# A Law of Next-Token Prediction in Large Language Models

Hangfeng He[*]        Weijie J. Su[†]

August 2024; Revised August 2025

## Abstract

Large language models (LLMs) have been widely employed across various application domains, yet their black-box nature poses significant challenges to understanding how these models process input data internally to make predictions. In this paper, we introduce a precise and quantitative law that governs the learning of contextualized token embeddings through intermediate layers in pre-trained LLMs for next-token prediction. Our findings reveal that each layer contributes equally to enhancing prediction accuracy, from the lowest to the highest layer—a universal phenomenon observed across a diverse array of open-source LLMs, irrespective of their architectures or pre-training data. We demonstrate that this law offers new perspectives and actionable insights to inform and guide practices in LLM development and applications, including model scaling, pre-training tasks, and interpretation.

## 1 Introduction

The rapid advancement of large language models (LLMs) has profoundly impacted various fields, including mathematical discovery [36], medical diagnosis [6], genomic research [21, 16], and education [44]. Despite their transformative and widespread adoption, the deployment of LLMs is often impeded by a lack of understanding of how these enormous, complex black-box models internally process data to generate predictions [34]. Without understanding the prediction mechanisms, practitioners face challenges in interpreting these predictions for decision-making. For LLM developers, this lack of transparency hinders the development of general and robust design principles for LLMs. Consequently, these challenges constrain the full realization of the potential offered by LLM methodologies.

The primary difficulty arises from the hierarchical nature of LLM architectures, such as Transformer [45], RWKV [31], and Mamba [14]. These architectures are composed of multiple layers, each corresponding to simple functions such as linear or quadratic transformations—better known as the attention mechanism—or a nonlinear combination of both. While the input and output can be observed and specified, the internal transformation of data by each layer becomes elusive due to the hierarchical composition. Specifically, in the case of generative pre-trained transformers (GPT) [32], it is unclear how the embeddings of the input text are progressively transformed into features across different layers for next-token prediction, where a token refers to a word or subword. This generative nature of LLMs introduces additional complexity compared to traditional classification

[*]University of Rochester. Email: `hangfeng.he@rochester.edu`.
[†]University of Pennsylvania. Email: `suw@wharton.upenn.edu`.

tasks in multilayer perceptrons (MLP): the vocabulary size typically exceeds the embedding dimension, and models undergo relatively few training epochs. The goal of this paper is to shed light on the inner workings of LLMs as they process token embeddings across all layers. In particular, we aim to identify universal and quantitative patterns that can provide useful principles and insights to refine training processes and enhance interpretability in LLMs.

In this paper, we introduce a quantitative and precise characterization of how LLMs learn contextualized token embeddings for next-token prediction across all layers. As illustrated in Fig. 1, our extensive experiments demonstrate that LLMs enhance their ability to predict the next token according to an exponential law, where each layer improves token prediction by approximately an equal multiplicative factor from the first layer to the last. We refer to this as the law of equi-learning. This law is consistently observed across a wide range of open-source LLMs, including those based on the Transformer architecture and more recent architectures like Mamba and RWKV. Specifically, our experiments reveal the emergence of this law in GPT-1 [32], GPT-2 [33], Llama-1 [42], Llama 2 and its fine-tuned variant Llama 2-Chat [43], Llama 3 and its instruction-fine-tuned version Llama 3 Instruct [9], Mistral 7B and its fine-tuned version Mistral 7B-Instruct [23], phi-1.5 [25], phi-2 [22], phi-3 [1], RWKV and its chat version RWKV-Raven [31], and Mamba [14].

While one might intuitively expect that different token embeddings would be progressively differentiated across the layers of LLMs (an example is shown in Fig. 2), it is remarkable that a universal and geometric law governing tens of thousands of tokens emerges in models of such immense complexity. Moreover, the equi-learning law is perhaps the simplest geometric pattern that could arise across intermediate layers, and what is striking is that this pattern is indeed the observed reality. Notably, the law emerges naturally during the training process without any explicit constraints designed to induce its appearance.

The equi-learning law suggests that every layer should be considered equally important in characterizing the formation of features from input embeddings. In particular, the quantitative nature of this law implies that the layer at the midpoint of the model is precisely where the LLM has achieved half of its overall capability in predicting the next token. This finding challenges the view that feature learning can be disproportionately attributed to certain layers over others [41, 26]. Moreover, the equi-learning law provides practical guidelines and insights into several empirical aspects of LLM training. For instance, this law enables a fine-grained understanding of how the overall capabilities of an LLM relate to its depth, leading to a more nuanced perspective on model scaling that goes beyond what is captured by test loss alone. Additionally, this law sheds light on the superiority of next-token prediction—the currently dominant pre-training task—over alternative training approaches employed in models such as BERT [8], RoBERTa [27] and T5 [35].

## 2  Main Results

LLMs are nonlinear models designed to predict the subsequent token given a sequence of preceding tokens. Formally, a model processes an input sequence of tokens $x_1, x_2, \ldots, x_t, \ldots$. Each token $x_t$ is initially mapped to a vector $\mathbf{h}_{t,0}$ in the embedding space. These embeddings subsequently undergo transformation through a hierarchical stack of model layers—typically comprising attention mechanisms and various operations—yielding a sequence of contextualized token embeddings $\mathbf{h}_{t,\ell}$ at each layer $1 \leq \ell \leq L$. Notably, $\mathbf{h}_{t,\ell}$ is computed using the outputs from the preceding layer $\{\mathbf{h}_{j,\ell-1} \mid 1 \leq j \leq t\}$, which, due to the autoregressive nature of LLMs, contains information exclusively from tokens at the current position and all preceding positions. Finally, the LLM uses the embedding from
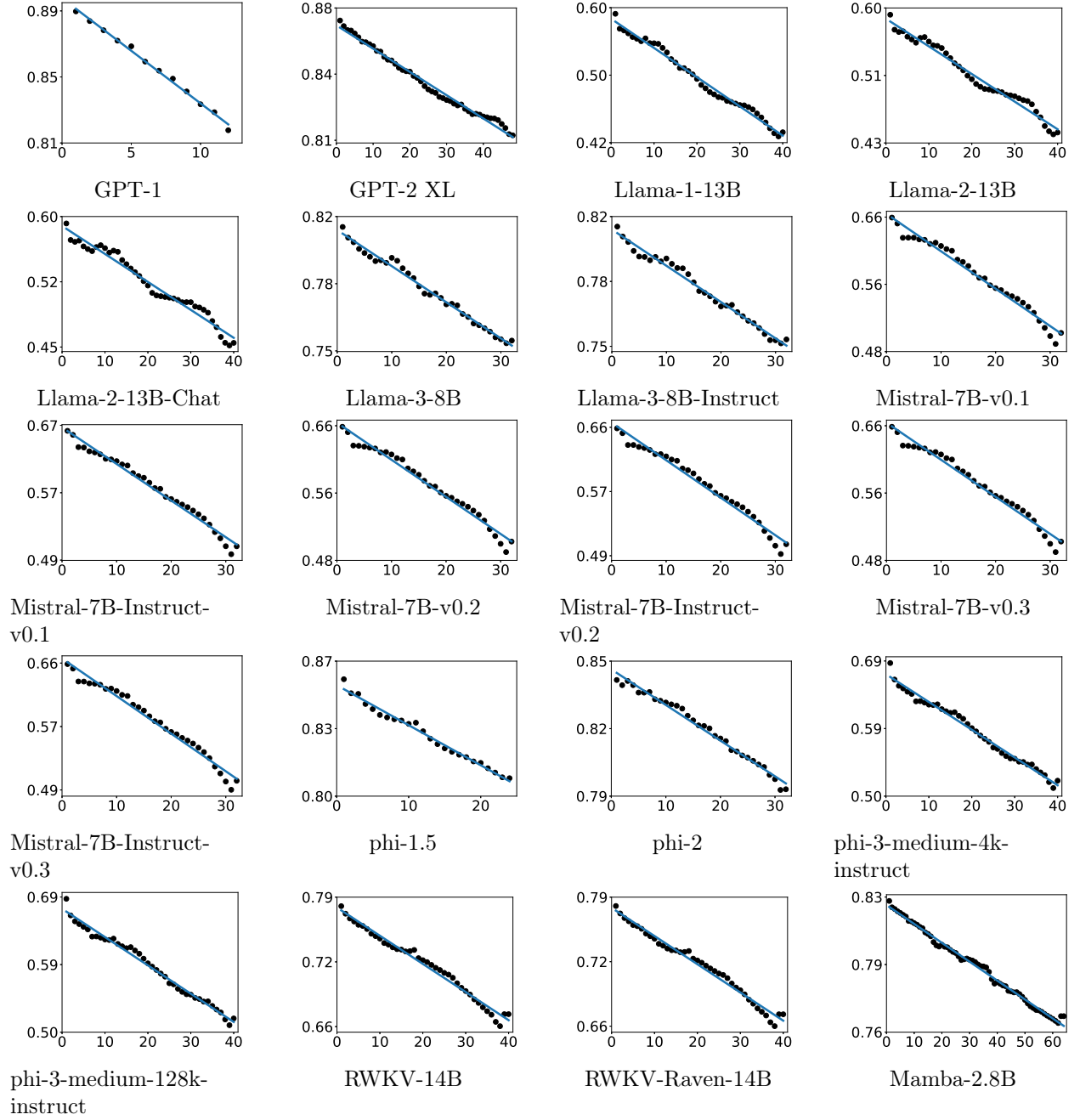
Fig. 1: The law of equi-learning in large language models. Throughout the paper, the x axis represents the layer index, and the y axis, on a logarithmic scale, represents the prediction residual (PR) defined in Eq. 2, unless otherwise specified. The Pearson correlation coefficients, by row first, are -0.997, -0.994, -0.994, -0.988; -0.983, -0.993, -0.992, -0.988; -0.991, -0.989, -0.988, -0.989; -0.988, -0.994, -0.993, -0.992; -0.992, -0.991, -0.991, -0.997. More details can be found in the Supplementary Materials.
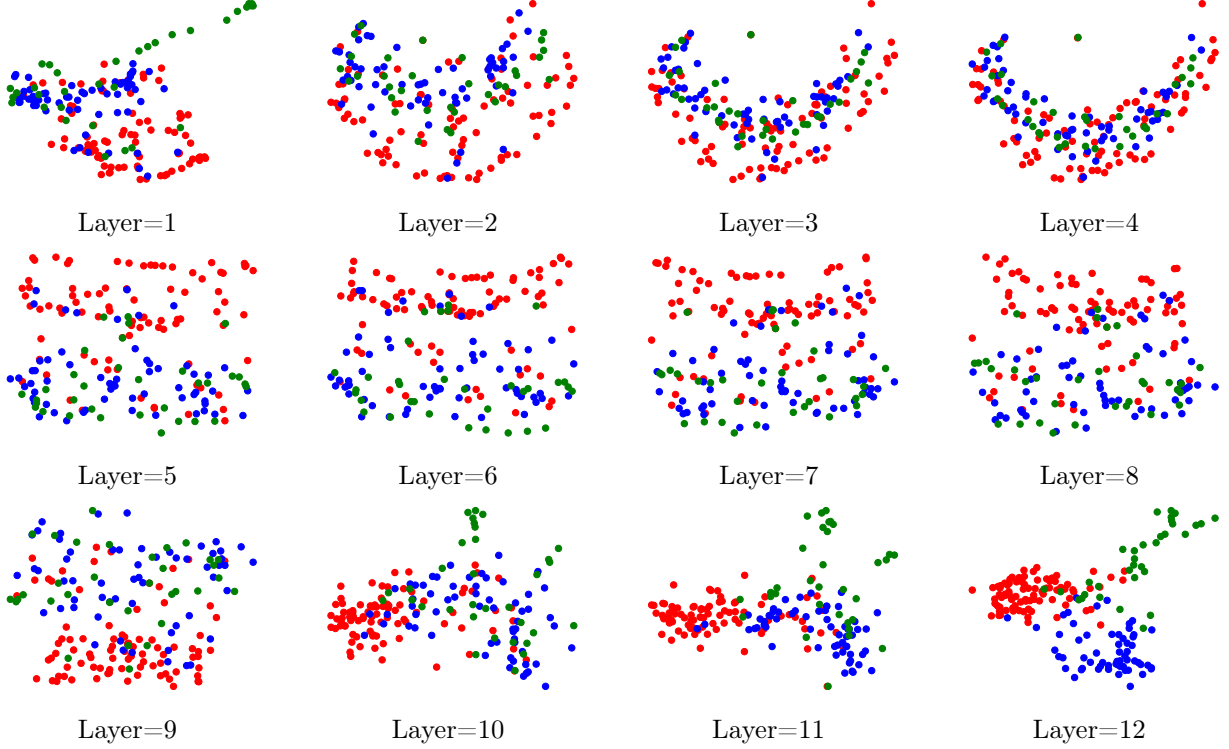
Fig. 2: Intermediate-layer contextualized token embeddings for `patients</w>` (red), `cells</w>` (blue), and `disorder</w>` (green) plotted on the plane of the first two principal components. The embeddings are extracted from the GPT-1 model using 200 sentences sampled from the MedRAG-Textbooks dataset [49] in the medicine domain. As the model progresses from lower to upper layers, the contextualized token embeddings exhibit a clear and progressively increasing separation. The x-axis and y-axis represent the first and second principal components, respectively. More details can be found in the Supplementary Materials.

the last layer, denoted $\mathbf{h}_{t,\text{last}} := \mathbf{h}_{t,L}$, to predict the subsequent discrete token $x_{t+1}$. Aggregating all such pairs across the entire training corpus yields the dataset $\mathcal{D} := \{(\mathbf{h}_{t,\text{last}}^s, x_{t+1}^s) \mid 1 \leq s \leq S\}$.

