{1}{d_{\text{vob}}} \boldsymbol{W}^{U} \left( \sum_{x \in (\mathcal{Z} \cap \mathcal{A})} \mathbb{P}_{\pi} \left( x \in \boldsymbol{X} \mid \alpha \in \boldsymbol{X} \right) \boldsymbol{W}_{x}^{E} + \boldsymbol{W}_{\alpha}^{E} \right) + \mathcal{O} \left( d_{\text{vob}}^{-2} \boldsymbol{W}^{U} \boldsymbol{W}_{\alpha}^{E} \right) \right) \\ = & \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \left( \boldsymbol{\phi}_{\alpha}^{y} + \frac{1}{d_{\text{vob}}} \boldsymbol{W}^{U} \boldsymbol{W}^{E} \boldsymbol{\phi}_{\alpha}^{\boldsymbol{X}} \right) + \boldsymbol{\eta}, \end{split}$$

where  $\boldsymbol{\eta} = \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \left( \frac{1}{d_{\text{vob}}} \left( \boldsymbol{W}^{U} \boldsymbol{W}_{\alpha}^{E} - \mathbf{1} \right) + \mathcal{O} \left( d_{\text{vob}}^{-2} \boldsymbol{W}^{U} \boldsymbol{W}_{\alpha}^{E} \right) \right)$  contains the higher-order term and the data independent term.

#### D.4 Proof of Corollary 2

*Proof.* Since the small initialization, we assume that the activation function can be approximated by the following form with the Weierstrass approximation theorem.

$$\sigma\left(\sum_{x \in \mathbf{X}} \mathbf{W}_{x}^{E}\right) = C_{0} + C_{1}\left(\sum_{x \in \mathbf{X}} \mathbf{W}_{x}^{E}\right) + C_{2}\left(\sum_{x \in \mathbf{X}} \mathbf{W}_{x}^{E}\right)^{\odot 2} + \epsilon.$$

With the loss of the generalization, we assume that  $C_0 = 0, C_1 = 1, C_2 = \frac{1}{2}$ . Then we have

$$\frac{d\boldsymbol{W}_{\alpha}^{E}}{dt} = \underbrace{\sum_{\nu \in \mathcal{V}} \frac{r_{\alpha,\nu}}{N_{\alpha,\nu}} \left( \boldsymbol{W}^{U,T} \boldsymbol{e}_{\nu} \right) \odot \sum_{i=1}^{N_{\alpha,\nu}} \left( \mathbf{1} + \sum_{x \in \boldsymbol{X}_{(\alpha,\nu)}^{i}} \boldsymbol{W}_{\boldsymbol{X}_{(\alpha,\nu)}^{i}}^{E} \right)}_{\boldsymbol{J}^{y}} - \underbrace{\frac{r_{\alpha}^{\text{in}}}{N_{\alpha}^{\text{in}}} \sum_{i=1}^{N_{\alpha}^{\text{in}}} \left( \mathbf{1} + \sum_{x \in \boldsymbol{X}_{\alpha}^{i}} \boldsymbol{W}_{\boldsymbol{X}_{\alpha}^{i}}^{E} \right) \odot \left( \boldsymbol{W}^{U,T} \boldsymbol{p}_{\alpha}^{i} \right)}_{\boldsymbol{I}^{p}}.$$

For the term  $J^y$  we have

$$\boldsymbol{J}^{y} = \boldsymbol{W}^{U,T} \sum_{\nu \in \mathcal{V}} r_{\alpha,\nu} \boldsymbol{e}_{\nu} + \sum_{\nu \in \mathcal{V}} r_{\alpha,\nu} \left( \boldsymbol{W}^{U,T} \boldsymbol{e}_{\nu} \right) \odot \left( \boldsymbol{W}_{\alpha}^{E} + \frac{1}{2N_{\alpha,\nu}} \sum_{i=1}^{2N_{\alpha,\nu}} \boldsymbol{W}_{x_{(\alpha,\nu)}^{i}}^{E} \right).$$

Let  $N \to \infty$ , we have that

$$\begin{split} \boldsymbol{J}^{y} &= \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \boldsymbol{\phi}_{\alpha}^{y} \odot \left( \boldsymbol{1} + \boldsymbol{W}_{\alpha}^{E} \right) + \sum_{\nu \in \mathcal{V}} \operatorname{diag} \left( \boldsymbol{W}_{\nu}^{U} \right) r_{\alpha,\nu} \sum_{x' \in \mathcal{V}} \mathbb{P} \left( x' \in \boldsymbol{X} \mid \alpha \in \boldsymbol{X}, y = \nu \right) \boldsymbol{W}_{x'}^{E} \\ &= \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \boldsymbol{\phi}_{\alpha}^{y} \odot \left( \boldsymbol{1} + \boldsymbol{W}_{\alpha}^{E} \right) + \sum_{\nu \in \mathcal{V}} \operatorname{diag} \left( r_{\alpha,\nu} \boldsymbol{W}_{\nu}^{U} \right) \boldsymbol{W}^{E} \left( \boldsymbol{\phi}_{\alpha}^{\boldsymbol{X}|y,T} \right)_{\nu}^{T} \\ &= \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \boldsymbol{\phi}_{\alpha}^{y} \odot \left( \boldsymbol{1} + \boldsymbol{W}_{\alpha}^{E} \right) + \mathbb{T} \cdot \left( \boldsymbol{\phi}_{\alpha}^{\boldsymbol{X}|y} \right)^{T}, \end{split}$$

