# THE GEOMETRY OF TOKENS IN INTERNAL REPRESENTATIONS OF LARGE LANGUAGE MODELS

**Karthik Viswanathan**[1,2*], **Yuri Gardinazzi**[2,3] , **Giada Panerai**[2],

**Alberto Cazzaniga**[2] and **Matteo Biagetti**[2*]

[1]University of Amsterdam, Amsterdam, the Netherlands
[2]Area Science Park, Trieste, Italy
[3]University of Trieste, Trieste, Italy

January 22, 2025

## ABSTRACT

We investigate the relationship between the geometry of token embeddings and their role in the next token prediction within transformer models. An important aspect of this connection uses the notion of empirical measure, which encodes the distribution of token point clouds across transformer layers and drives the evolution of token representations in the mean-field interacting picture. We use metrics such as intrinsic dimension, neighborhood overlap, and cosine similarity to observationally probe these empirical measures across layers. To validate our approach, we compare these metrics to a dataset where the tokens are shuffled, which disrupts the syntactic and semantic structure. Our findings reveal a correlation between the geometric properties of token embeddings and the cross-entropy loss of next token predictions, implying that prompts with higher loss values have tokens represented in higher-dimensional spaces.

## 1 Introduction

In the context of interpretability of transformer models, a set of analytic approaches have been developed with the goal of modeling transformer architectures as dynamical systems of particles [1, 2, 3]. In this perspective, the transformers are viewed as evolving a mean-field interacting particle system where the evolution of tokens across layers is controlled by their empirical measure[2] [4]. Under a set of assumptions such as time-independent weights, this interpretation is used to show that tokens tend to cluster in the later layers [5]. This clustering behavior can be associated with the empirically observed rank collapse phenomenon in transformer models [6, 7, 8, 9, 10].

An important insight from [2] in the context of next token prediction is that the output measure of tokens encodes the probability distribution of the next token, and its clustering indicates a small number of possible outcomes. A complementary perspective to the evolution of token representations across layers can be gained by studying the latent predictions of transformer models [11] from the perspective of iterative inference [12] which indicates that the probabilities of the next tokens are incrementally updated layer by layer. The work by [13] suggests that causal LLMs appear to develop a reasonably accurate prediction regarding the next token in the middle layers, with subsequent layers refining these predictions. This means we should expect the empirical measures of the internal layers to reflect this trend, i.e. a rapid change of the empirical measure in the early layers and a more refined change towards the later layers. Since the latent predictions are obtained by unembedding the residual stream [14], and our method probes

---

[*]Correspondence: k.viswanathan@uva.nl, matteo.biagetti@areasciencepark.it
[2]In this context, the empirical measure and the output measure are used to characterize the distribution of tokens in the internal layers and the output layer respetively

the geometric properties of the residual stream, we can expect the statistical properties (e.g. entropy) of the latent prediction probabilities to be encoded in the geometry of the internal representations of the tokens.

In this work, we combine these viewpoints to examine the empirical measure of the internal layers from a geometric perspective. To observationally probe the empirical measure, we draw inspiration from previous works using intrinsic dimension and neighbourhood overlap to study the geometry of internal representations [15, 16, 17, 18, 19, 20, 21, 22]. In these works, an important difference is that point clouds are built as a collection of prompts represented as a single point (the last token), thereby lacking a direct link to the empirical measure of tokens within a prompt. Additionally, we also calculate cosine similarity as a general probe of pairwise relations among tokens.

To test how the geometric properties of token representations change as a function of the model's internal dynamics, we probe it in a regime where the syntactic and semantic structures of the prompts are disrupted through systematic token shuffling. Our analysis achieves these main results:

- **Token-Level Intrinsic Dimension and Cosine Similarity**: Section 4. We observe that the intrinsic dimension (ID) of token representations generally exhibits a peak, whose height increases with the degree of token shuffling. This peak is located at early to middle layers of the models. On the other hand, cosine similarity among tokens increases with shuffling, indicating increased alignment of token vectors.

- **Neighborhood Overlap Consistency**: Section 4. The neighborhood overlap (NO) metric shows that token relationships around the ID peak become less consistent as the amount of shuffling increases. This highlights that structured data retains more coherent token neighborhoods through the model layers than shuffled data.

- **Correlation with Model Loss**: Section 5. We find a statistical relation between the geometry of tokens and the probability distribution of the next token: the intrinsic dimension of the token representations across hidden layers is correlated to the average cross-entropy loss of the next token probability distribution for a given prompt. This suggests that prompts with a higher cross-entropy loss have token representations lying in higher dimensional manifolds.

## 2 Related work

**Lenses for Mechanistic Interpretability in Transformers.** Mechanistic interpretability in transformers explores how transformer models encode and utilize information, focusing on semantic and algorithmic interpretations. Semantic interpretation investigates what latent properties are learned by models and how individual neurons may code for specific concepts [23]. Structural probing [24, 25, 26] and dictionary learning [27, 28, 29] offer insights into how features are represented and reconstructed in transformer architectures. Relevant to this work, is the approach of the logit lens [13]. This method offers insight into a model's predictive process by applying the final classification layer, which converts the residual stream activation into logits/vocabulary space, to intermediate residual stream activations. This reveals how prediction confidence evolves throughout the computational stages. This is feasible because transformers typically construct their predictions iteratively across layers [30]. Building on this concept, the tuned lens [11] employs affine probes to translate internal representations into probability distributions over the vocabulary. Similarly, the Future Lens [31] examines how individual representations encode information about forthcoming tokens.

**Analytic Approaches to Transformer Models.** Recent analytical works [3, 5] indicate that analyzing geometric properties of token representations and their dynamics can offer meaningful insights into how transformers function. [5] introduced the novel perspective of viewing the evolution of tokens in the transformer layers as particles in a dynamical system. They predict clustering behavior in transformer models in a simplified setting which was later extended to include causally masked attention [32]. [3] adopts the above perspective and examines particle geometry in the presence of MLP layers. This perspective not only offers insights into the geometric dynamics of tokens but also addresses the trainability of transformers based on initialization hyperparameters, including the strength of attentional and MLP residual connections.

Further studies [33, 34] theoretically investigate the expressive power of transformers as maps from arbitrary input measures to output measures and prove the appearance of dynamic metastability, i.e. the particles cluster in the infinite time limit but they resemble a configuration of several clusters for a long period of time. This behavior aligns more closely with practical observations than the clustering dynamics. This analytical framework highlights the significance of studying the distribution of the internal representations of the tokens (referred to as the *empirical measure*) by i) suggesting a relation between the empirical measure to the next token prediction loss [2] ii) understanding the role of the empirical measure in governing the token dynamics [4].

**Geometric Approaches to Transformer Models.** The manifold hypothesis posits that real-world high-dimensional data often lie on or near a lower-dimensional manifold within the high-dimensional space [35]. The dimension of this approximating manifold is usually named the *intrinsic dimension* of the data. Several studies have demonstrated that

the intrinsic dimension of data representations in deep networks shows a remarkable dynamic range, characterized by distinct phases of expansion and contraction [15, 16, 17]. Data manifolds created by internal representations in deep networks have been also explored from the perspective of neuroscience and statistical mechanics [36, 37]. In LLMs, a geometric analysis of representations has uncovered a rich set of phenomena. Geometric properties, such as intrinsic dimension and the composition of nearest neighbors, evolve throughout the network's sequence of internal layers. These changes mark distinct phases in the model's operation, signal the localization of semantic information [18, 19]. [38] analyze the intrinsic dimension by considering all tokens to reveal semantic correlations in images and text inside deep neural networks. While the aforementioned works analyze internal representations in linguistic processing, the geometry of context embeddings has been linked to language statistics [39] and used to highlight differences between real and artificial data [40].

## 3 Method

Transformer models take as input a sequence of vectors embedded in $d$-dimensions of varying length $N$, $\{x_i\}_{i \in [N]} \in \mathbb{R}^{d \times N}$. Each element of the sequence is called a *token*, while the entire sequence is a *prompt*. A transformer is then a sequence of maps:

$$\{x_i(1)\}_{i \in [N]} \rightarrow \{x_i(2)\}_{i \in [N]} \cdots \rightarrow \{x_i(N_{\text{layers}})\}_{i \in [N]}, \tag{1}$$

where $x_i(\ell) \in \mathbb{R}^{d \times N}$ represents the $i$-th token at layer $\ell$, $N_{\text{layers}}$ the total number of model layers and $N$ is the number of tokens.