To assess the capability of the LLM in predicting the next token, we evaluate how well a linear regression model fits on the dataset $\mathcal{D}$. For this purpose, we identify $x$ with its index in the token vocabulary. Let $\hat{x}_{\text{next}} = \mathbf{w} \cdot \mathbf{h} + b$ denote the least-squares fit on $\mathcal{D}$. This suggests using the following metric, which we term the prediction residual (PR), to quantify the LLM's next-token prediction capability:

$$\text{PR} := \frac{\sum (x_{\text{next}} - \hat{x}_{\text{next}})^2}{\sum (x_{\text{next}} - \bar{x}_{\text{next}})^2}, \tag{1}$$

where the sum is over all $x_{\text{next}} = x_{t+1}^s$, $\hat{x}_{\text{next}} = \mathbf{w} \cdot \mathbf{h}_{t,\text{last}}^s + b$, and $\bar{x}_{\text{next}}$ represents the mean of all $x_{t+1}^s$. In statistical terms, this measure is known as the fraction of variance unexplained, equivalent to one minus the coefficient of determination [47]. It serves as a canonical metric for evaluating the proportion of variance in the dependent variable that remains unaccounted for by the independent variables. A high PR value indicates limited predictive power of the token embeddings, while a low value suggests strong predictive capability. Thus, PR inherently captures how effectively the

token embeddings explain next-token prediction[1], aligning with the linear probing paradigm widely employed to analyze the structural properties of contextualized representations in LLMs [18, 26].

To investigate how the predictive power of an LLM with depth $L$ evolves across its layers, we calculate the PR for the next-token prediction task at each intermediate layer. Let $\mathrm{PR}_l$ denote this value for the $l$-th layer, where $1 \leq l \leq L$. Specifically, instead of using the last-layer embedding $\mathbf{h}_{t,\mathrm{last}} \equiv \mathbf{h}_{t,L}$, we use the embedding of the current token at layer $\ell$, denoted $\mathbf{h}_{t,\ell}$, to predict the next token when computing $\mathrm{PR}_\ell$. The law of equi-learning (see Fig. 1) states that the dynamics of predictive power across layers follow the relationship

$$\mathrm{PR}_l \approx \rho^{l-1} \times \mathrm{PR}_1$$

for some decay ratio $0 < \rho < 1$. Since the token embeddings at the input layer (the 0-th layer) are not yet contextualized, they are not included. This implies that $\log \mathrm{PR}_{l+1} - \log \mathrm{PR}_l \approx -\log \frac{1}{\rho}$, indicating a roughly constant reduction in the logarithm of the PR value across each layer, hence the name equi-learning. The Pearson correlation coefficients between the logarithm of the PR value and the layer index range from $-0.983$ to $-0.997$ in Fig. 1. For GPT-1, for example, the decay ratio $\rho$ is approximately 0.993 (as shown in the top-left plot of Fig. 1). In general, $\rho$ depends on various factors, including the model architecture, pre-training data, model depth, feature dimension, and pre-training time.

This law provides what may be the first precise geometric characterization of the learning process for contextualized token embeddings within the intermediate layers of LLMs. Notably, the pre-training objective focuses solely on the last-layer embeddings, aiming to minimize a loss function associated with the last-layer value of PR. Surprisingly, this training dynamics inherently ensure that each layer contributes equally, rather than allowing some layers to carry a disproportionate amount of the workload.

While related phenomena have been observed in MLPs for classification tasks [17, 30], the law presented in this work makes distinct contributions through its focus on LLMs trained for next-token prediction. This task fundamentally differs from classification in several key aspects: it processes sequential data with progressively increasing text lengths and operates within an inherently probabilistic framework without ground truth labels for token predictions. Furthermore, whereas previous studies primarily examined the terminal phase of MLP training, our work demonstrates that the law of equi-learning manifests well before the terminal phase in LLMs, which typically undergo very few epochs of training. This early emergence is particularly noteworthy given the substantial complexity of Transformer-based architectures compared to MLPs. Importantly, this law illuminates the internal mechanisms of the sophisticated LLM architectures [48].

The universality of this law is demonstrated by its emergence across a diverse range of open-source LLMs, spanning different architectures, model sizes, and pre-training datasets. Our investigation employs a comprehensive collection of probing datasets to calculate PR and evaluate the law of equi-learning, including BookCorpus [52], C4 [35], OpenWebText [13], Wikipedia [11], peS2o [39], The Pile [12], Redpajama [7], and OSCAR [40].

**When does the law emerge?** To deepen our understanding of the law's dynamics during training, we investigate the effects of three key factors—training steps, training epochs, and data repetition—on its emergence throughout the process. As illustrated in Fig. 3, the progression of this law as a function of training steps closely resembles the behavior of the equi-separation law observed in MLPs [17]. At model initialization, the PR of contextualized token embeddings for

---

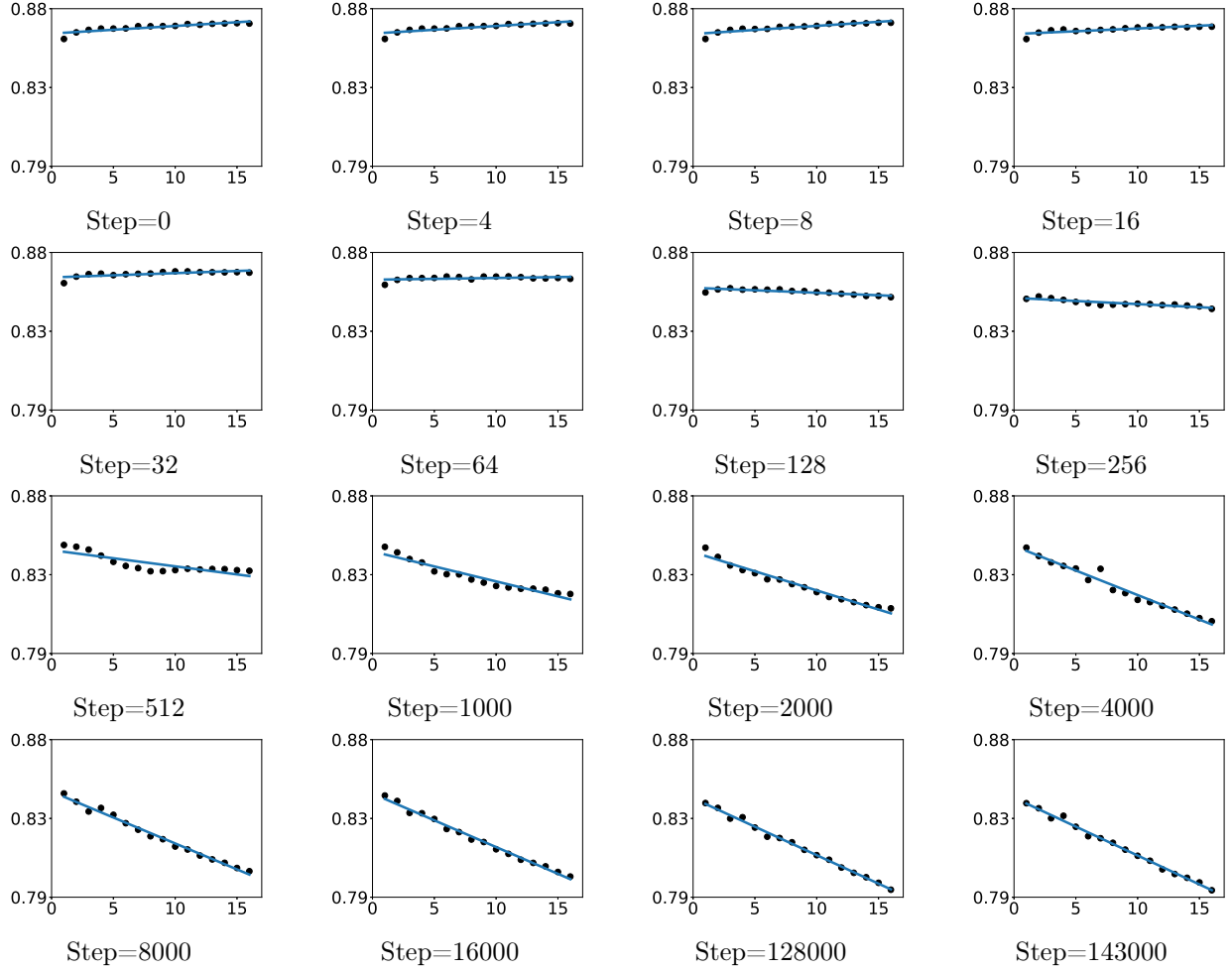[1]See more elaboration on this metric in Section 4.

Fig. 3: Pythia-1B [3] trained on the Pile dataset for various training steps, using a batch size of 2 million, with the total number of steps reaching 143,000. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2. Refer to Fig. S7 in the Supplementary Materials for an enlarged version at initialization (Step=0).

next-token prediction may exhibit an upward trend from lower to higher layers. However, after a certain amount of training (e.g., around 8,000 steps), the law of equi-learning becomes apparent. Beyond this point, the decay ratio continues to decrease until convergence, with the law consistently manifesting during this phase.
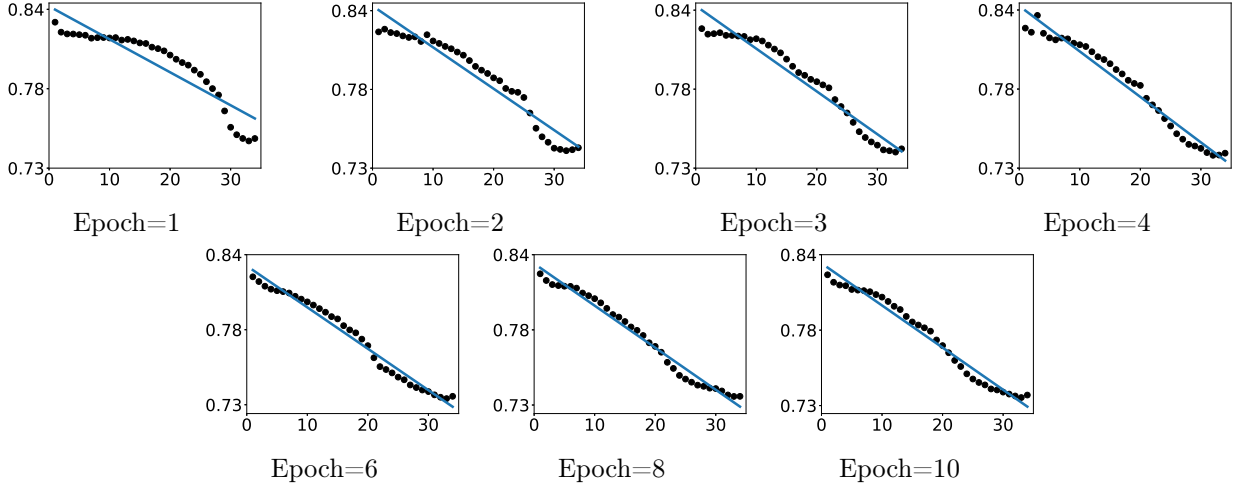