where  $\mathbb{T} \in \mathbb{R}^{d \times d_{\mathrm{vob}} \times d_{\mathrm{vob}}}$ ,  $\mathbb{T}_{:,:,\nu} = r_{\alpha,\nu} \mathrm{diag}\left( \boldsymbol{W}_{\nu}^{U} \right) \boldsymbol{W}^{E}$  for  $\nu \in \mathcal{V}$  and 0 otherwise.

Similarly, for the term  $J^p$ , we have that

$$\boldsymbol{J}^{p} = \boldsymbol{W}^{U,T} r_{\alpha}^{\text{in}} \left( \frac{1}{d_{\text{vob}}} \boldsymbol{1} - \frac{1}{d_{\text{vob}}} \boldsymbol{W}^{U} \left( \left( \boldsymbol{W}^{E} - \text{diag} \left( \boldsymbol{W}^{U,T} \boldsymbol{1} \right) \right) \boldsymbol{\phi}_{\alpha}^{\boldsymbol{X}} + \boldsymbol{W}_{\alpha}^{E} \right) + \boldsymbol{\epsilon} \right),$$

where  $\epsilon = \mathcal{O}\left(\frac{1}{d_{\mathrm{vob}}^2} m{W}^U m{W}_{\alpha}^E\right)$ . Then we have that

$$\frac{d\boldsymbol{W}_{\alpha}^{E}}{dt} = \mathbb{T} \cdot \left(\boldsymbol{\phi}_{\alpha}^{\boldsymbol{X}|y}\right)^{T} + \boldsymbol{\eta}_{\boldsymbol{\phi}_{\alpha}^{y}} + \frac{1}{d_{\text{vob}}} \boldsymbol{\eta}_{\boldsymbol{\phi}_{\alpha}^{\boldsymbol{X}}},$$

where  $m{\eta}_{m{\phi}_{lpha}^y} = m{W}^{U,T} r_{lpha}^{ ext{in}} m{\phi}_{lpha}^y \odot \left( \mathbf{1} + m{W}_{lpha}^E \right), m{\eta}_{m{\phi}_{lpha}^X} = d_{ ext{vob}} m{J}^p.$ 

## D.5 Proof of Corollary 3

*Proof.* With Proposition 2, we have that

$$\begin{split} \frac{d\boldsymbol{W}_{\nu}^{U}}{dt} &= \frac{r_{\nu}^{\text{out}}}{N_{\nu}^{\text{out}}} \sum_{i=1}^{N_{\nu}^{\text{out}}} \left( \sum_{x \in \boldsymbol{X}_{(\cdot,\nu)}^{i}} \boldsymbol{W}_{x}^{E} \right)^{T} - \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{p}^{i,\nu} \left( \sum_{x \in \boldsymbol{X}^{i}} \boldsymbol{W}_{x}^{E} \right)^{T} \\ &= Lr_{\nu}^{\text{out}} \sum_{x \in \mathcal{V}} \mathbb{P}_{\pi} \left( x \in \boldsymbol{X} \mid y = \nu \right) \boldsymbol{W}_{x}^{E,T} - L \sum_{x \in \mathcal{V}} \mathbb{E}_{\pi} \left[ \boldsymbol{p}^{\nu} \mid x \in \boldsymbol{X} \right] \boldsymbol{W}_{x}^{E,T} \\ &= Lr_{\nu}^{\text{out}} \left( \boldsymbol{W}^{E} \boldsymbol{\varphi}_{\nu}^{\boldsymbol{X}} \right)^{T} - \boldsymbol{\eta}, \end{split}$$

where  $\boldsymbol{\eta} = L\left(\boldsymbol{W}^{E}\mathbb{E}_{\pi}\left[\boldsymbol{p} \mid x \in \boldsymbol{X}\right]\right)^{T}$ .