In transformer models, prompts can vary based on the specific application, representing protein sequences, image pixels, or text sentences. In this study, we focus on causal language models and use sentences as our input prompts, though the technique can be extended to other input types as well. The prompt size can significantly vary depending on the dataset considered: sentences can be $\mathcal{O}(10)$ - $\mathcal{O}(1000)$ tokens long. Given that our goal is to study and interpret the geometrical behavior at the token level across model layers, we select prompts with a sufficient number of tokens, i.e. $N \geq 1024$ tokens, to ensure reliable estimates of our observables.

**Empirical measure.** Given $n$ points at positions $x_1, \ldots, x_n \in \mathbb{R}^d$ (a point cloud), their empirical measure is the probability measure $\mu = \frac{1}{n} \sum_{j=1}^{n} \delta_{x_j}$, i.e., the empirical measure encodes the distribution of points in the embedding space. In the context of transformers [2], the empirical measure characterizes the distribution of the tokens at each layer of the sequence 1. The empirical measure for the last layer is the *output* measure. The dynamical evolution of tokens in this framework, as described by Equation (1) in [4], indicates that the change in the token representation of token $i$ is controlled by a layer-dependent kernel $K_\ell$ and depends purely on the current token representation $x_i(\ell)$ and the empirical measure[3]. To probe the empirical measure across layers, we use cosine similarity, intrinsic dimension, and neighborhood overlap, as defined below.

**Intrinsic Dimension.** A substantial body of literature focuses on developing precise estimators for the intrinsic dimension of manifolds [41]. In particular, nearest-neighbors-based algorithms are robust to high dimensionality and capture the non-linear structure of the manifold. In addition, it has been argued that a scale-sensitive algorithm can provide a stable estimation of the dimension, as it allows us to find the range of scale where the dimension is constant.

GRIDE [42] is a likelihood-based ID estimator that estimates the intrinsic dimension $\hat{d}(n_1, n_2)$ using the ratios $\dot{\mu} = \mu_{i,n_1,n_2} = \frac{r_{i,n_2}}{r_{i,n_1}}$, where $r_{i,k}$ is the Euclidean distance between point $i$ and its $k$-th nearest neighbour and $1 \leq n_1 < n_2$. Under the assumption of local uniform density, the distribution of $\mu_{i,n_1,n_2}$ is given by,

$$f_{\mu_{i,n_1,n_2}}(\dot{\mu}, d) = \frac{d\left(\dot{\mu}^d - 1\right)^{n_2 - n_1 - 1}}{\dot{\mu}^{(n_2 - 1)d + 1} B\left(n_2 - n_1, n_1\right)}, \quad \dot{\mu} > 1 \tag{2}$$

where $B(\cdot, \cdot)$ is the beta function. The ID estimate $\hat{d}(n_1, n_2)$ is obtained by maximizing the above likelihood with respect to $d$ assuming that the ratios $\mu_{i,n_1,n_2}$ are independent for different points. The conventional choice for the GRIDE algorithm is to set $n_2 = 2n_1$ and examine the variation of $\hat{d}$ for $n_2 \in \{2, 4, 8..\}$, where the parameter $n_2$ is known as the range scaling parameter. In this work, we mainly work with range scaling = 2 (unless explicitly mentioned), which is related to the TWO-NN estimator [41]:

$$\hat{d}_{\text{TWO-NN}} = \frac{N - 1}{\sum_i^N \log\left(\mu_{i,1,2}\right)}, \tag{3}$$

---

[3]the dynamics of a token $i$ depends on the position of all the tokens $x_j(\ell)$ but not on their labels, which is an assumption in the mean-field interacting particle framework.

**Algorithm 1** Shuffling algorithm

**Require:** $tokens, S$        ▷ $S$ is the shuffle index
**Output:** $permutedTokens$
    $nBlocks \leftarrow 4^S$
    $n \leftarrow tokens.length()$
    $B \leftarrow \lceil n/nBlocks \rceil$        ▷ $B$ is the block size
    $blocks \leftarrow \mathrm{splitInBlocks}(tokens, nBlocks, B)$        ▷ Split list into $nBlocks$ sublists of size $B$
    $permutedBlocks \leftarrow \mathrm{randomPermutation}(blocks)$
    $permutedTokens \leftarrow \mathrm{mergeBlocks}(permutedBlocks)$

| This | paper | is | titled | : | the | geometry | of | tokens | in | internal | representations | of | large | language | models | S=0 |

| This | paper | is | titled | : | the | geometry | of | tokens | in | internal | representations | of | large | language | models | S=1 |

↓ shuffling

| : | the | geometry | of | of | large | language | models | tokens | in | internal | representations | This | paper | is | titled |

| This | paper | is | titled | : | the | geometry | of | tokens | in | internal | representations | of | large | language | models | S=2 |

↓ shuffling

| geometry | tokens | language | large | of | in | representations | of | This | paper | : | is | models | titled | internal | the |

Figure 1: **The shuffling algorithm with an example.** Top Panel: Algorithmic description of the shuffling procedure described in Section 4. Bottom Panel: An example of the shuffling algorithm using $N = 16$ tokens. The first row ($S = 0$) corresponds to the unshuffled sequence. When $S = 1$, the tokens are split into $4^1$ blocks first and then, the blocks are shuffled. The last row $S = 2$ shows the fully shuffled case where the tokens are randomly permuted.

The above equation relates the intrinsic dimension to the generic ratios $\mu_{i,1,2}$ thereby implying that there is an inverse relation between the dimension estimate and the generic ratios ($\mu_{i,1,2}$). Equation 3 indicates that a higher dimensional estimate implies a lower ($\mu_{i,1,2}$) on average. [4]

**Neighborhood Overlap.** The neighborhood overlap $\chi_k^{l,m}$ was introduced in [43] to measure similarity between representations in different layers $\ell, m$ at a given scale $k$. Given the representations of $N$ tokens in layers $\ell$ and $m$, we can define $\chi_k^{\ell,m}$ as

$$\chi_k^{\ell,m} = \frac{1}{N} \sum_i \frac{1}{k} \sum_{j \in \mathcal{N}_k^\ell(i)} \mathbb{I}\left(j \in \mathcal{N}_k^m(i)\right) \tag{4}$$

where $\mathcal{N}_k^\ell(i)$ is the set of $k$-nearest neighbors of a token $i$ in layer $\ell$. Intuitively, it measures the average number of shared $k$-nearest neighbors in layers $\ell, m$. In our context, we set $m = \ell + 1$, i.e. we calculate the neighborhood overlap between adjacent layers. By doing so, we measure the change in pairwise relations among tokens between successive layers.

## 3.1 Models and Datasets

**Models.** In this work, we analyze 3 different pre-trained decoder-only LLMs: Llama 3 8B [44], Mistral 7B [45], Pythia 6.9B[46], each of them having 32 hidden layers and a hidden dimension of 4096. For brevity, we call them LLAMA, MISTRAL, and PYTHIA from now on. In the plots, layer 0 represents the embedding layer, with the hidden layers starting from layer 1. We extract internal representations from these models using the HuggingFace Transformers

---

[4]Note that the main assumption of the estimator is local homogeneity (Poisson distributed points within the local neighborhood of the point), which is generally true on a wide range of datasets.

library[5]. The token representations, stored in the `hidden_state` variable, corresponds to the representations in the residual stream [14] after one attention and one MLP update. In the models considered for analysis, layer normalization is applied before self-attention and MLP sublayers. LLAMA and MISTRAL add the self-attention outputs to the residual stream before the MLP whereas PYTHIA adds the self-attention and MLP sublayer outputs to the residual stream in parallel.

**Datasets.** As a dataset representative of text in an extensive way, we use the Pile dataset which comprehends text from 22 different sources [47]. For computational reasons, we opted for the reduced size version Pile-10K [48]. We further filter only prompts of sequence length $N \geq 1024$ according to the tokenization schemes of all the above models. This choice ensures a reliable ID estimate. This results in 2244 prompts after filtering. We truncate the prompts by keeping the first $N = 1024$ tokens to eliminate the length-induced bias of our ID estimates if it were to be present.

## 4 The geometry of shuffled and unshuffled prompts

Evaluating geometric observables at the token level directly probes the model's internal dynamics. As a way of quantifying geometric changes, we compare in-distribution data to various levels of token shuffling. By progressively disrupting the syntactic and semantic structure while preserving unigram frequency distribution, we observe the incremental effects on our observables across layers.