Fig. 4: A 2.8B GPT-2 model pre-trained on a 4 billion token subset of C4 over multiple epochs. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2.
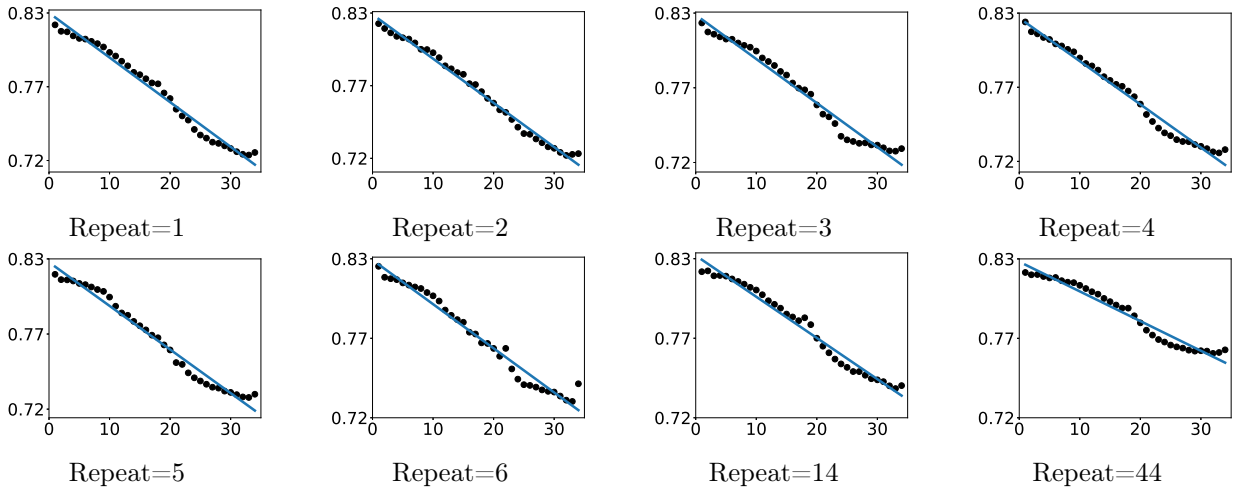


Fig. 5: A 2.8B GPT-2 model pre-trained on varying numbers of unique tokens from the C4 dataset, with a constant total of 55 billion training tokens. Note that the number of unique tokens equals 55 billion divided by the number of repeats. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2.

Although LLMs are typically trained over a limited number of epochs or even a single epoch, we also explore the training dynamics in data-constrained regimes, anticipating that the availability of text data may soon be limited by the finite volume of Internet content [28]. Specifically, we
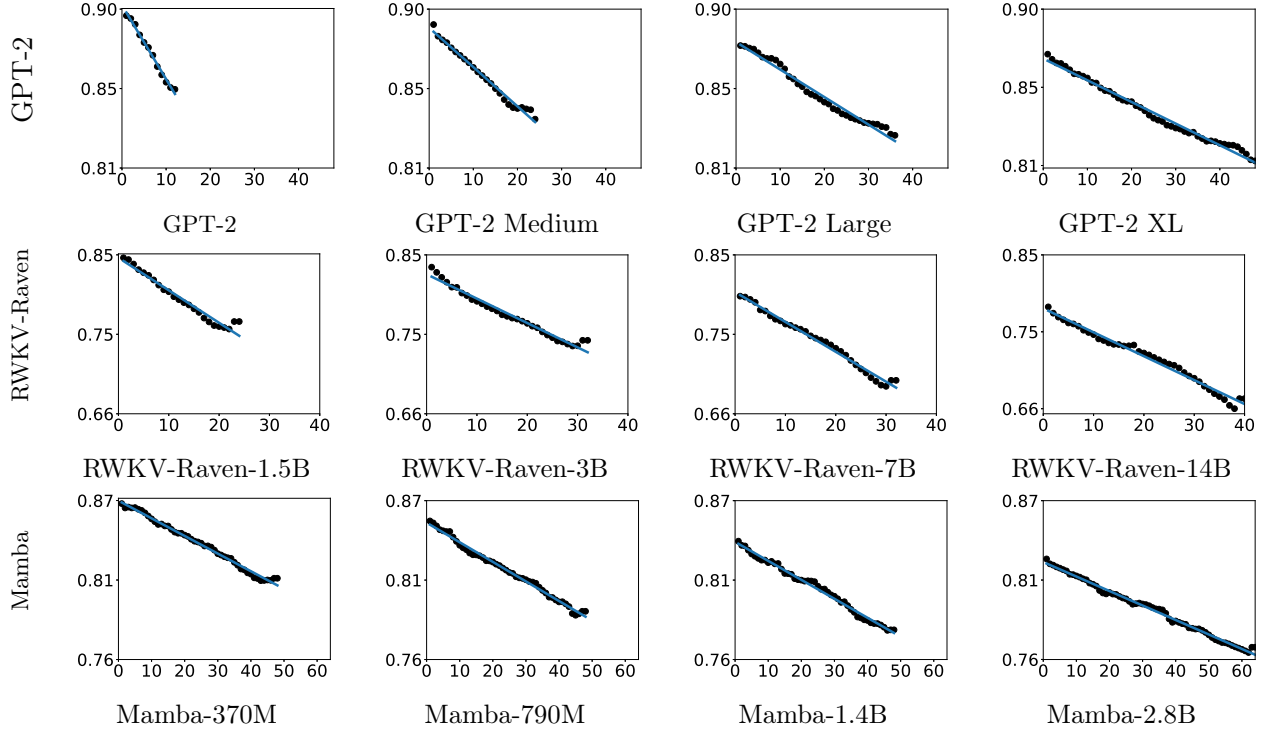
Fig. 6: Illustration of the law of equi-learning with varying model sizes. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2. Please note that the x axis and y axis share the same ranges within the same model series (same row).

utilize various pre-trained 2.8B GPT-2 models released by [28] to examine the impact of training epochs and data repetition. As depicted in Fig.4 and Fig.5, the law of equi-learning emerges as long as the number of epochs is not too small and the number of repeats is not excessively high. Considering these three factors—training steps, training epochs, and data repetition—We observe that a sufficient total number of tokens facilitates the emergence of the equi-learning law, provided that the number of unique tokens is adequate. This finding serves as a crucial condition for training effective LLMs and offers insights into optimizing the training process.

## 3 Perspectives from the Law

The universality of the equi-learning law provides fine-grained perspectives that are applicable to the practical development of LLMs. These perspectives offer new insights into the training processes of LLMs and contribute to advancing transparency in these black-box models. We illustrate the impact of this law on key aspects such as model scaling, pre-training tasks, and information flow. Additional findings, including the impact of pre-training data quality, are discussed in the Supplementary Materials.

**Model scaling.** It is well established that larger models generally yield improved performance [24]. In this part, we explore the impact of model scaling on the observed law. As depicted in Fig. 6, our analysis reveals four consistent trends across three model series, each based on a distinct network architecture: (1) larger models exhibit lower PR values for first-layer token embeddings;

(2) larger models demonstrate smaller PR values at the final layer, indicating enhanced next-token prediction capabilities; (3) larger models demonstrate an increased layer-wise decay ratio ($\rho$); and (4) larger models exhibit a reduced overall decay ratio, defined as the PR of last-layer token embeddings divided by the PR of first-layer token embeddings. The first trend may be attributed to the increased size or improved quality of first-layer token embeddings. The other trends suggest that while larger models typically display enhanced feature learning capabilities, resulting in more refined last-layer token embeddings, their performance at the individual layer level may be less effective compared to that of smaller models. These findings indicate that the law of equi-learning provides a more nuanced perspective on LLM behavior, particularly in the context of model scaling, offering insights that extend beyond the limitations of test loss alone.

**Pre-training task.** In addition to the prevalent next-token prediction task, various other pre-training tasks have been employed in the development of LLMs, including masked language modeling [8] and span corruption [35] (for clarity, these three tasks are abbreviated as NTP, MLM, and SC, respectively). As shown in Fig. 7, different probing tasks are used to analyze the contextualized token embeddings of BERT [8], RoBERTa [27], and T5 [35]. Initially, we employ the mainstream NTP task to examine BERT, RoBERTa, and T5. However, the law of equi-learning does not emerge, possibly due to differences in their pre-training tasks compared to NTP. We then apply their respective pre-training tasks—MLM for BERT and RoBERTa, and SC for T5—to probe the models. For BERT and RoBERTa, under MLM, the law of equi-learning appears but is noisy, likely due to the complexity of the token replacement strategy used during pre-training. Specifically, 15% of tokens were masked and replaced with 80% [MASK] tokens, 10% random tokens, and 10% unchanged tokens for the purpose of predicting the original tokens. For T5, even with SC, the law does not emerge, which may be attributed to its encoder-decoder architecture. In SC, the decoder input often lacks natural coherence and relies heavily on the encoder input, while the decoder's cross-attention layer might further complicate the learning of contextualized token embeddings. These findings suggest that the choice of pre-training task is critical for the emergence of the law of equi-learning, with more naturalistic tasks potentially facilitating its appearance. This may provide supporting evidence for the superiority of NTP, the currently dominant pre-training task.

**Information flow.** Controlling the flow of information is crucial for effective feature learning in sequential data [19, 5]. To elucidate the dynamics of information flow in LLMs across different architectures, we analyze the contextualized token embeddings of the current token ($x_t$) at each layer and their ability to predict other tokens within the sequence, ranging from the previous token ($x_{t-1}$) to the next-next token ($x_{t+2}$), including the default next token ($x_{t+1}$). Specifically, the embedding of the current token at layer $\ell$, denoted as $\mathbf{h}_{t,\ell}$, is used to predict not only the next token ($x_{t+1}$), but also the previous token ($x_{t-1}$), the current token ($x_t$), and the next next token ($x_{t+2}$) when computing $\mathrm{PR}_\ell$. As depicted in Fig. 8, our results reveal a clear pattern: following the learning of contextualized token embeddings across layers, LLMs exhibit a tendency to forget prior information, including the current token itself, as indicated by the positive correlations between PR and layer index. In contrast, the prediction of future tokens improves, evidenced by negative correlations between PR and layer index. These findings suggest that as the learning of contextualized token embeddings progresses from lower to higher layers, LLMs increasingly discard historical information while simultaneously enhancing their predictive capabilities for upcoming tokens.
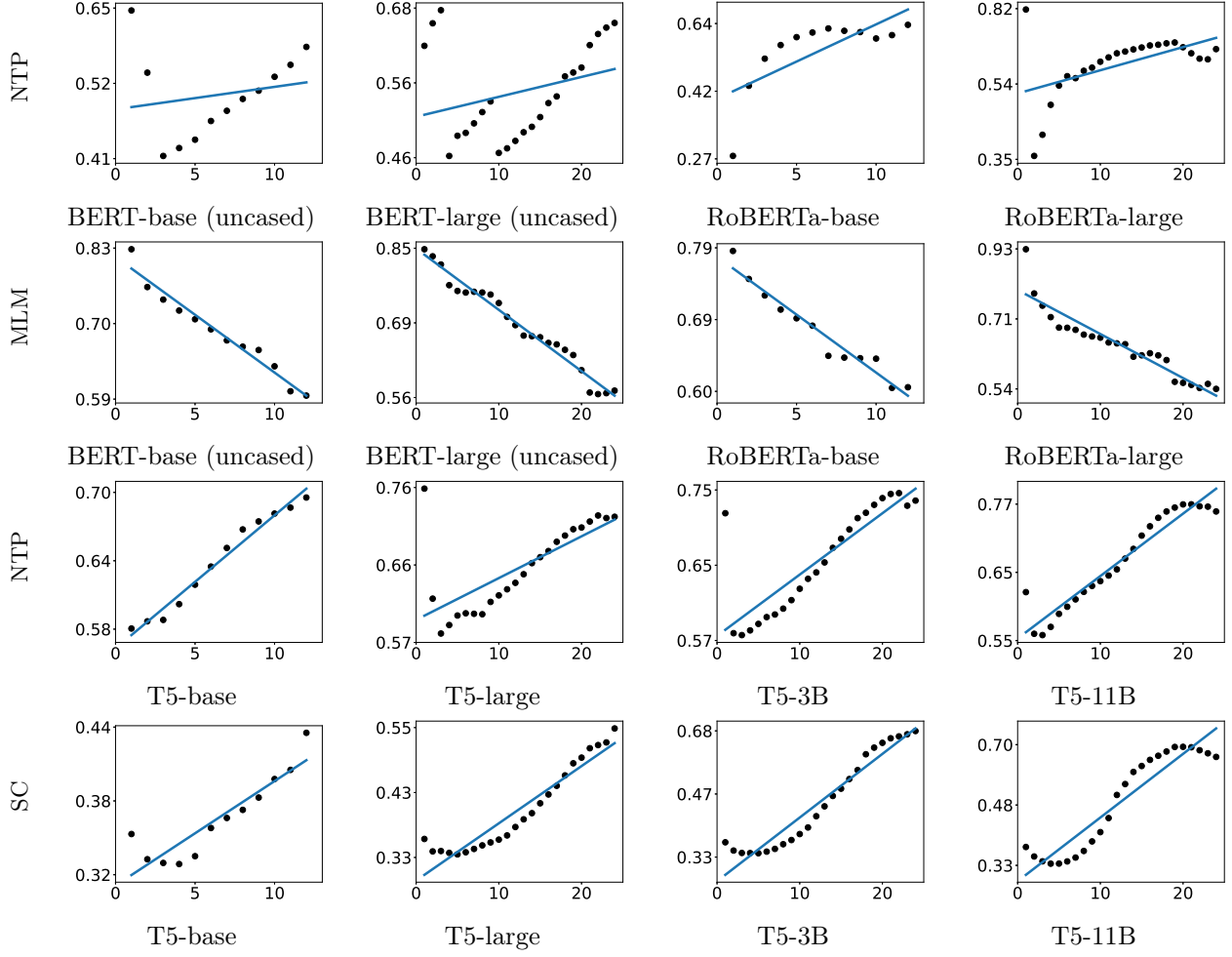
9