## D.6 Proof of Corollary 4

*Proof.* The next-token-prediction training loss could be formulated as

$$\ell^{i} = \frac{1}{L} \sum_{t=1}^{L-1} \text{CrossEntropy} \left( F_{\text{lan}} \left( \boldsymbol{X}_{:t} \right) ; \boldsymbol{e}_{\boldsymbol{X}_{t+1}} \right).$$

So we have that

$$\frac{\partial \ell^{i}}{\partial \boldsymbol{W}_{s}^{E}} = \frac{1}{L} \sum_{t=1}^{L-1} \boldsymbol{W}^{U,T} \left( \boldsymbol{p}_{t}^{i} - \boldsymbol{e}_{\boldsymbol{X}_{t+1}^{i}} \right) \odot \left( \delta_{\boldsymbol{X}_{t}^{i}=s} \boldsymbol{1} + \tilde{F}^{(1)} \left( \boldsymbol{X}_{:t}^{i} \right) \right).$$

Furthermore, we have that

$$\begin{split} \frac{d\boldsymbol{W}_{s}^{E}}{dt} = & \frac{1}{NL} \sum_{i=1}^{N} \sum_{t=1}^{L-1} \boldsymbol{W}^{U,T} \left( \boldsymbol{e}_{\boldsymbol{X}_{t+1}^{i}} - \boldsymbol{p}_{t}^{i} \right) \odot \left( \delta_{\boldsymbol{X}_{t}^{i}=s} \boldsymbol{1} + \tilde{F}^{(1)} \left( \boldsymbol{X}_{:t}^{i} \right) \right) \\ = & \frac{1}{NL} \boldsymbol{W}^{U,T} \sum_{i=1}^{N} \sum_{t=1}^{L-1} \delta_{\boldsymbol{X}_{t}^{i}=s} \boldsymbol{e}_{\boldsymbol{X}_{t+1}^{i}} + \frac{1}{NL} \boldsymbol{W}^{U,T} \sum_{i=1}^{N} \sum_{t=1}^{L-1} \boldsymbol{e}_{\boldsymbol{X}_{t+1}^{i}} \odot \tilde{F}^{(1)} \left( \boldsymbol{X}_{:t}^{i} \right) \\ & - \frac{1}{NL} \sum_{i=1}^{N} \sum_{t=1}^{L-1} \boldsymbol{W}^{U,T} \boldsymbol{p}_{t}^{i} \odot \left( \delta_{\boldsymbol{X}_{t}^{i}=s} \boldsymbol{1} + \tilde{F}^{(1)} \left( \boldsymbol{X}_{:t}^{i} \right) \right). \end{split}$$

Since the small initialization, assuming that  $||\boldsymbol{W}||_{\infty} = \mathcal{O}\left(d^{-\gamma}\right)$  for any trainable parameter matrix  $\boldsymbol{W}$ , we have that  $||\tilde{F}^{(1)}\left(\boldsymbol{X}_{:t}^{i}\right)||_{\infty} = \mathcal{O}\left(d^{1-2\gamma}\right)$  in the initial stage. Let  $N \to \infty$ , we have that

$$\frac{d\mathbf{W}_{s}^{E}}{dt} = r_{s}^{\text{in}}\mathbf{W}^{U,T} \left( \boldsymbol{\phi}_{s}^{\text{next}} - \boldsymbol{\eta}^{E} \right),$$

where  $m{\eta}^E = \sum_{t=1}^{L-1} \mathbb{E}_{\pi} \left[ m{p} \mid m{X}_t = s 
ight] + \mathcal{O} \left( d^{1-2\gamma} m{\phi}_s^{ ext{next}} 
ight)$ . Similarly, we have that

$$\frac{d\boldsymbol{W}_{s}^{U}}{dt} = \frac{1}{NL} \sum_{i=1}^{N} \sum_{t=1}^{L-1} \left( \delta_{\boldsymbol{X}_{t+1}^{i}=s} - \boldsymbol{p}_{\boldsymbol{X}_{:t}^{i}}^{i,s} \right) \left( \boldsymbol{W}_{\boldsymbol{X}_{t}^{i}}^{E,T} + \tilde{F} \left( \boldsymbol{X}_{:t}^{i} \right)^{T} \right),$$

where  $p_{m{X}_{:t}^i}^{i,s}$  means the s-th element of the output probability with input sequence  $m{X}_{:t}^i$ . Let  $N o \infty$ , we have

$$rac{doldsymbol{W}_{s}^{U}}{dt} = r_{s}^{ ext{out}} \left( oldsymbol{W}^{E} oldsymbol{arphi}_{s}^{ ext{pre}} 
ight)^{T} + oldsymbol{\eta}^{U},$$

where 
$$oldsymbol{\eta}^U = \sum_{t=1}^{L-1} \mathbb{E}_{\pi} \left[ oldsymbol{p}_{oldsymbol{X}_{:t}}^{E,T} oldsymbol{W}_{oldsymbol{X}_{t}}^{E,T} \right] + \mathcal{O}\left(r_{s}^{\mathrm{out}} d^{1-2\gamma} \left(oldsymbol{W}^{E} oldsymbol{arphi}_{s}^{\mathrm{pre}} 
ight)^{T} \right).$$