**Shuffling method.** We define the shuffling of tokens in the following way: given a prompt with $N$ tokens, $X = \{x_i\}_{i \in [N]}$, we split the sequence into $nBlocks$ blocks of size $B$ such that $nBlocks \times B = N$ and take one random permutation of the blocks, as schematically presented in Figure 1. In our experiments, we choose $nBlocks = 4^S$, where $S$ is the shuffle index where the shuffle index $S$ quantifies the degree of shuffling. Note that the $S = 0$ represents the unshuffled state and the shuffle index for the fully shuffled case ($\hat{S}$) corresponds to the value of $S$ when the number of tokens $N = 4^{\hat{S}}$. In Figure 1, we have $\hat{S} = 2$ since we consider 16 tokens, whereas in the experiments, we have $\hat{S} = 5$ because we have $1024 = 4^5$ tokens.

We show two main results: i) the effect of various degrees of shuffling on our metrics for a single, randomly extracted prompt and ii) the qualitative behavior of the unshuffled and the fully shuffled prompts on average. For the former observable, we consider the $3218^{\text{th}}$ prompt from the Pile-10K dataset, with the Pile set name: *ArXiv*. This prompt is shuffled to six different levels labeled by $(S = 0, 1, \ldots, 5)$. We study the representations of this prompt in LLAMA[6]. To understand the qualitative differences between the shuffled and unshuffled prompts across the dataset, we find the averages of the geometric quantities (cosine similarity, ID and NO) over 2244 prompts.

### 4.1 Cosine Similarity

As a first step into investigating the geometry of internal representations at the token level, we compute the cosine similarity among tokens for each layer. In Figure 2, we show the average cosine similarity for different levels of shuffling as a function of model layers on a single prompt (Left Panel) and for the average over all prompts (Right Panel) for the LLAMA model. We can see that the cosine similarity increases with increasing shuffling and increasing layers. This implies that tokens are distributed along the same direction towards the last layers. For the structured prompts, the average cosine similarity is closer to zero, indicating that their directions are more orthogonal.

These results seem related to earlier works: [49] computes the average cosine similarity of randomly sampled words from BERT, GPT-2 and ELMo models across layers, finding high cosine similarity, in agreement with our shuffled case. In [50] average cosine similarity was also computed on pre-trained text transformers, finding an average value of $\approx 0.5$ in the last layer.

### 4.2 Intrinsic Dimension

Next, we examine the intrinsic dimension profile of tokens as a function of layers. Figure 3 displays the ID calculated for a range scaling of 2 for LLAMA. The Left Panel shows the ID profile of a single prompt at various levels of shuffling, while the Right Panel presents the average ID across 2244 prompts for both fully shuffled and structured cases. In all scenarios, we observe a peak in ID in the early to middle layers. Additionally, the height of this peak increases with the degree of shuffling, indicating a relation between the two. Previous work focusing on studies of the

---

[5]The library is available at https://huggingface.co/docs/transformers/en/index

[6]The qualitative behavior discussed in this section holds in general for other prompts and models. We show this in the case of intrinsic dimension by looking at the ID profile of other prompts using LLAMA (Figure 8) and the ID profile of $3218^{\text{th}}$ prompt in other models (Figure 9)
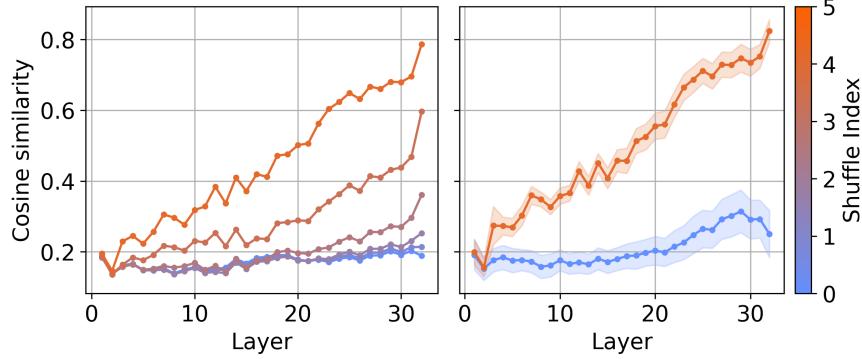
Figure 2: **Average Cosine Similarity.** Left Panel: average cosine similarity among tokens for a single prompt as a function of model layers. Right Panel: average cosine similarity averaged over 2244 prompts as a function of layers for the full shuffle ($S = 5$) and the structured case ($S = 0$). The color bar indicates the shuffle index $S$. The shaded regions indicate the standard deviation from the mean. All curves have been calculated for the LLAMA model.

geometry of internal representations at the prompt-level have investigated similar metrics. In [20] it was found that the ID calculated from the last token representations of shuffled prompts have a lower intrinsic dimension than the unshuffled prompts, which is different than the behavior observed at the token-level. We explore this difference in the effect of shuffling on the prompt and token level ID profiles in more detail in Section D.



Figure 3: **Intrinsic Dimension.** Left Panel: intrinsic dimension for a single random prompt as a function of model layers. Right Panel: intrinsic dimension averaged over 2244 prompts as a function of layers for the full shuffle ($S = 5$) and the structured case ($S = 0$). The shaded regions indicate the standard deviation from the mean. The color bar indicates the shuffle index $S$. All curves have been calculated for the LLAMA model.

**Distribution of tokens at the ID peak.**

We consider the relation in equation 3 between ID and the generic ratios $r_{i,2}/r_{i,1}$, i.e. the ratio of the distance of the second neighbor over the one of the first neighbor to the $i$-th point. According to equation 3, as ID grows we expect the ratio to tend to unity on average, implying that the first two nearest neighbors are roughly at equal distance from the reference token. On the other hand, if ID decreases we expect the two nearest neighbors to be at more varying distances. Therefore, a higher ID at the peak means that the nearest neighbors tend to be more equidistant for the shuffled prompts.

Additionally, we examine the angular distribution of nearest neighbors [51], as it offers a complementary perspective to the previous discussion, which is exclusively based on the distances to the nearest neighbors. Hence, we compute the cosine similarity between $x_{i,1} - x_i$ and $x_{i,2} - x_i$ for each token $i$ to determine the distribution of the angle formed by the first two nearest neighbors centered at token $i$.[7] We visualize this in Figure 4. In the left panel, we show the histogram of the angles between $x_{i,1} - x_i$ and $x_{i,2} - x_i$ for each token $i$ of a random prompt, at layer 10 of LLAMA,

---

[7]In this paragraph, $x_{i,k}$ denotes the $k^{\text{th}}$ nearest token to token $i$ and $r_{i,k}$ is the distance between token $i$ and its $k^{\text{th}}$ nearest token.

i.e. around the ID peak. In the right panel, we show the histogram of means over 2244 prompts for the full shuffle and structured cases. The distributions of mean angles differ between the two cases, with the mean angle between the nearest tokens being closer to 60 degrees for shuffled prompts. Combined with the earlier observation that the ratio $r_{i,2}/r_{i,1}$ is closer to unity, this suggests that the triangle formed by $x_i$, $x_{i,1}$ and $x_{i,2}$ is more equilateral in the full shuffle case at the ID peak. These findings suggest a distinguishable arrangement of tokens for shuffled prompts that deserve further investigation in future work.



Figure 4: **Angle distribution between nearest neighbors.** Left Panel: histogram of the angles between the first and second nearest neighbor at layer 10 of the LLAMA model for a single prompt for the full shuffle case and structured case. The dotted vertical lines indicate the average angle between the nearest neighbors in both cases. Right Panel: histogram of the average angle between the first and second nearest neighbor at layer 10 of the LLAMA model in the fully shuffled (orange) and structured case (blue). The histograms are computed from 2244 prompts in each case.

### 4.3 Neighborhood overlap

We compute the neighborhood overlap at $k_{NN} = 2$ as a function of layers for the LLAMA model. We choose $k_{NN} = 2$ because we would like to examine a similar range of scales with ID computed using GRIDE at range scaling = 2. As a consistency check, we also calculate NO from $k_{NN} = 1$ to $k_{NN} = 6$ finding similar results (see Figure 11 in Appendix A). In Figure 5, we show a random prompt for different levels of shuffling (left panel) and the average over all prompts for the full shuffle and the structured case (right panel). The NO of the shuffled cases is lower than structured case around the layers corresponding to the ID peak, while being statistically similar away from the peak.