Fig. 7: Different probing tasks are used for BERT (uncased), RoBERTa, and T5. Under the next-token prediction (NTP) task, the law of equi-learning does not appear because these models are pre-trained on different tasks. For BERT and RoBERTa, under their pre-training task, masked language modeling (MLM), the law of equi-learning appears but is somewhat noisy. In contrast, for T5, even with its pre-training task of span corruption (SC), the law of equi-learning does not appear. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2. Note that the prediction target in PR is changed from the next token to the masked token for MLM and the corrupted span for SC.
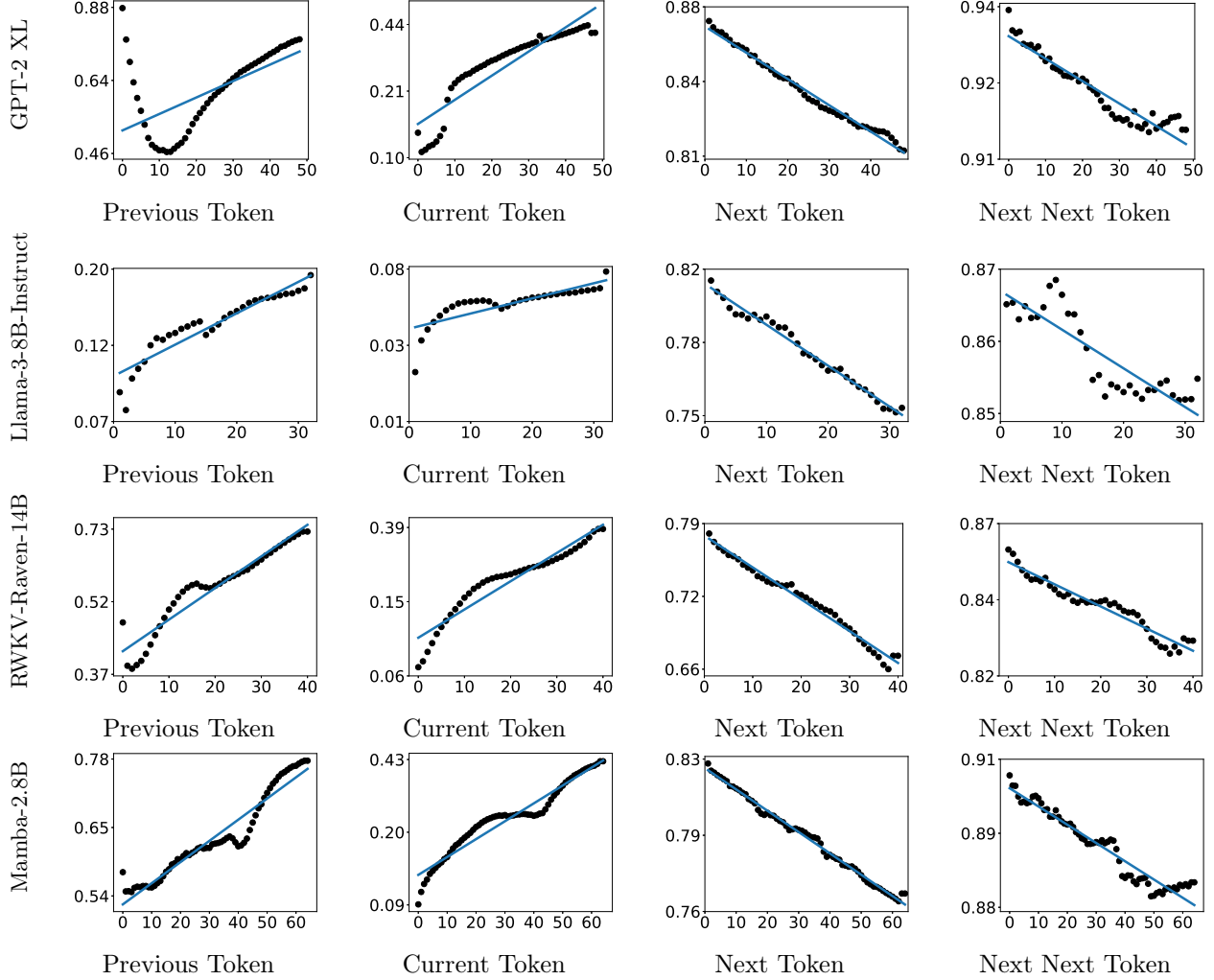
Fig. 8: The contextualized token embeddings of the current token $(x_t)$ at each layer are utilized to predict various tokens in the sequence, ranging from previous token $(x_{t-1})$ to next next token $(x_{t+2})$, including the default next token $(x_{t+1})$. It is observed that LLMs tend to forget previous information (positive correlations between the PR and the layer index), including the current token, and improve their prediction of future tokens (negative correlations between the PR and the layer index) after the learning of contextualized token embeddings across layers. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2. Note that the prediction target in PR is changed from the next token to the previous token, current token, or next next token in settings other than next-token prediction.

## 4   Discussion

Despite the extensive research on structures within pre-trained LLMs, many have found it challenging to identify precise, quantitative laws governing their internal dynamics. In this work, we challenge this view by introducing the law of equi-learning, which describes how contextualized token embeddings evolve from the first to the last layer. This law, which is both quantitative and

11

precise, has been consistently observed across various architectures, including Transformer, Mamba, and RWKV. Its emergence provides crucial insights into the training and interpretation of LLMs, offering new perspectives that deepen our understanding of their internal mechanisms.

The significance of the equi-learning law lies in its potential to refine the development and application of LLMs. An open question is how the decay ratio $\rho$ depends on factors such as model depth and pre-training data. Understanding this dependence could lead to the development of more efficient LLMs by minimizing $\rho^{L-1}$, the overall decay ratio, in the equi-learning law. The emergence of this law specifically under the PR metric defined in (2), rather than alternative metrics (see Fig. S8 in the Supplementary Materials), warrants further investigation. We hypothesize that this specificity may arise from the PR metric's incorporation of token indices derived from byte pair encoding [37], which might capture structural information beyond simple classification metrics. Nevertheless, we believe that the law could be extended beyond the PR metric, though this remains an avenue for future research. The law's eventual emergence also suggests the possibility of setting different learning rates across layers to accelerate convergence to the equilibrium described by the law. Moreover, preserving the equi-learning law during model pruning and fine-tuning may yield practical benefits, potentially through the preservation of certain weights or the use of techniques like LoRA [20].

A central open question is to uncover the mechanism underlying the equi-learning law. Related phenomena have been analytically derived in deep linear networks (DLNs) with linearly separable data [46], but the strong assumptions—such as the absence of nonlinearity and simplified data—limit their relevance to real-world LLMs. A spring-block analogy has been proposed to illustrate a similar law in deep neural networks (DNNs) [38]. However, this framework remains primarily heuristic, as it does not establish a concrete correspondence between fundamental physical elements—such as elastic potential energy and friction—and specific components within DNNs. Furthermore, it does not address the substantial architectural and functional disparities between DNNs and modern LLMs. A formal derivation of the equi-learning law in LLMs will likely require a deeper understanding of model dynamics and the structure of data, aligned with recent geometric approaches to analyzing DNNs and LLMs [30, 48, 4, 50], which we leave for future work.

Broadly, the equi-learning law could be leveraged in transfer learning, particularly when the bottom layers are frozen while the top layers are re-trained to adapt models to new domains. Additionally, our preliminary results in the Supplementary Materials suggest that higher-quality pre-training data may require high-quality probing data to facilitate the emergence of the equi-learning law. This finding implies the potential for using this law to improve the evaluation of LLM capabilities across different tasks. For interpretability, a welcome advance would be the development of new methodologies that consider the collective contributions of all layers, rather than just a few, in interpreting the predictions of LLMs.

## Acknowledgments

# References

[1] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.

[2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

[3] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR, 2023.

[4] Kwan Ho Ryan Chan, Yaodong Yu, Chong You, Haozhi Qi, John Wright, and Yi Ma. Redunet: A white-box deep network from the principle of maximizing rate reduction. *Journal of machine learning research*, 23(114):1–103, 2022.

[5] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NIPS 2014 Workshop on Deep Learning, December 2014*, 2014.

[6] Jan Clusmann, Fiona R Kolbinger, Hannah Sophie Muti, Zunamys I Carrero, Jan-Niklas Eckardt, Narmin Ghaffari Laleh, Chiara Maria Lavinia Löffler, Sophie-Caroline Schwarzkopf, Michaela Unger, Gregory P Veldhuizen, et al. The future landscape of large language models in medicine. *Communications medicine*, 3(1):141, 2023.

[7] Together Computer. Redpajama: An open source recipe to reproduce llama training dataset, 2023.

[8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.

[9] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[10] Cong Fang, Hangfeng He, Qi Long, and Weijie J Su. Exploring deep neural networks via layer-peeled model: Minority collapse in imbalanced training. *Proceedings of the National Academy of Sciences*, 118(43):e2103091118, 2021.

[11] Wikimedia Foundation. Wikimedia downloads.

[12] 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.

[13] Aaron Gokaslan, Vanya Cohen, Ellie Pavlick, and Stefanie Tellex. Openwebtext corpus. http://Skylion007.github.io/OpenWebTextCorpus, 2019.

[14] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In *First Conference on Language Modeling*, 2024.

[15] Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré, Adam Chilton, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel Rockmore, Diego Zambrano, et al. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. *Advances in Neural Information Processing Systems*, 36, 2024.

[16] Minsheng Hao, Jing Gong, Xin Zeng, Chiming Liu, Yucheng Guo, Xingyi Cheng, Taifeng Wang, Jianzhu Ma, Xuegong Zhang, and Le Song. Large-scale foundation model on single-cell transcriptomics. *Nature Methods*, pages 1–11, 2024.

[17] Hangfeng He and Weijie J Su. A law of data separation in deep learning. *Proceedings of the National Academy of Sciences*, 120(36):e2221704120, 2023.

[18] John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, 2019.

[19] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[20] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.

[21] Kaixuan Huang, Yuanhao Qu, Henry Cousins, William A Johnson, Di Yin, Mihir Shah, Denny Zhou, Russ Altman, Mengdi Wang, and Le Cong. CRISPR-GPT:: An LLM agent for automated design of gene-editing experiments. *arXiv preprint arXiv:2404.18021*, 2024.