## 5 Intrinsic dimension is correlated with the model's loss

In the previous section, we probed the empirical measure across layers through geometric quantities like intrinsic dimension, neighborhood overlap, and cosine similarity. This section analyzes model behavior connecting our observations with next-token predictions. Specifically, we examine the correlation between the intrinsic dimension and the average cross-entropy loss of the next token predictions.

Given a prompt $X$ consisting of tokens $(x_1, x_2, ...x_N)$ and the model's next token prediction $p_\theta$ over a vocabulary $\mathcal{V}$, the average cross-entropy loss[8] is

$$\text{loss}(X) = -\frac{1}{N} \sum_i^N \log p_\theta (x_i \mid x_{<i}) \qquad (5)$$

where $\log p_\theta (x_i \mid x_{<i})$ is the log-likelihood of the $i^{\text{th}}$ token conditioned on the preceding tokens $(x_{<i})$.

To quantitatively examine the correlation between the intrinsic dimension and the average cross-entropy, we use the Pearson correlation coefficient $(\rho)$, defined as the ratio between the covariance of two variables and the product of their standard deviations. We compute the Pearson correlation between the average cross-entropy loss and the logarithm of ID across layers for the population of 2244 prompts across different models and show the result in Figure 6.

---

[8]This quantity is referred to by various names in the literature, including average surprisal, log perplexity, and average next-token prediction error, among others.
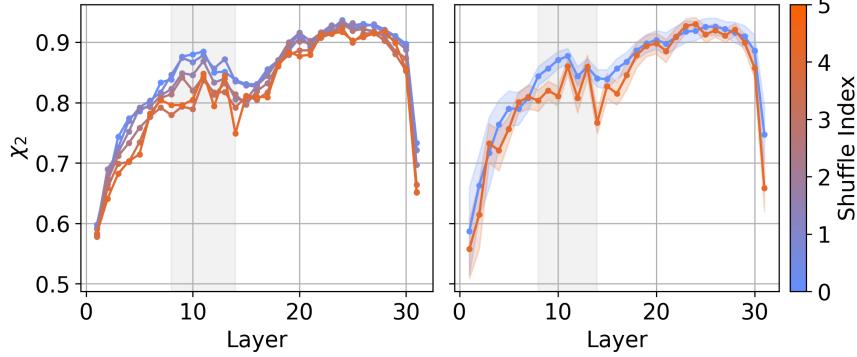
Figure 5: **Neighborhood Overlap.** Left Panel: neighborhood overlap for a single random prompt as a function of model layers for $k_{\mathrm{NN}} = 2$. The colorbar indicates the shuffle index $S$. Right Panel: neighborhood overlap averaged over 2244 prompts as a function of layers for the full shuffle ($S = 5$) and the structured case ($S = 0$). The shaded regions indicate the standard deviation from the mean and the grey region indicates the region around the ID peak when the shuffled prompts have a lower NO than the structured prompts. All curves have been calculated for the LLAMA model.



Figure 6: **Correlation between intrinsic dimension and the average cross-entropy loss.** Pearson coefficient between the logarithm of the intrinsic dimension and model loss for different models as a function of layers. The shaded regions indicate the standard deviation from the mean. The three curves correspond to LLAMA (orange), MISTRAL (magenta), and PYTHIA (blue). The $p$-values for the Pearson coefficients in this plot are below $0.01$ except for the last layer in PYTHIA.

All three models have a high correlation, particularly around the ID peak. The connection between the cross-entropy loss and ID was discussed in [19] where the correlation was calculated between the peak ID of the dataset of the last token representations and the log of dataset perplexity in Fig. 2 of [19]. However, we get a correlation in a similar spirit at a finer level since it reveals a correlation at the level of individual prompts (more details on the comparison in Sec. D).

## 5.1 Understanding the correlation between intrinsic dimension and average cross-entropy loss

In Figure 6, we observe a correlation between a geometric quantity (the ID of internal representations) and an information-theoretic quantity (the cross-entropy loss). This transition from geometry to information-theoretic perspective occurs at the softmax layer between the last layer representations and the next token predictions. In this section, we understand the relationship between the ID at the last layer and the cross-entropy loss in more detail. The positive correlation between these two quantities can be explained through the following steps:

1. **Unembedding Tokens to Logits**: We expect the ID of the last layer to be strongly correlated to the ID of the logits since the unembedding is a linear transformation. This is confirmed by the Pearson coefficient of $\rho = 0.96$ between the $\log ID$ of the last layer and the logits.

2. **Logits to Contextual Entropy**: With this step, we relate the geometric perspective to the information-theoretic perspective. Typically, a softmax layer converts the logits to next token prediction probabilities, $p_\theta (v \mid x_{<i})$. From this, one can define the contextual entropy

$$H (x_{<i}) = - \sum_{v \in \mathcal{V}} p_\theta (v \mid x_{<i}) \log p_\theta (v \mid x_{<i}) = \mathbb{E}_{v \sim p(\cdot \mid \boldsymbol{x}_{<i})} [- \log p_\theta (v \mid x_{<i})] \tag{6}$$

where $(x_{<i})$ is the *context*. We can average this quantity over all the tokens in a prompt to obtain the *average contextual entropy* $\mathcal{H}(X)$:

$$\mathcal{H}(X) = \frac{1}{N} \sum_{i=1}^{N} H (x_{<i}) = - \frac{1}{N} \sum_{i=1}^{N} \sum_{v \in \mathcal{V}} p_\theta (v \mid x_{<i}) \log p_\theta (v \mid x_{<i}) \tag{7}$$

We empirically show a correlation between the logarithm of logits ID and the contextual entropy in the LLAMA model by observing a Pearson correlation of $\rho = 0.43$, as shown in the left panel of Figure 7.

3. **Contextual Entropy $\sim$ Cross-Entropy Loss**: Equation 6 shows that the contextual entropy is the expected value of the cross-entropy loss, with the expectation computed using the next token probabilities $p_\theta$. When we consider a large number of tokens in the prompts, we expect the contextual entropy to be almost equal to the cross-entropy loss of the next token predictions when averaged over all the tokens. This can be seen empirically in the right panel of Figure 7.
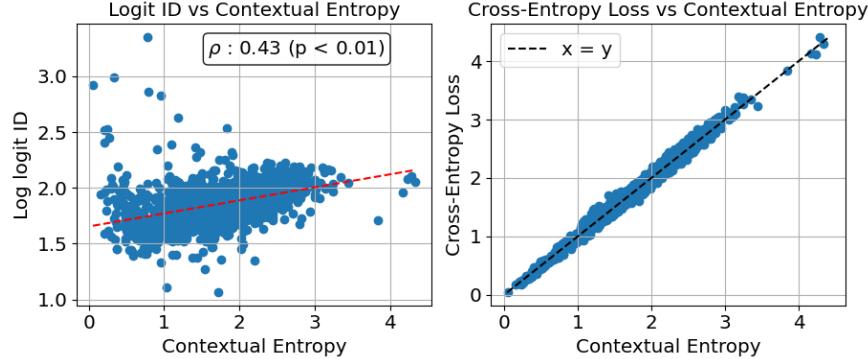


Figure 7: **Correlating intrinsic dimension at the last layer to cross-entropy loss.** The points in the following plots are calculated using the 2244 prompts considered in this paper for the LLAMA model. - (a) Left Panel: analysis of the correlation between the logits ID at scaling = 2 (refer to Figure 18 in the Appendix for scaling = 4, 8) and the contextual entropy to the average contextual entropy and (b) Right Panel: comparing the average contextual entropy to the average cross-entropy loss.

Given that the second step of the relation above is non-trivial, we analyze how the softmax layer connects the geometric properties of logits to the statistical properties of the resulting next-token prediction probabilities more in-depth. We wish to demonstrate that the correlation between logit IDs and the contextual entropy (entropy of next token predictions) is a fundamental property of the softmax layer, making this relationship more general.

Given $\mathbf{z} = (z_1, z_2, ... z_{|\mathcal{V}|})$ where $z_\alpha \in \mathbb{R}$, as the input to a softmax layer, the associated entropy of the probability distribution generated by the softmax operation is

$$S(\mathbf{z}) = - \sum_{\alpha=1}^{|\mathcal{V}|} p(\mathbf{z})_\alpha \log p(\mathbf{z})_\alpha = \left( \log \sum_{\alpha=1}^{|\mathcal{V}|} e^{z_\alpha} - \frac{\sum_{\alpha=1}^{|\mathcal{V}|} z_\alpha e^{z_\alpha}}{\sum_{\alpha=1}^{|\mathcal{V}|} e^{z_\alpha}} \right), \tag{8}$$

where $p(\mathbf{z})_\alpha = \frac{e^{z_\alpha}}{\sum_{\beta=1}^{|\mathcal{V}|} e^{z_\beta}}$ is the probability of the $\alpha^{\text{th}}$ word in a vocabulary with $|\mathcal{V}|$ entries.

When the next token predictions are obtained using the softmax activation function, the contextual entropy reduces to the above expression, which we refer to as the softmax entropy[9].

From Equations 7 and 8, we see that the average contextual entropy for a prompt $X$ is the average of the softmax entropy of the corresponding logits -

$$\mathcal{H}(X) = \frac{1}{N} \sum_{i=1}^{N} S(\mathbf{z}_i). \tag{9}$$

These relations suggest that the underlying manifold on which the logits lie plays a role in the evaluation of the entropy. Given this manifold $\mathcal{M}$ the expected value of the softmax entropy is given by[10]

$$\langle S \rangle_{\mathcal{M}} = \frac{1}{\mathrm{Vol}(\mathcal{M})} \int_{\mathcal{M}} d\mathbf{z} \, S(\mathbf{z}). \tag{10}$$

From what is observed empirically, we expect that the dimension of the manifold $\mathcal{D}_{\mathcal{M}}$, typically much smaller than $|\mathcal{V}|$, should play a role in the integral. We consider the toy example example when the logits have $\mathcal{D}_{\mathcal{M}}$ activated entries that are independently and uniformly drawn from $[0, 1]$ and the other entries are large negative numbers, i.e. $z_\alpha \sim \mathcal{U}(0, 1), 1 \leq \alpha \leq \mathcal{D}_{\mathcal{M}}$ and $z_\alpha = -\infty, \mathcal{D}_{\mathcal{M}} < \alpha \leq |\mathcal{V}|$. This results in evaluating the integral on a unit box in $\mathcal{D}_{\mathcal{M}}$ dimensions ($[0, 1]^{\mathcal{D}_{\mathcal{M}}}$) because the inactive entries do not contribute towards the softmax entropy.

$$\langle S \rangle_{[0,1]^{\mathcal{D}_{\mathcal{M}}}} = \int_{[0,1]^{\mathcal{D}_{\mathcal{M}}}} d\mathbf{z} \, S(\mathbf{z}) = \int_0^1 \cdots \int_0^1 \left( \prod_{\alpha=1}^{\mathcal{D}_{\mathcal{M}}} dz_\alpha \right) \left( \log \sum_{\alpha=1}^{\mathcal{D}_{\mathcal{M}}} e^{z_\alpha} - \frac{\sum_{\alpha=1}^{\mathcal{D}_{\mathcal{M}}} z_\alpha e^{z_\alpha}}{\sum_{\alpha=1}^{\mathcal{D}_{\mathcal{M}}} e^{z_\alpha}} \right), \tag{11}$$

We evaluate this integral using numerical integration and compare the result with $\log \mathcal{D}_{\mathcal{M}}$

finding good agreement. In this toy example, we have related an information-theoretic quantity, the expected softmax entropy ($\langle S \rangle_{[0,1]^{\mathcal{D}_{\mathcal{M}}}}$) to a geometric quantity $\mathcal{D}_{\mathcal{M}}$, the dimension of the unit box, where we have

$$\langle S \rangle_{[0,1]^{\mathcal{D}_{\mathcal{M}}}} \sim \log \mathcal{D}_{\mathcal{M}}. \tag{12}$$

Another example involving logits with $\mathcal{D}_{\mathcal{M}}$ activated entries occurs when the next-token probabilities are uniformly distributed over the probability simplex $\Delta_{\mathcal{D}_{\mathcal{M}}}$. In this case, $p(\mathbf{z}) \in \Delta_{\mathcal{D}_{\mathcal{M}}}$ is drawn from the Dirichlet distribution with $\boldsymbol{\alpha} = 1$, where the expected entropy [52, 53] is given by

$$\langle S \rangle_{\Delta_{\mathcal{D}_{\mathcal{M}}}} = \psi(\mathcal{D}_{\mathcal{M}} + 1) - \psi(2) = \sum_{k=1}^{\mathcal{D}_{\mathcal{M}}} \frac{1}{k} - 1 \tag{13}$$

where $\psi$ is the Digamma function. From the above relation and using bounds on the harmonic number [54], it can be shown that

$$\left( \log \mathcal{D}_{\mathcal{M}} - \frac{1}{2} \right) < \langle S \rangle_{\Delta_{\mathcal{D}_{\mathcal{M}}}} \leq \log \mathcal{D}_{\mathcal{M}} \tag{14}$$

and in the asymptotic limit,

$$\lim_{\mathcal{D}_{\mathcal{M}} \to \infty} \langle S \rangle_{\Delta_{\mathcal{D}_{\mathcal{M}}}} = \log \mathcal{D}_{\mathcal{M}} + \gamma - 1 \sim \log \mathcal{D}_{\mathcal{M}} - 0.42 \tag{15}$$

where $\gamma$ is the Euler-Mascheroni constant. Similar to the unit box, we observe a $\log \mathcal{D}_{\mathcal{M}}$ dependence for the expected entropy in the probability simplex $\Delta_{\mathcal{D}_{\mathcal{M}}}$. While this relation might not hold for a generic manifold, it would be worth investigating this in more detail. We reserve this for future work.

## 6 Conclusions

The primary aim of this study was to connect different approaches to the interpretability of LLMs. Our strategy towards this goal was to examine the geometric structure of token-level representations across the layers of these models and to

---

[9] We use $S(\mathbf{z})$ to denote the softmax entropy that is defined at the level of logits and $H(x_{<i})$ to denote the contextual entropy which is more generically defined at the level of tokens. For clarity, we use the Greek letters to indicate the index in vocabulary and the Roman letters to indicate indices of tokens in a prompt.

[10] For simplicity, we assume a uniform distribution of $\mathbf{z}$ on the manifold $\mathcal{M}$ here. This definition can be extended to a manifold when $\mathbf{z}$ is distributed according to an arbitrary probability density function on $\mathcal{M}$.

relate it to the probability distribution of the next token prediction. We employed three key metrics: cosine similarity, intrinsic dimension, and neighborhood overlap, to capture different aspects of this geometric structure. Our findings revealed that the intrinsic dimension of token representations peaks in the early to middle layers, with higher peaks in shuffled data, i.e. when syntactic and semantic structures are disrupted. Additionally, cosine similarity among tokens increases with shuffling, suggesting greater alignment of token vectors. The neighborhood overlap metric showed that structured data maintains more coherent token neighborhoods across layers, while increased shuffling reduces this consistency, reflecting the model's sensitivity to the input structure. We observe these features consistently across different models. All these analyses converge into the key finding of this paper, which is the correlation of the ID of token representations to the model's cross-entropy loss, implying that ID could be an important metric for evaluating model performance across different models.

This correlation should be notably significant during the training process, especially in the context of developmental interpretability [55]. As demonstrated at the prompt level in previous research [20], and confirmed by our findings at the token level (Appendix E) ID remains largely constant across layers and low in the early training stages, but it increases as training progresses. At the token level, we observe that the ID tends to rise due to enhanced model expressivity, while there is also a tendency for ID to decline as the minimization of loss improves. Indeed, as seen in Figure 21 in Appendix, ID initially rises and then shows a slight decrease after checkpoint 64K (this behavior was observed in [56]). We believe it would be intriguing to explore these aspects in greater depth, but we defer this investigation to future work.

Experiments could be improved in several directions: first, we computed our observables at low ranges of nearest neighbors. For a more holistic approach, a multiscale analysis can reveal further relations among these observables. Secondly, the differences in distribution patterns for structured versus shuffled data, as suggested by cosine similarity and ID studies, might encode essential information on how tokens are distributed in space in the two cases. It is interesting to consider other geometric observables and understand their relation to the next token probabilities. These targeted explorations could provide practical applications for the design and training of LLMs, potentially leading to more interpretable and efficient models. While we show that the geometry of tokens encodes the next token prediction loss, we also potentially provide an unsupervised tool to understand how the model processes a given prompt.

### Reproducibility

The experiments were run on an NVIDIA H100 GPU with 94 GB memory. All the results contained in this work are reproducible using the repository found at link: https://github.com/RitAreaSciencePark/token_geometry.

### Acknowledgments

### References

[1] J. Vuckovic, A. Baratin, and R. T. des Combes, "A mathematical theory of attention," 2020. https://arxiv.org/abs/2007.02876.