[22] Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Sebastien Bubeck, Caio César Teodoro Mendes, Weizhu Chen, Allie Del Giorno, Ronen Eldan, Sivakanth Gopi, et al. Phi-2: The surprising power of small language models. *Microsoft Research Blog*, 2023.

[23] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7B. *arXiv preprint arXiv:2310.06825*, 2023.

[24] 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.

[25] Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023.

[26] Nelson F Liu, Matt Gardner, Yonatan Belinkov, Matthew E Peters, and Noah A Smith. Linguistic knowledge and transferability of contextual representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1073–1094, 2019.

[27] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

[28] Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language models. *Advances in Neural Information Processing Systems*, 36, 2024.

[29] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.

[30] Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020.

[31] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Leon Derczynski, et al. RWKV: Reinventing RNNs for the transformer era. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14048–14077, 2023.

[32] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training, 2018.

[33] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners, 2019.

[34] Adityanarayanan Radhakrishnan, Daniel Beaglehole, Parthe Pandit, and Mikhail Belkin. Mechanism for feature learning in neural networks and backpropagation-free machine learning models. *Science*, 383(6690):1461–1467, 2024.

[35] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.

[36] Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.

[37] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, 2016.

[38] Cheng Shi, Liming Pan, and Ivan Dokmanić. A spring-block theory of feature learning in deep neural networks. *arXiv preprint arXiv:2407.19353*, 2024.

[39] Luca Soldaini and Kyle Lo. peS2o (Pretraining Efficiently on S2ORC) Dataset. Technical report, Allen Institute for AI, 2023. ODC-By, `https://github.com/allenai/pes2o`.

[40] Pedro Ortiz Suarez, Laurent Romary, and Benoît Sagot. A monolingual approach to contextualized word embeddings for mid-resource languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1703–1714, 2020.

[41] Ian Tenney, Dipanjan Das, and Ellie Pavlick. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, 2019.

[42] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[43] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[44] Xinming Tu, James Zou, Weijie Su, and Linjun Zhang. What should data science education do with large language models. *Harvard Data Science Review*, 6(1), jan 19 2024. https://hdsr.mitpress.mit.edu/pub/pqiufdew.

[45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[46] Peng Wang, Xiao Li, Can Yaras, Zhihui Zhu, Laura Balzano, Wei Hu, and Qing Qu. Understanding deep representation learning via layerwise feature compression and discrimination. *arXiv preprint arXiv:2311.02960*, 2023.

[47] Sanford Weisberg. *Applied Linear Regression*. Wiley Series in Probability and Statistics. Wiley, 2005.

[48] Robert Wu and Vardan Papyan. Linguistic collapse: Neural collapse in (large) language models. *Advances in Neural Information Processing Systems*, 37:137432–137473, 2024.

[49] Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. Benchmarking retrieval-augmented generation for medicine. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 6233–6251, 2024.

[50] Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Hao Bai, Yuexiang Zhai, Benjamin D Haeffele, and Yi Ma. White-box transformers via sparse rate reduction: compression is all there is? *Journal of Machine Learning Research*, 25(300):1–128, 2024.

[51] Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.

[52] Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27, 2015.

# Supplementary Text

This section outlines the general experimental setup, distinctive configurations for main text figures, and additional results, with comprehensive details available in our code repository[2].

## General Setup

In this subsection, we detail the general experimental setup utilized throughout this study.

**LLMs.** In this study, we focus on autoregressive LLMs, where the objective is to predict the subsequent token in a sequence, constrained to attend solely to preceding tokens. Formally, the model takes as input a sequence of discrete tokens $x_1, x_2, \ldots, x_t \in \mathcal{V}^t$, where $\mathcal{V}$ denotes the vocabulary specific to the model. These tokens are obtained via a tokenizer applied to raw text; for instance, the sentence *"We love Physics."* is tokenized by the Llama 3 tokenizer into four tokens: "We", " love", " Physics", and "." Vocabulary sizes vary across models (e.g., 32K for Llama 2 and 128K for Llama 3). Each token $x_i$ is mapped to an initial embedding vector $\mathbf{h}_{i,0}$ via a learned embedding table. These embeddings are then propagated through a stack of model layers—typically transformer layers—resulting in a set of contextualized token embeddings $\mathbf{h}_{i,\ell}$ at each layer $1 \leq \ell \leq L$. Due to the autoregressive nature of the model, the contextualized token embeddings $\mathbf{h}_{i,\ell}$ is computed using only information from position $i$ and all positions preceding it, specifically from the previous layer's outputs $\{\mathbf{h}_{j,\ell-1} \mid 1 \leq j \leq i\}$. This ensures that the model does not access future tokens. The last-layer embedding of the current token, $\mathbf{h}_{t,L}$ (also denoted $\mathbf{h}_{t,\text{last}}$), is used to predict the next token $x_{t+1}$ by projecting it into vocabulary space and applying a softmax to produce a probability distribution over the vocabulary. This autoregressive decoding procedure is iteratively applied at each position $1 \leq t \leq T-1$, where $T$ denotes the length of the input token sequence.

**Prediction residual (PR).** To assess the capability of contextualized token embeddings in predicting the next token, we calculate the fraction of variance unexplained (FVU) by a linear regression model predicting the next token. Notably, this metric is similar to the concurrent measure proposed by [38], specifically the Root Mean Squared Error (RMSE) of the optimal linear regressor based on the features in MLPs. In their study, RMSE was shown to reproduce a similar noise–nonlinearity phase diagram in MLP training under regression, extending phenomena originally observed in classification to the regression setting. This finding supports the potential of PR as a meaningful measure of representation quality.

**Layer normalization.** In pre-layer normalization (pre-LN) models, default initialized layer normalization is applied to all layers except the last-layer token embeddings. This is because layer normalization is moved to the input of each sub-block, with an additional layer normalization added after the final self-attention block [33]. Given that different models utilize distinct forms of layer normalization—such as LayerNorm [2] in GPT-2 and RMSNorm [51] in Llama-1—we will apply the specific layer normalization technique used by each model to normalize its contextualized token embeddings. Throughout this paper, nearly all models are pre-LN models, with the exceptions of GPT-1, BERT, and RoBERTa.

**Probing datasets.** For our experiments, we consider eight distinct probing datasets: BookCorpus, C4, OpenWebText, Wikipedia, peS2o, Pile, Redpajama, and OSCAR. We sampled sentences based on their average length, extracting $3,000$ sentences from BookCorpus, 200 from C4, and 100 from each of the remaining datasets. For consistency, sentences were truncated to a maximum

---

[2]Our code is publicly available at `https://github.com/HornHehhf/LLM-ELL`.

length of 512 tokens across all datasets, except for BookCorpus. Unless otherwise specified, we will utilize the probing dataset that exhibits the strongest Pearson correlation in next-token prediction for each LLM.

## Detailed Experimental Settings

In this subsection, we show detailed experimental settings for the figures in the main text.

**Large language models.** As depicted in Fig. 1, the law of equi-learning is observed in various open-source large language models. Please note that phi-1 was excluded from our analysis, as it is a LLM specifically designed for code. For each model, we evaluate the largest size that can be executed on a local machine equipped with two L40S GPUs, each possessing 48 GB of memory. The corresponding model sizes are as follows: 117M for GPT-1, 1.5B for GPT-2, 13B for Llama-1, 13B for Llama 2 and Llama 2-Chat, 8B for Llama 3 and Llama 3 Instruct, 7B for three versions of Mistral 7B and Mistral 7B-Instruct (v0.1, v0.2, v0.3), 1.3B for phi-1.5, 2.7B for phi-2, 14B for phi-3 (phi-3-medium) with varying context lengths (4K, 128K), 14B for RWKV and RWKV-Raven, and 2.8B for Mamba. The corresponding probing datasets used are as follows: BookCorpus for GPT-1, BookCorpus for GPT-2, peS2o for Llama-1, peS2o for Llama 2 and Llama 2-Chat, BookCorpus for Llama 3 and Llama 3 Instruct, C4 for three versions of Mistral 7B and Mistral 7B-Instruct (v0.1, v0.2, v0.3), BookCorpus for phi-1.5, BookCorpus for phi-2, C4 for phi-3 (phi-3-medium) with varying context lengths (4K, 128K), C4 for RWKV and RWKV-Raven, and BookCorpus for Mamba. Notably, reinforcement learning from human feedback (RLHF) [29] does not significantly impact the law, as demonstrated by Llama-2-13B, Llama-3-8B, Mistral-7B-v0.1, Mistral-7B-v0.2, Mistral-7B-v0.3, and RWKV, along with their fine-tuned versions, as illustrated in Fig. 1.

**Visualization.** Figs. S1, S2, and S3 depict the visualization of contextualized token embeddings across various layers of the GPT-1 model, using $3,000$ sentences from the BookCorpus (consistent with the setting in Fig.1). We examine three token pairs: `they</w>` vs. `them</w>`, `have</w>` vs. `had</w>`, and `are</w>` vs. `is</w>`. Principal component analysis (PCA) is applied to project these contextualized token embeddings onto a two-dimensional plane. The results show a distinct and progressive separation of contextualized token embeddings within each pair as one moves from the lower to the upper layers of the model. Moreover, as illustrated in Fig. S4, the law of equi-learning manifests when evaluating the GPT-1 model across distinct domains. Specifically, we analyzed 200 sentences from MedRAG-Textbooks [49] for the medicine domain, 1000 sentences from LegalBench [15] for the law domain, and 500 sentences from US-Congressional-Speeches[3] for the politics domain. Consistent with other probing datasets, the number of sampled sentences was determined by their average length, and sentences were truncated to a maximum length of 512 tokens. For each domain, we visualized the embeddings of token sets across various layers of the GPT-1 model—`patients</w>`, `cells</w>`, and `disorder</w>` for medicine; `law</w>` and `policy</w>` for law; and `president</w>` and `country</w>` for politics. Similarly, PCA was employed to project these contextualized token embeddings onto a two-dimensional plane. As depicted in Fig. 2, Fig. S5, and Fig. S6, the results reveal a distinct and progressively increasing separation of contextualized token embeddings within each token set as the model transitions from lower to upper layers.

**Training dynamics.** As illustrated in Figs. 3, 4, and 5, we examine the evolution of the observed law throughout the training process. Specifically, Fig. 3 depicts the PR of contextualiezed token embeddings at each layer of Pythia-1B at various training steps (an enlarged version at

---

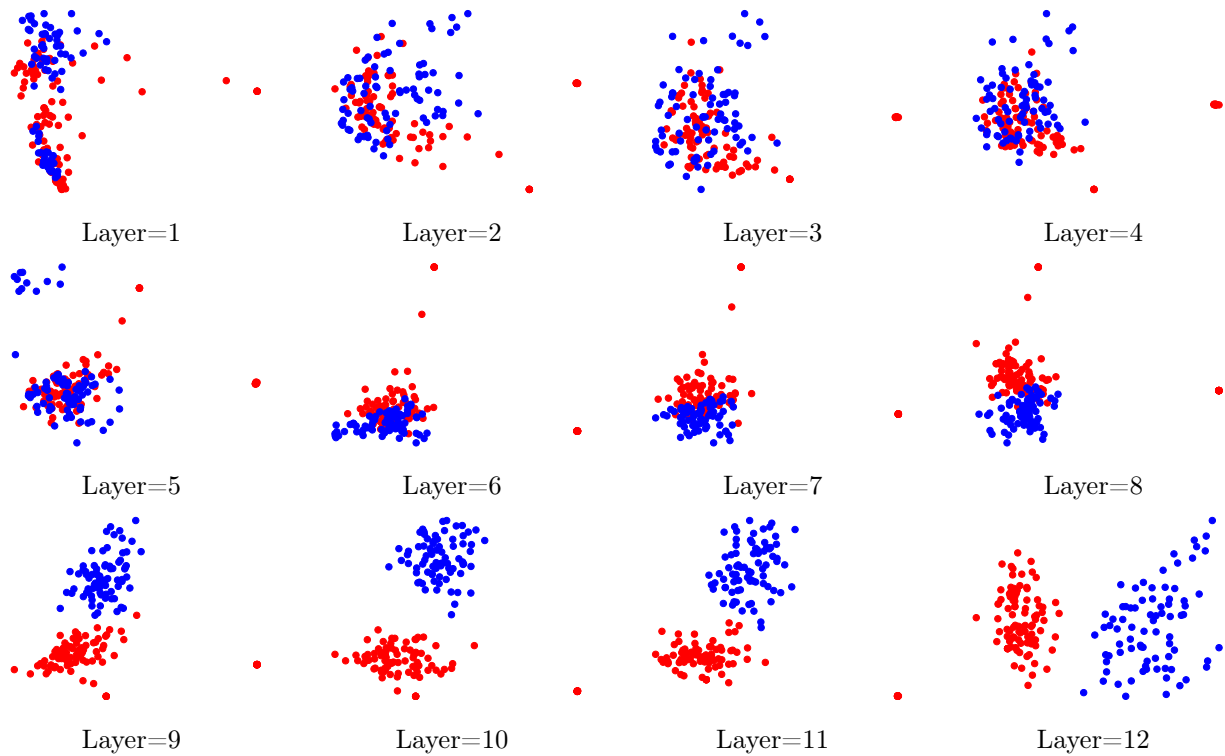[3]See more in `https://huggingface.co/datasets/Eugleo/us-congressional-speeches`.