[2] B. Geshkovski, C. Letrouit, Y. Polyanskiy, and P. Rigollet, "A mathematical perspective on transformers," 2024. https://arxiv.org/abs/2312.10794.

[3] A. Cowsik, T. Nebabu, X.-L. Qi, and S. Ganguli, "Geometric dynamics of signal propagation predict trainability of transformers," 2024. https://arxiv.org/abs/2403.02579.

[4] A. Agrachev and C. Letrouit, "Generic controllability of equivariant systems and applications to particle systems and neural networks," 2024. https://arxiv.org/abs/2404.08289.

[5] B. Geshkovski, C. Letrouit, Y. Polyanskiy, and P. Rigollet, "The emergence of clusters in self-attention dynamics," in *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, eds., vol. 36, pp. 57026–57037. Curran Associates, Inc., 2023. https://proceedings.neurips.cc/paper_files/paper/2023/file/b2b3e1d9840eba17ad9bbf073e009afe-Paper-Conference.pdf.

[6] S. Anagnostidis, L. Biggio, L. Noci, A. Orvieto, S. P. Singh, and A. Lucchi, "Signal propagation in transformers: Theoretical perspectives and the role of rank collapse," in *Advances in Neural Information Processing Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, eds. 2022. https://openreview.net/forum?id=FxVH7iToXS.

[7] H. Shi, J. GAO, H. Xu, X. Liang, Z. Li, L. Kong, S. M. S. Lee, and J. Kwok, "Revisiting over-smoothing in BERT from the perspective of graph," in *International Conference on Learning Representations*. 2022. https://openreview.net/forum?id=dUV91uaXm3.

[8] X. Wu, A. Ajorlou, Z. Wu, and A. Jadbabaie, "Demystifying oversmoothing in attention-based graph neural networks," in *Thirty-seventh Conference on Neural Information Processing Systems*. 2023. https://openreview.net/forum?id=Kg65qieiuB.

[9] B. He, J. Martens, G. Zhang, A. Botev, A. Brock, S. L. Smith, and Y. W. Teh, "Deep transformers without shortcuts: Modifying self-attention for faithful signal propagation," in *The Eleventh International Conference on Learning Representations*. 2023. https://openreview.net/forum?id=NPrsUQgMjKK.

[10] X. Wu, A. Ajorlou, Y. Wang, S. Jegelka, and A. Jadbabaie, "On the role of attention masks and layernorm in transformers," in *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. 2024. https://openreview.net/forum?id=lIH6oCdppg.

[11] N. Belrose, Z. Furman, L. Smith, D. Halawi, I. Ostrovsky, L. McKinney, S. Biderman, and J. Steinhardt, "Eliciting latent predictions from transformers with the tuned lens," 2023. https://arxiv.org/abs/2303.08112.

[12] S. Jastrzebski, D. Arpit, N. Ballas, V. Verma, T. Che, and Y. Bengio, "Residual connections encourage iterative inference," in *International Conference on Learning Representations*. 2018. https://openreview.net/forum?id=SJa9iHgAZ.

[13] nostalgebraist, "interpreting gpt: the logit lens," *LessWrong* (2020). https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens.

[14] N. Elhage, N. Nanda, C. Olsson, T. Henighan, N. Joseph, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, N. DasSarma, D. Drain, D. Ganguli, Z. Hatfield-Dodds, D. Hernandez, A. Jones, J. Kernion, L. Lovitt, K. Ndousse, D. Amodei, T. Brown, J. Clark, J. Kaplan, S. McCandlish, and C. Olah, "A mathematical framework for transformer circuits," *Transformer Circuits Thread* (2021) . https://transformer-circuits.pub/2021/framework/index.html.

[15] A. Ansuini, A. Laio, J. H. Macke, and D. Zoccolan, "Intrinsic dimension of data representations in deep neural networks," in *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. Curran Associates Inc., Red Hook, NY, USA, 2019.

[16] D. Doimo, A. Glielmo, A. Ansuini, and A. Laio, "Hierarchical nucleation in deep neural networks," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, eds., vol. 33, pp. 7526–7536. Curran Associates, Inc., 2020.

[17] P. Pope, C. Zhu, A. Abdelkader, M. Goldblum, and T. Goldstein, "The intrinsic dimension of images and its impact on learning," in *International Conference on Learning Representations*. 2021. https://openreview.net/forum?id=XJk19XzGq2J.

[18] L. Valeriani, D. Doimo, F. Cuturello, A. Laio, A. Ansuini, and A. Cazzaniga, "The geometry of hidden representations of large transformer models," in *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, eds., vol. 36, pp. 51234–51252. Curran Associates, Inc., 2023. https://proceedings.neurips.cc/paper_files/paper/2023/file/a0e66093d7168b40246af1cddc025daa-Paper-Conference.pdf.

[19] E. Cheng, C. Kervadec, and M. Baroni, "Bridging information-theoretic and geometric compression in language models," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, H. Bouamor, J. Pino, and K. Bali, eds., pp. 12397–12420. Association for Computational Linguistics, Singapore, Dec., 2023. https://aclanthology.org/2023.emnlp-main.762.

[20] E. Cheng, D. Doimo, C. Kervadec, I. Macocco, J. Yu, A. Laio, and M. Baroni, "Emergence of a high-dimensional abstraction phase in language transformers," 2024. https://arxiv.org/abs/2405.15471.

[21] D. Doimo, A. P. Serra, A. ansuini, and A. Cazzaniga, "The representation landscape of few-shot learning and fine-tuning in large language models," in *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. 2024. https://openreview.net/forum?id=nmUkwoOHFO.

[22] Y. Gardinazzi, G. Panerai, K. Viswanathan, A. Ansuini, A. Cazzaniga, and M. Biagetti, "Persistent topological features in large language models," 2024. https://arxiv.org/abs/2410.11042.

[23] J. B. Hamrick and S. Mohamed, "Levels of analysis for machine learning," *CoRR* **abs/2004.05107** (2020) , 2004.05107. https://arxiv.org/abs/2004.05107.

[24] A. Rogers, O. Kovaleva, and A. Rumshisky, "A primer in BERTology: What we know about how BERT works," *Transactions of the Association for Computational Linguistics* **8** (2020) 842–866. https://aclanthology.org/2020.tacl-1.54.

[25] Y. Belinkov, "Probing classifiers: Promises, shortcomings, and advances," *Computational Linguistics* **48** no. 1, (Mar., 2022) 207–219. https://aclanthology.org/2022.cl-1.7.

[26] Y. Belinkov, S. Gehrmann, and E. Pavlick, "Interpretability and analysis in neural NLP," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, A. Savary and Y. Zhang, eds., pp. 1–5. Association for Computational Linguistics, Online, July, 2020. https://aclanthology.org/2020.acl-tutorials.1.

[27] M. S. Lewicki and T. J. Sejnowski, "Learning Overcomplete Representations," *Neural Computation* **12** no. 2, (Feb., 2000) 337–365. https://doi.org/10.1162/089976600300015826. _eprint: https://direct.mit.edu/neco/article-pdf/12/2/337/814391/089976600300015826.pdf.

[28] H. Lee, A. Battle, R. Raina, and A. Ng, "Efficient sparse coding algorithms," in *Advances in Neural Information Processing Systems*, B. Schölkopf, J. Platt, and T. Hoffman, eds., vol. 19. MIT Press, 2006. https://proceedings.neurips.cc/paper_files/paper/2006/file/2d71b2ae158c7c5912cc0bbde2bb9d95-Paper.pdf.

[29] M. Faruqui, Y. Tsvetkov, D. Yogatama, C. Dyer, and N. A. Smith, "Sparse overcomplete word vector representations," in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, C. Zong and M. Strube, eds., pp. 1491–1500. Association for Computational Linguistics, Beijing, China, July, 2015. https://aclanthology.org/P15-1144.

[30] M. Geva, A. Caciularu, K. Wang, and Y. Goldberg, "Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Y. Goldberg, Z. Kozareva, and Y. Zhang, eds., pp. 30–45. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, Dec., 2022. https://aclanthology.org/2022.emnlp-main.3.

[31] K. Pal, J. Sun, A. Yuan, B. Wallace, and D. Bau, "Future lens: Anticipating subsequent tokens from a single hidden state," in *Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL)*, J. Jiang, D. Reitter, and S. Deng, eds., pp. 548–560. Association for Computational Linguistics, Singapore, Dec., 2023. https://aclanthology.org/2023.conll-1.37.