Fig. S1: Intermediate-layer contextualized token embeddings for they</w> (red) and them</w> (blue) plotted on the plane of the first two principal components. The x-axis and y-axis represent the first and second principal components, respectively.
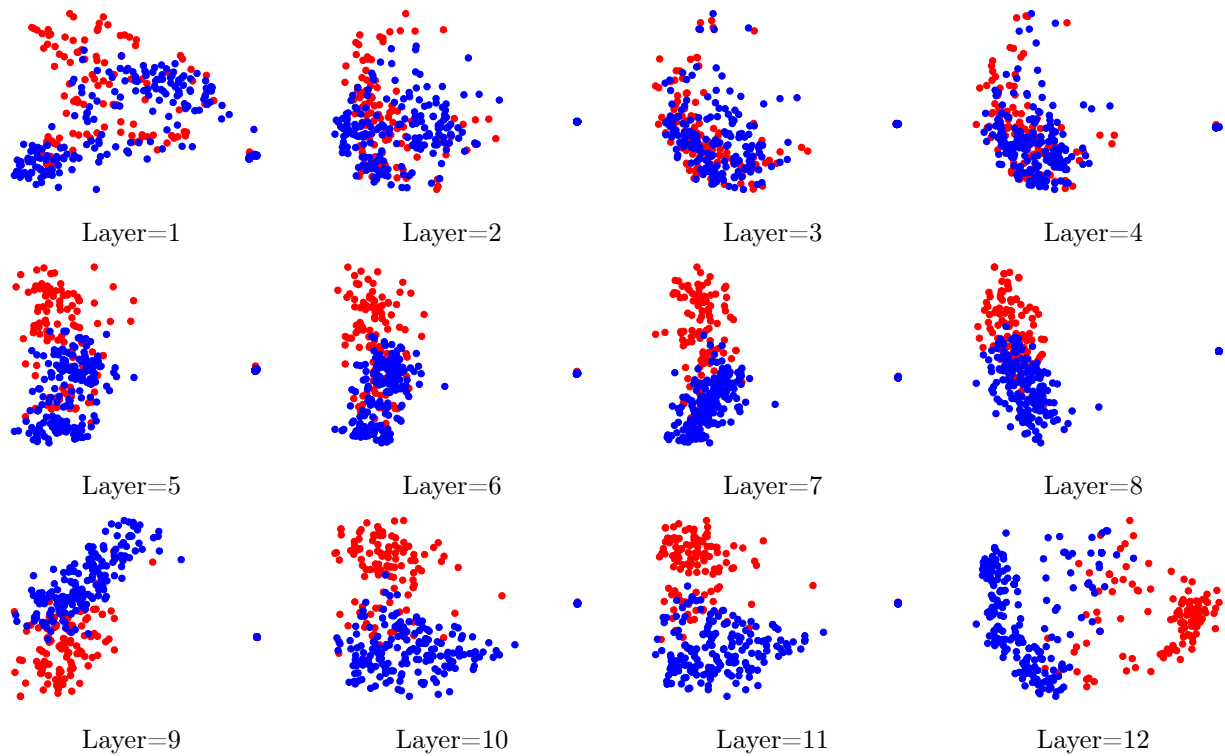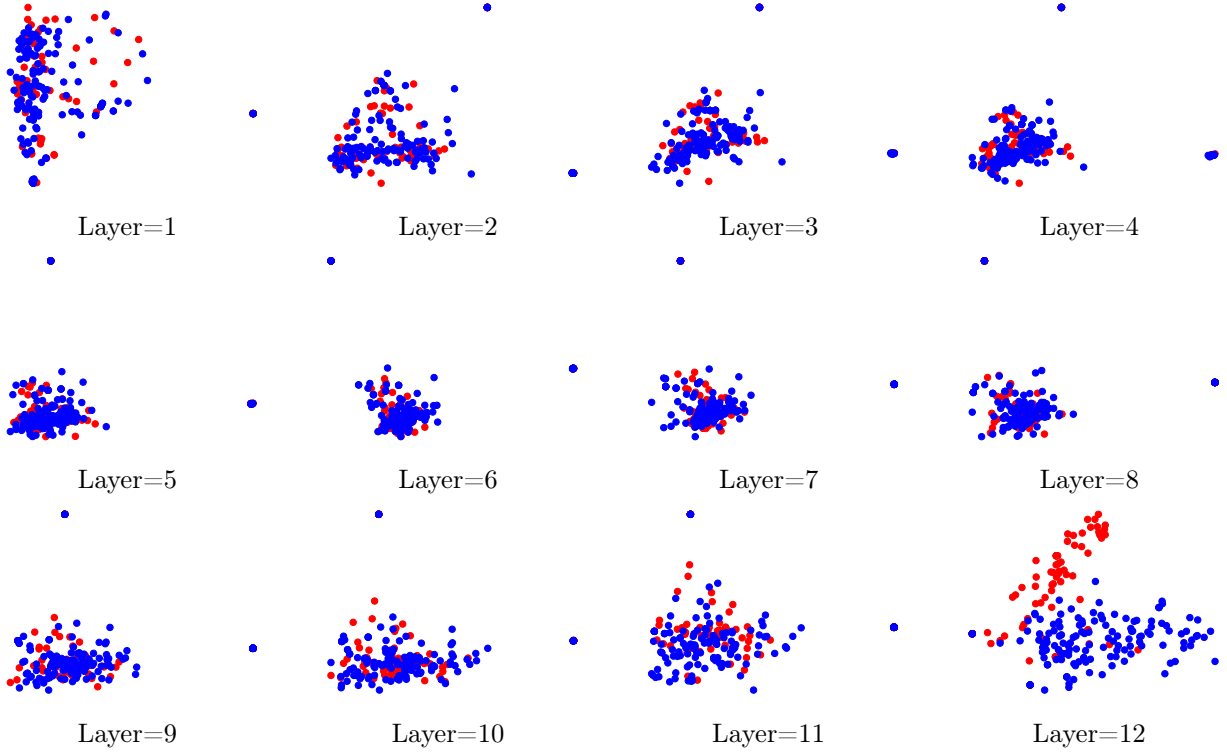
Fig. S2: Intermediate-layer contextualized token embeddings for have</w> (red) and had</w> (blue) plotted on the plane of the first two principal components. The x-axis and y-axis represent the first and second principal components, respectively.

Fig. S3: Intermediate-layer contextualized token embeddings for `are</w>` (red) and `is</w>` (blue) plotted on the plane of the first two principal components. The x-axis and y-axis represent the first and second principal components, respectively.
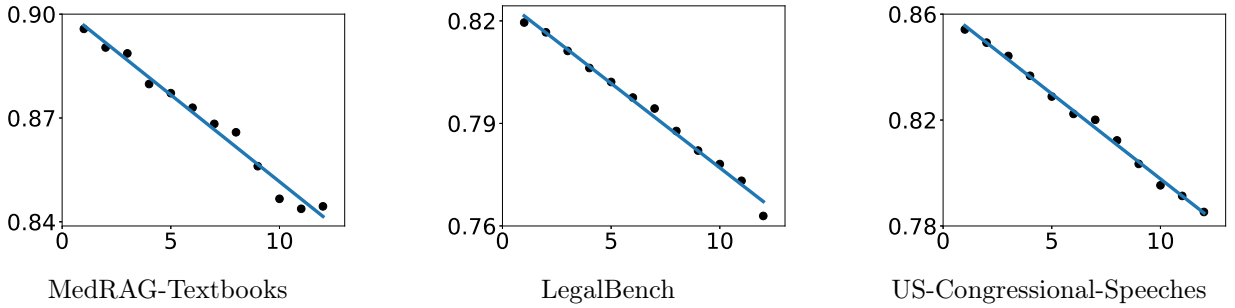


Fig. S4: The law of equi-learning emerges when GPT-1 is evaluated across different domains, including medicine (MedRAG-Textbooks), law (LegalBench), and politics (US-Congressional-Speeches), with corresponding Pearson correlation coefficients of $-0.990$, $-0.995$, and $-0.998$, respectively. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2.
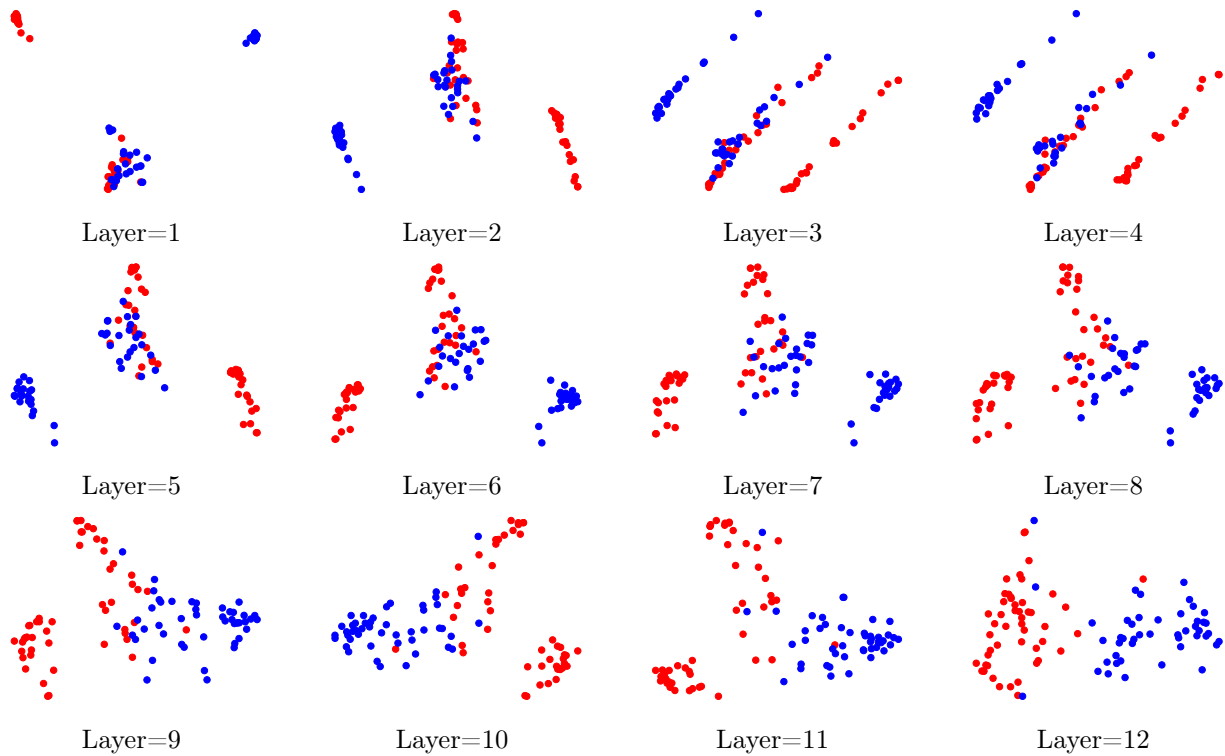
Fig. S5: Intermediate-layer contextualized token embeddings for law</w> (red) and policy</w> (blue) plotted on the plane of the first two principal components. The x-axis and y-axis represent the first and second principal components, respectively.
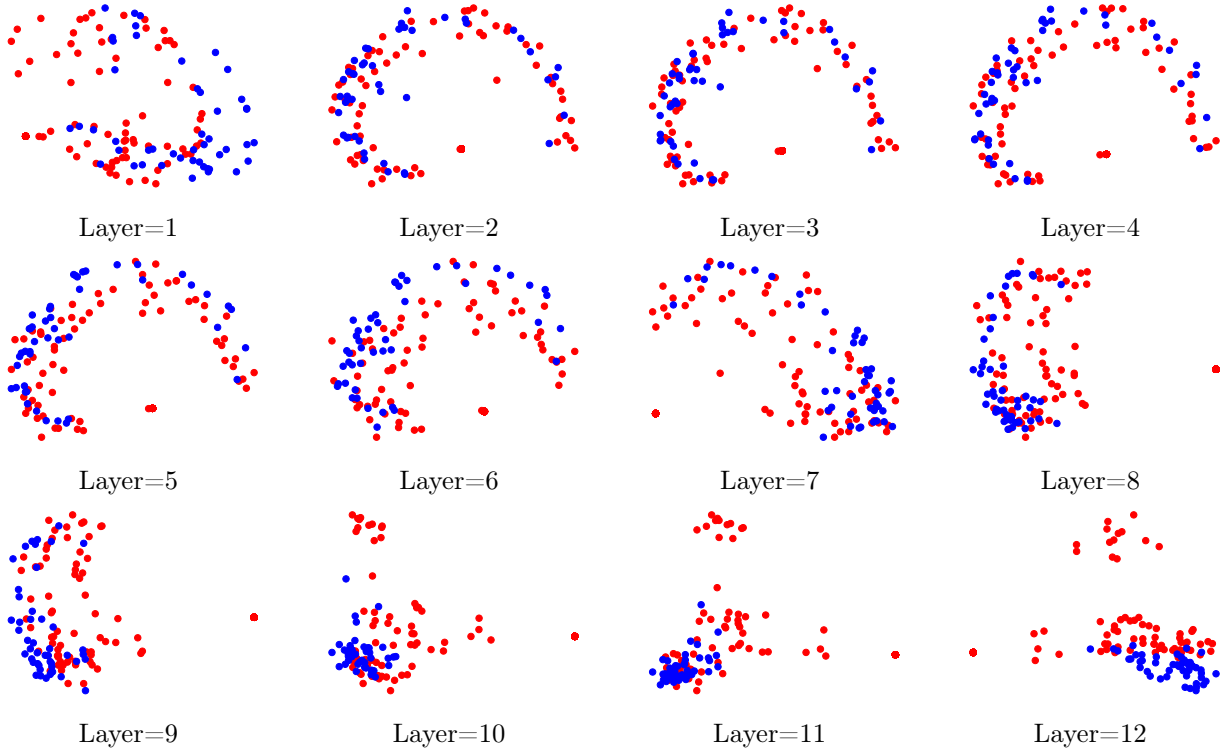
Fig. S6: Intermediate-layer contextualized token embeddings for `president</w>` (red) and `country</w>` (blue) plotted on the plane of the first two principal components. The x-axis and y-axis represent the first and second principal components, respectively.
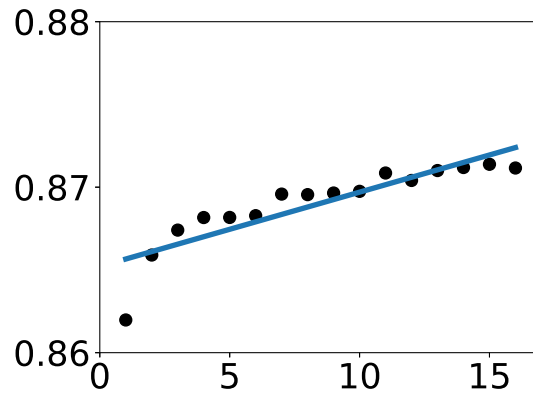


Fig. S7: An enlarged version of Pythia-1B at initialization (Step=0). The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2.

initialization (Step=0) is in Fig. S7), using BookCorpus as the probing dataset. Additionally, Figs. 4 and 5 present the PR of contextualized token embeddings at each layer of pre-trained 2.8B GPT-2 models, as released by [28], across different training epochs and data repetitions, with OpenWebText serving as the probing dataset.

**Model scaling.** As illustrated in Fig. 6, the contextualized token embeddings are compared across different model sizes within the same model series. We analyze three distinct model series—GPT-2, RWKV-Raven, and Mamba—each with at least four different model sizes. The largest version of each series is depicted in Fig. 1.

**Pre-training task.** As illustrated in Fig. 7, different probing tasks are employed to analyze the contextualized token embeddings of BERT, RoBERTa, and T5. Notably, the probing datasets used are peS2o for BERT and RoBERTa, and C4 for T5. Since masked language modeling (MLM) and span corruption (SC) mask or corrupt only 15% of tokens, leading to a corresponding reduction to 15% of the effective examples compared to NTP, we multiply the size of the probing datasets by 7 for MLM and SC.

**Information flow.** As illustrated in Fig. 8, the contextualized token embeddings of the current token $(x_t)$ at each layer are leveraged to predict various tokens in the sequence, spanning previous token $(x_{t-1})$ to next next token $(x_{t+2})$, including the default next token $(x_{t+1})$. For clarity, we present four models selected from those shown in Fig. 1: GPT2-XL, Llama-3-8B-Instruct, RWKV-Raven-14B, and Mamba-2.8B. For simplicity, we present an explicit formulation for predicting the previous token, $x_{t-1}$; analogous formulas for predicting other tokens can be obtained by appropriately adjusting the prediction target in $\text{PR}^{\text{prev}}$ and $\text{PR}^{\text{prev}}_{\ell}$. Specifically, the model uses the last-layer embedding of the current token, denoted as $\mathbf{h}_{t,\text{last}}$, to predict the *previous* token in the sequence. Let $x^s_{t-1}$ represent the immediately preceding token that the LLM aims to predict based on the first $t$ tokens, $x^s_1, x^s_2, \ldots, x^s_t$, through the last-layer embedding of the current token $\mathbf{h}^s_{t,\text{last}}$, for an index $1 \leq s \leq S$ over all training corpus. This process forms a dataset $\mathcal{D}_{\text{prev}} := \{(\mathbf{h}^s_{t,\text{last}}, x^s_{t-1}) \mid 1 \leq s \leq S\}$. To assess the capability of the LLM's current token embeddings in predicting the previous token, we evaluate how well a linear regression model fits on the dataset $\mathcal{D}_{\text{prev}}$. For this purpose, we identify $x$ with its index in the token vocabulary. Let $\hat{x}_{\text{prev}} = \mathbf{w} \cdot \mathbf{h} + b$ denote the least-squares fit on $\mathcal{D}_{\text{prev}}$. This suggests using the following metric to quantify the LLM's previous-token prediction capability:

$$\text{PR}^{\text{prev}} := \frac{\sum(x_{\text{prev}} - \hat{x}_{\text{prev}})^2}{\sum(x_{\text{prev}} - \bar{x}_{\text{prev}})^2}, \tag{2}$$

where the sum is over all $x_{\text{prev}} = x^s_{t-1}$, $\hat{x}_{\text{prev}} = \mathbf{w} \cdot \mathbf{h}^s_{t,\text{last}} + b$, and $\bar{x}_{\text{prev}}$ represents the mean of all $x^s_{t-1}$. To investigate how the predictive power of an LLM with depth $L$ evolves across its layers, we calculate the $\text{PR}^{\text{prev}}$ for the previous-token prediction task at each intermediate layer. Let $\text{PR}^{\text{prev}}_l$ denote this value for the $l$-th layer, where $1 \leq l \leq L$. Specifically, instead of using the last-layer embedding $\mathbf{h}_{t,\text{last}} \equiv \mathbf{h}_{t,L}$, we use the embedding of the current token at layer $\ell$, denoted $\mathbf{h}_{t,\ell}$, to predict the previous token when computing $\text{PR}^{\text{prev}}_{\ell}$.

## Additional Results

**Measure analysis.** Fig. S8 illustrates that the widely used metric of separation fuzziness, commonly applied to assess deep learning features in classification tasks [30, 10, 17], is inadequate for the emergence of the law of equi-learning in LLMs. This inadequacy may stem from the larger token vocabulary size compared to the embedding dimension and the presence of very similar or even identical contexts followed by different tokens in natural language data [48]. For simplicity,
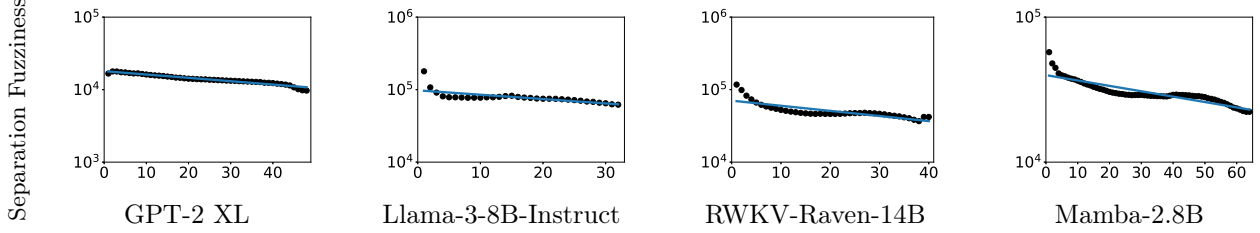
Fig. S8: With the measure of separation fuzziness, the law of equi-learning is not clear though the decreasing trend is still observed. The x-axis denotes the layer index, while the y-axis (log scale) shows the separation fuzziness.
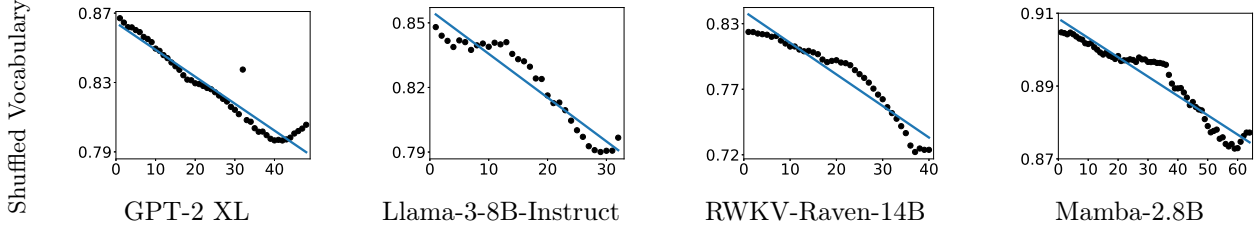


Fig. S9: When the vocabulary is shuffled, the law of equi-learning is not very clear though the decreasing trend is still observed. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2.

we selected four models from Fig. 1: GPT-2 XL, Llama-3-8B-Instruct, RWKV-Raven-14B, and Mamba-2.8B, and utilized separation fuzziness to evaluate the quality of contextualized token embeddings. Furthermore, as demonstrated in Fig. S9, the law of equi-learning becomes obscured when the vocabulary is shuffled, resulting in differing token index orders. This observation suggests that the widely adopted byte pair encoding algorithm [37] in tokenizers can produce a meaningful token index order, which is critical to the emergence of the law of equi-learning. Further investigation is needed and is deferred to future work.

**Layer normalization analysis.** As illustrated in Fig. S10, layer normalization plays a critical role in pre-LN models, as the absence of additional layer normalization makes the law noisier. Notably, this normalization effect can also be attained through standardization, specifically $z = \frac{x-\mu}{\sigma}$ across the embedding dimension. For simplicity, we selected four models from Fig. 1—GPT-2 XL, Llama-3-8B-Instruct, RWKV-Raven-14B, and Mamba-2.8B—and analyzed the impact of layer normalization. This finding aligns with observations related to the equi-separation law in MLPs [17], where batch normalization is crucial for its emergence. It is important to note that batch normalization does not impact the PR of contextualized token embeddings in our case.