[32] N. Karagodin, Y. Polyanskiy, and P. Rigollet, "Clustering in causal attention masking," in *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. 2024. https://openreview.net/forum?id=OiVxYf9trg.

[33] B. Geshkovski, P. Rigollet, and D. Ruiz-Balet, "Measure-to-measure interpolation using transformers," 2024. https://arxiv.org/abs/2411.04551.

[34] B. Geshkovski, H. Koubbi, Y. Polyanskiy, and P. Rigollet, "Dynamic metastability in the self-attention model," 2024. https://arxiv.org/abs/2410.06833.

[35] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep learning*, vol. 1. MIT Press, 2016.

[36] S. Chung, D. D. Lee, and H. Sompolinsky, "Classification and geometry of general perceptual manifolds," *Physical Review X* **8** no. 3, (2018) 031003.

[37] U. Cohen, S. Chung, D. D. Lee, and H. Sompolinsky, "Separability and geometry of object manifolds in deep neural networks," *Nature communications* **11** no. 1, (2020) 746.

[38] S. Acevedo, A. Rodriguez, and A. Laio, "Unsupervised detection of semantic correlations in big data," 2024. `https://arxiv.org/abs/2411.02126`.

[39] Y. Zhao, T. Behnia, V. Vakilian, and C. Thrampoulidis, "Implicit geometry of next-token prediction: From language sparsity patterns to model representations," 2024. `https://arxiv.org/abs/2408.15417`.

[40] E. Tulchinskii, K. Kuznetsov, K. Laida, D. Cherniavskii, S. Nikolenko, E. Burnaev, S. Barannikov, and I. Piontkovskaya, "Intrinsic dimension estimation for robust detection of AI-generated texts," in *Thirty-seventh Conference on Neural Information Processing Systems*. 2023. `https://openreview.net/forum?id=8uOZ0kNji6`.

[41] E. Facco, M. d'Errico, A. Rodriguez, and A. Laio, "Estimating the intrinsic dimension of datasets by a minimal neighborhood information," *Scientific Reports* **7** no. 1, (Sep, 2017) 12140.

[42] F. Denti, D. Doimo, A. Laio, and A. Mira, "Distributional results for model-based intrinsic dimension estimators," 2021. `https://arxiv.org/abs/2104.13832`.

[43] D. Doimo, A. Glielmo, A. Ansuini, and A. Laio, "Hierarchical nucleation in deep neural networks," in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20. Curran Associates Inc., Red Hook, NY, USA, 2020.

[44] Meta, "Introducing meta llama 3: The most capable openly available llm to date," 2024. `https://ai.meta.com/blog/meta-llama-3/`.

[45] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. de las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L. R. Lavaud, M.-A. Lachaux, P. Stock, T. L. Scao, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed, "Mistral 7b," 2023. `https://arxiv.org/abs/2310.06825`.

[46] S. Biderman, H. Schoelkopf, Q. Anthony, H. Bradley, K. O'Brien, E. Hallahan, M. A. Khan, S. Purohit, U. S. Prashanth, E. Raff, *et al.*, "Pythia: a suite for analyzing large language models across training and scaling," in *Proceedings of the 40th International Conference on Machine Learning*, pp. 2397–2430. 2023.

[47] L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite, N. Nabeshima, *et al.*, "The Pile: An 800GB dataset of diverse text for language modeling," *arXiv preprint arXiv:2101.00027* (2020) .

[48] N. Nanda, "Pile-10k dataset," 2022. `https://huggingface.co/datasets/NeelNanda/pile-10k`.

[49] K. Ethayarajh, "How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, K. Inui, J. Jiang, V. Ng, and X. Wan, eds., pp. 55–65. Association for Computational Linguistics, Hong Kong, China, Nov., 2019. `https://aclanthology.org/D19-1006`.

[50] V. W. Liang, Y. Zhang, Y. Kwon, S. Yeung, and J. Y. Zou, "Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning," *Advances in Neural Information Processing Systems* **35** (2022) 17612–17625.

[51] E. Thordsen and E. Schubert, "Abid: Angle based intrinsic dimensionality," in *Similarity Search and Applications: 13th International Conference, SISAP 2020, Copenhagen, Denmark, September 30 – October 2, 2020, Proceedings*, p. 218–232. Springer-Verlag, Berlin, Heidelberg, 2020. `https://doi.org/10.1007/978-3-030-60936-8_17`.

[52] D. H. Wolpert and D. R. Wolf, "Estimating functions of probability distributions from a finite set of samples," *Phys. Rev. E* **52** (Dec, 1995) 6841–6854. `https://link.aps.org/doi/10.1103/PhysRevE.52.6841`.

[53] I. Nemenman, F. Shafee, and W. Bialek, "Entropy and inference, revisited," in *Proceedings of the 15th International Conference on Neural Information Processing Systems: Natural and Synthetic*, NIPS'01, p. 471–478. MIT Press, Cambridge, MA, USA, 2001.

[54] M. V. (https://math.stackexchange.com/users/218419/mark viola), "On the harmonic number ($h_n$) upper and lower "classical" bounds: which of those is closest to $h_n$?" Mathematics stack exchange. `https://math.stackexchange.com/q/2534095`. URL:https://math.stackexchange.com/q/2534095 (version: 2017-11-23).

[55] J. Hoogland, G. Wang, M. Farrugia-Roberts, L. Carroll, S. Wei, and D. Murfet, "The developmental landscape of in-context learning," 2024. `https://arxiv.org/abs/2402.02364`.

[56] A. Razzhigaev, M. Mikhalchuk, E. Goncharova, I. Oseledets, D. Dimitrov, and A. Kuznetsov, "The shape of learning: Anisotropy and intrinsic dimensions in transformer-based models," in *Findings of the Association for Computational Linguistics: EACL 2024*, Y. Graham and M. Purver, eds., pp. 868–874. Association for Computational Linguistics, St. Julian's, Malta, Mar., 2024.
`https://aclanthology.org/2024.findings-eacl.58`.

[57] R. Antonello and E. Cheng, "Evidence from fmri supports a two-phase abstraction process in language models," in *UniReps: 2nd Edition of the Workshop on Unifying Representations in Neural Models*.

[58] R. Sarfati, T. J. B. Liu, N. Boullé, and C. J. Earls, "Lines of thought in large language models," 2024.
`https://arxiv.org/abs/2410.01545`.

[59] K. Johnsson, C. Soneson, and M. Fontes, "Low bias local intrinsic dimension estimation from expected simplex skewness," *IEEE transactions on pattern analysis and machine intelligence* **37** no. 1, (Jan, 2015) 196–202.

## A Consistency Checks for the Shuffle Experiment

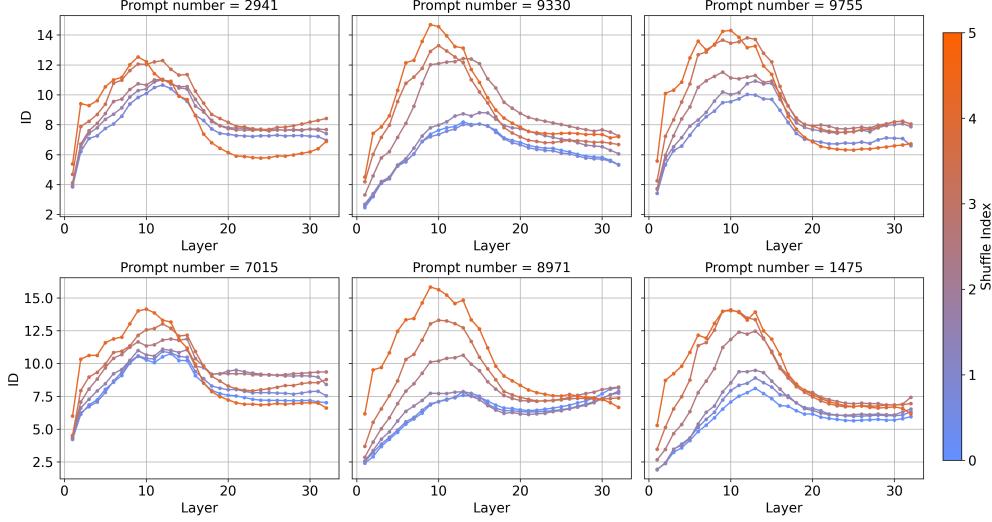In this section, we show the consistency of the results that were discussed in Section 4.