**Probing data analysis.** Table S1 highlights that the probing datasets exhibiting the strongest Pearson correlations for various LLMs in Fig.1 are BookCorpus, C4, and peS2o among the eight evaluated datasets. This trend may be attributed to the high quality of these datasets and their resemblance to the pre-training data used for these models. Notably, the optimal probing dataset among the eight evaluated transitions from peS2o to BookCorpus as models progress from Llama-1 and Llama 2 (including its chat variant) to Llama 3 (and its instruct variant). This shift suggests that higher-quality pre-training data may require correspondingly higher-quality probing datasets to facilitate the emergence of the law of equi-learning, likely reflecting the impact of carefully
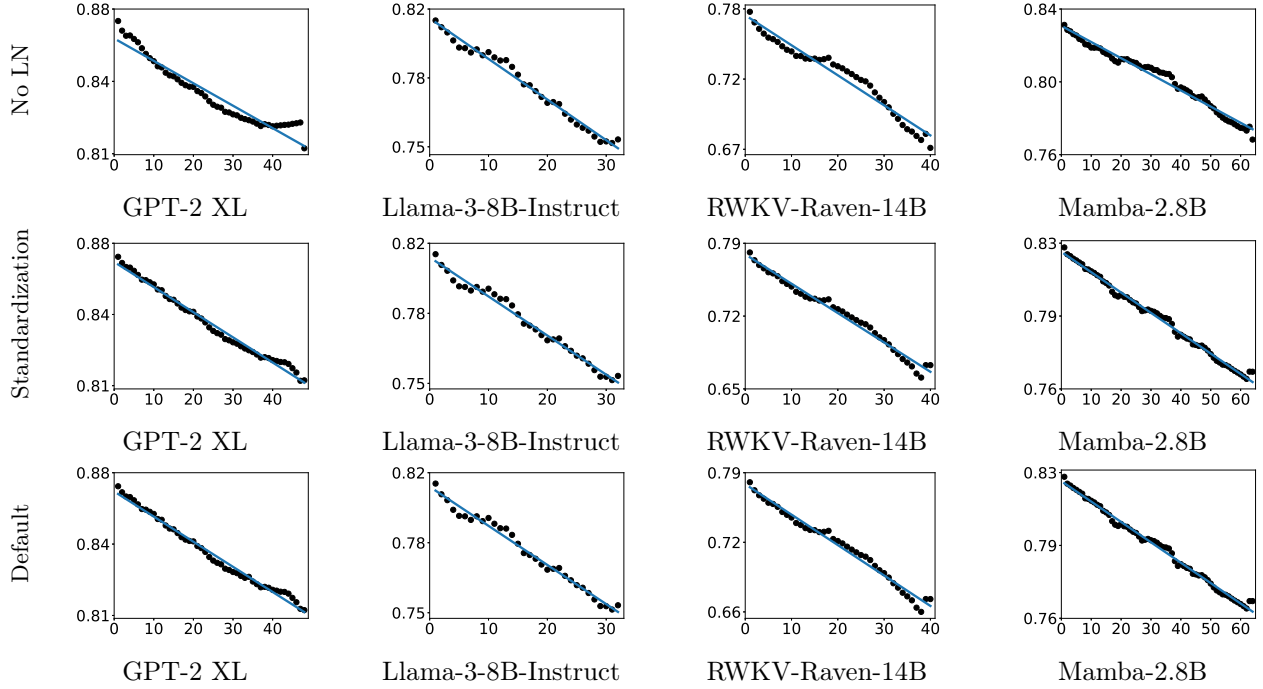
Fig. S10: For pre-LN models, omitting the additional initialized layer normalization (no LN) for contextualized token embeddings leads to a noisier manifestation of the law of equi-learning compared to the default approach with initialized layer normalization. This normalization effect can also be achieved through standardization, rather than relying on initialized layer normalization. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2.
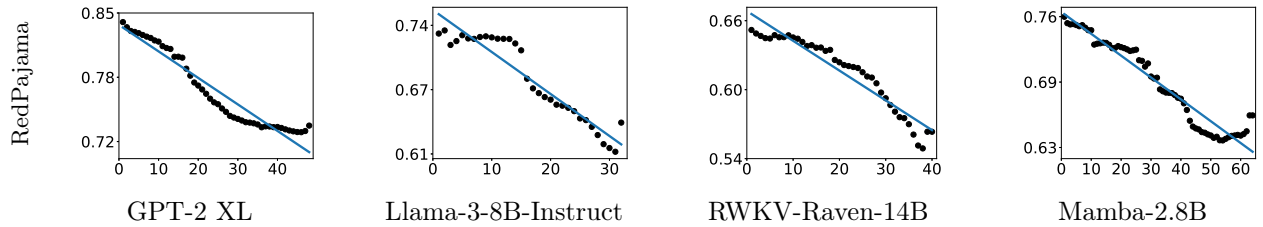


Fig. S11: With unsuitable probing dataset (e.g., RedPajama here), the law of equi-learning is not clear, though a descending trend is still observed. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2.

|  | BookCorpus | C4 | OpenWebText | Wikipedia | peS2o | Pile | Redpajama | OSCAR |
|---|---|---|---|---|---|---|---|---|
| GPT-1 | **-0.997** | -0.951 | -0.959 | -0.943 | -0.972 | -0.917 | -0.884 | -0.969 |
| GPT-2 XL | **-0.994** | -0.960 | -0.982 | -0.962 | -0.963 | -0.965 | -0.962 | -0.957 |
| Llama-1-13B | -0.956 | -0.993 | -0.990 | -0.986 | **-0.994** | -0.992 | -0.993 | -0.993 |
| Llama-2-13B | -0.913 | -0.985 | -0.977 | -0.966 | **-0.988** | -0.973 | -0.983 | -0.985 |
| Llama-2-13B-Chat | -0.879 | -0.984 | -0.967 | -0.962 | **-0.983** | -0.964 | -0.978 | -0.980 |
| Llama-3-8B | **-0.993** | -0.981 | -0.977 | -0.979 | -0.938 | -0.941 | -0.948 | -0.969 |
| Llama-3-8B-Instruct | **-0.992** | -0.981 | -0.974 | -0.979 | -0.940 | -0.940 | -0.948 | -0.967 |
| Mistral-7B-v0.1 | -0.874 | **-0.988** | -0.950 | -0.968 | -0.956 | -0.941 | -0.940 | -0.979 |
| Mistral-7B-Instruct-v0.1 | -0.850 | **-0.991** | -0.958 | -0.971 | -0.961 | -0.939 | -0.948 | -0.985 |
| Mistral-7B-v0.2 | -0.863 | **-0.989** | -0.952 | -0.967 | -0.953 | -0.936 | -0.941 | -0.979 |
| Mistral-7B-Instruct-v0.2 | -0.874 | **-0.988** | -0.952 | -0.967 | -0.958 | -0.931 | -0.936 | -0.979 |
| Mistral-7B-v0.3 | -0.863 | **-0.989** | -0.952 | -0.967 | -0.953 | -0.936 | -0.941 | -0.979 |
| Mistral-7B-Instruct-v0.3 | -0.863 | **-0.988** | -0.951 | -0.966 | -0.951 | -0.932 | -0.938 | -0.979 |
| phi-1.5 | **-0.994** | -0.967 | -0.986 | -0.976 | -0.974 | -0.958 | -0.967 | -0.971 |
| phi-2 | **-0.993** | -0.984 | -0.978 | -0.983 | -0.890 | -0.954 | -0.930 | -0.987 |
| phi-3-medium-4k-instruct | -0.894 | **-0.992** | -0.961 | -0.955 | -0.977 | -0.972 | -0.975 | -0.991 |
| phi-3-medium-128k-instruct | -0.902 | **-0.992** | -0.962 | -0.959 | -0.975 | -0.972 | -0.974 | -0.991 |
| RWKV-14B | -0.984 | **-0.991** | -0.983 | -0.987 | -0.967 | -0.984 | -0.952 | -0.983 |
| RWKV-Raven-14B | -0.984 | **-0.991** | -0.982 | -0.984 | -0.964 | -0.984 | -0.952 | -0.981 |
| Mamba-2.8B | **-0.997** | -0.981 | -0.989 | -0.987 | -0.992 | -0.981 | -0.966 | -0.986 |

Table S1: Pearson correlation coefficients between the logarithm of PR values (Eq. 2) and layer indices for all LLMs in Fig. 1, evaluated across eight probing datasets. For each LLM, the probing dataset with the strongest Pearson correlation is highlighted in bold.

designed pre-processing and curation pipelines used to enhance the pre-training data quality in Llama 3 [9]. Furthermore, as demonstrated in Fig. S11, the observed law's clarity diminishes when an inappropriate probing dataset is selected, although a general descending trend remains evident. For simplicity, we selected four models from Fig. 1—GPT-2 XL, Llama-3-8B-Instruct, RWKV-Raven-14B, and Mamba-2.8B—and presented the PR of their token embeddings using the RedPajama dataset as the probing dataset. These findings underscore the critical importance of selecting appropriate probing data for the emergence of the law of equi-learning.

**Pre-training data analysis.** As illustrated in Fig. S12, GPT-2 models pre-trained on two different datasets, C4 and OSCAR, exhibit distinct behaviors regarding the emergence of the law of equi-learning when evaluated on the same probing dataset (i.e., OpenWebText). Specifically, models pre-trained with OSCAR demonstrate more noise in the law's emergence compared to those pre-trained with C4. This discrepancy is likely attributable to the higher noise levels in OSCAR, stemming from its less stringent deduplication. These models were released by [28]. Our findings underscore the significant impact of pre-training data quality on the emergence of the law of equi-learning, suggesting that higher quality pre-training data may result in a more pronounced manifestation of this law.
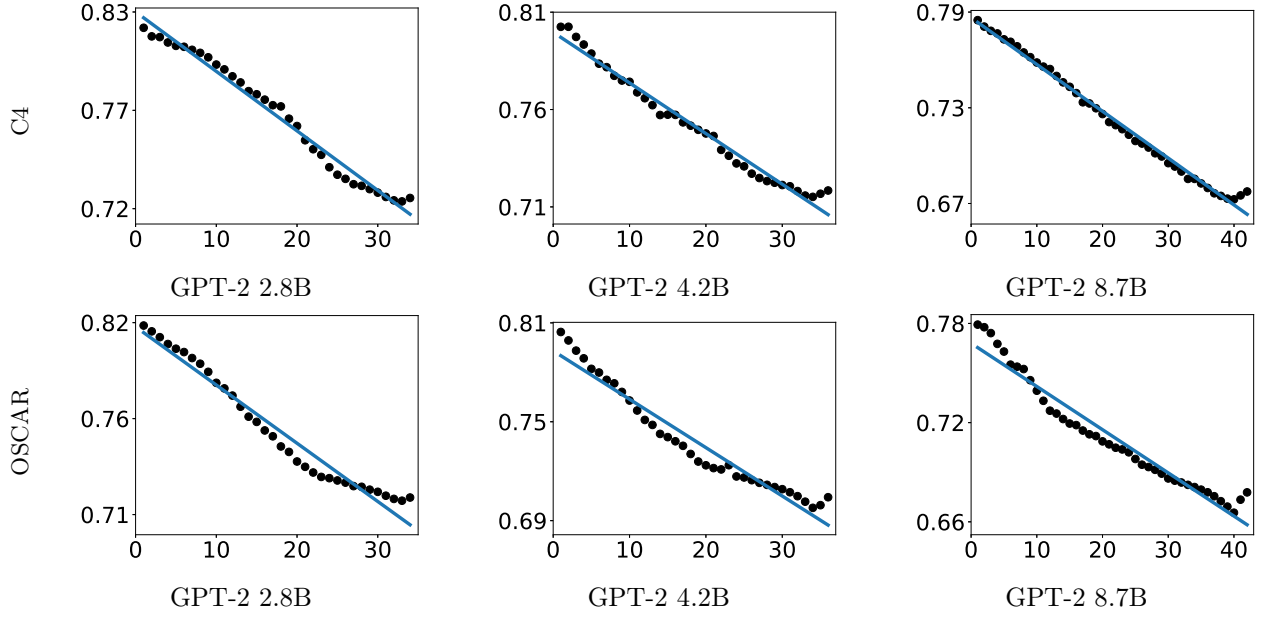
Fig. S12: Pre-training data can affect how the law of equi-learning behaves. The same GPT-2 models pre-trained on C4 and OSCAR are probed with the OpenWebText dataset. The x-axis denotes the layer index, while the y-axis (log scale) shows the prediction residual (PR) as defined in Eq. 2.