Figure 8: **ID profiles of 6 random prompts for LLAMA.** The prompts are taken from the filtered version of Pile described in the dataset section 3.1, but the prompt numbers refer to the Pile-10K dataset. The ID profiles are calculated using GRIDE at scaling = 2.
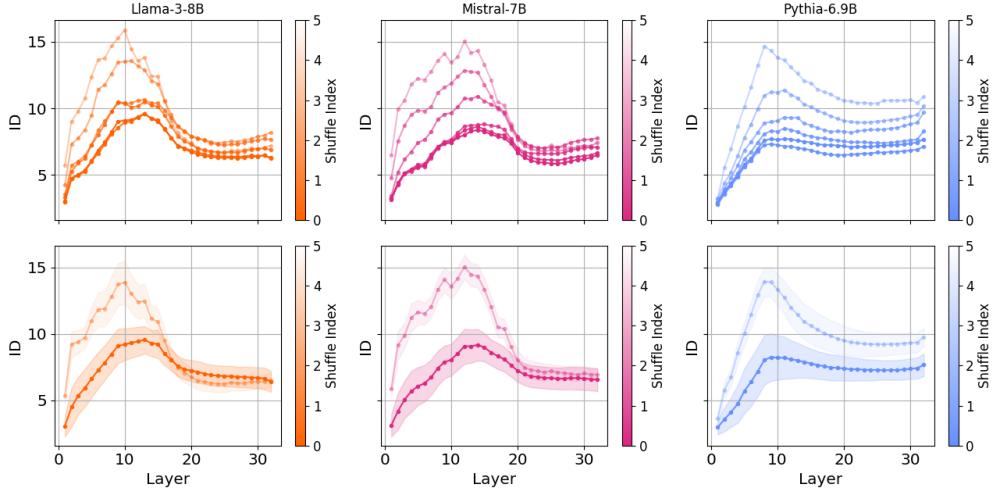


Figure 9: **ID profiles for shuffling for different models.** (a) ID profiles for prompt number 3218 from Pile-10K for different models across different levels of shuffling. Lighter colors represent a higher shuffle index, and darker colors indicate a more structured prompt, and (b) ID profiles for Pile-10K prompt number, averaged over 50 prompts, for both structured and fully shuffled cases. Lighter colors indicate higher shuffle indices and darker colors represent a more structured prompt. The shaded regions show the standard deviation from the mean.

**Intrinsic Dimension.** For the case of intrinsic dimension, we show the ID profiles of 6 random prompts sampled. It can be seen from Fig. 8 that the shuffled ID (orange) peak is always higher than the structured ID peak (blue) even though the degree of difference varies across prompts. We also verify that this behavior is consistent across models in Fig. 9.
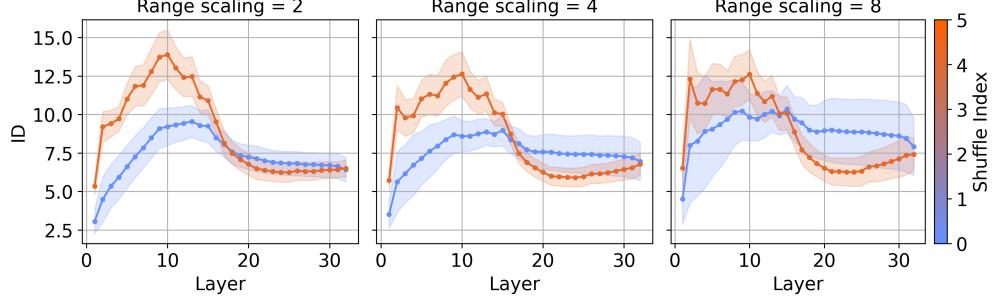
Figure 10: **Structured vs Shuffled ID for different range scalings.** Intrinsic dimension at scaling = $2, 4, 8$ as a function of layer for the full shuffle and structured case for the average over all the prompts for LLAMA.

**Neighborhood Overlap.** We compute the neighborhood overlap for the average over all the prompts of the full shuffle and the structured case using $k_{NN} = 1$ to $k_{NN} = 6$ in Figure 11. Results are consistent with what was discussed in 4.3.
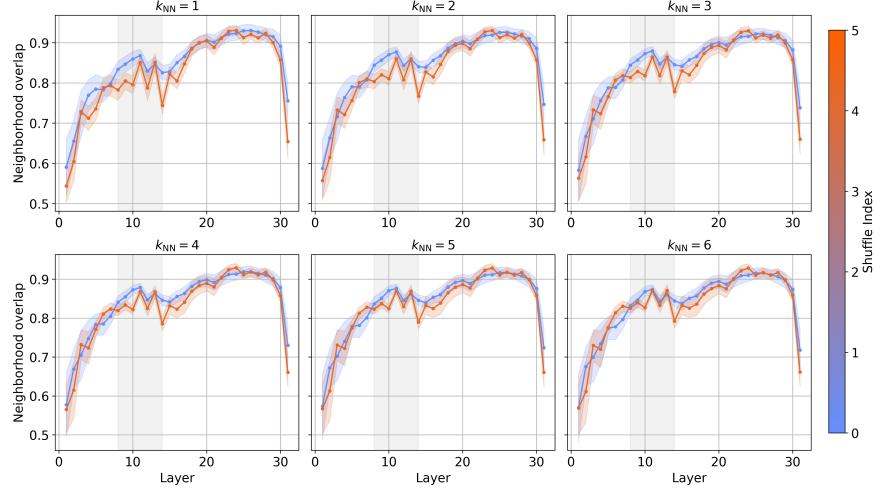


Figure 11: **Structured vs Shuffled NO for different $k_{NN}$.** Neighborhood overlap at $k_{NN} = 1$ to $k_{NN} = 6$ as a function of layer for the full shuffle and structured case for the average over all the prompts for LLAMA.

## B   Token Geometry of Prompts in Different Models

In the previous section, we noted that the geometry of internal representations is highly sensitive to shuffled inputs. Having focused on the representations from LLAMA model, we now extend our analysis to include two additional models: MISTRAL and PYTHIA. As described in Section 3.1, we note that PYTHIA was trained entirely on the Pile dataset. Hence, the dataset we consider for experiments, Pile-10K, is a subset of the same dataset on which PYTHIA was trained. While we do not know on which datasets LLAMA and MISTRAL were trained, we can assume that, if present, Pile was not the only dataset used. Therefore, we might expect PYTHIA to have a mildly different signature on our observables compared to MISTRAL and LLAMA. According to what we found in the previous section, we might expect a lower ID peak and a higher NO for PYTHIA.

### B.1   Intrinsic Dimension

We check the ID behaviour for LLAMA, MISTRAL and PYTHIA as a function of layers in Figure 12. On the left panel, we have the ID curve for a random prompt, while on the right panel, we show the mean ID profile across 2244 prompts.

We observe that PYTHIA has a lower ID peak

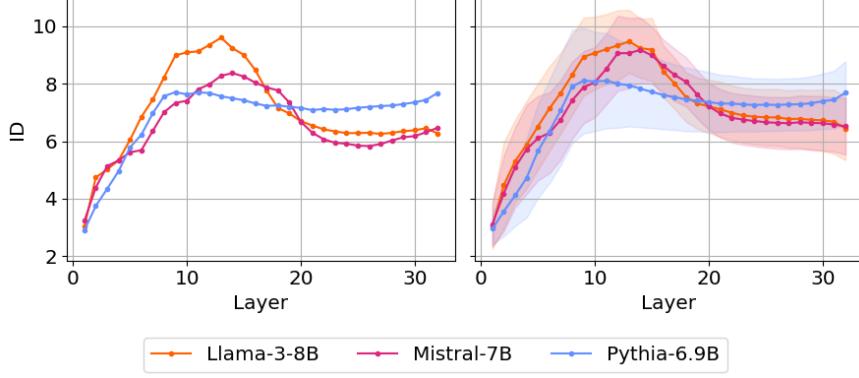on average than the other two models, though the significance is low.



Figure 12: **Intrinsic Dimension.** Left Panel: intrinsic dimension for a single prompt as a function of layers. Right Panel: intrinsic dimension averaged over 2244 prompts as a function of layers. The shaded regions indicate standard deviation from the mean. The curves correspond to LLAMA (orange), MISTRAL (magenta) and PYTHIA (blue).

## B.2   Neighborhood Overlap.

Similarly, we calculate NO and show it in Figure 13 as a function of layers for a random prompt (Left Panel) and the average over 2244 prompts (Right Panel).
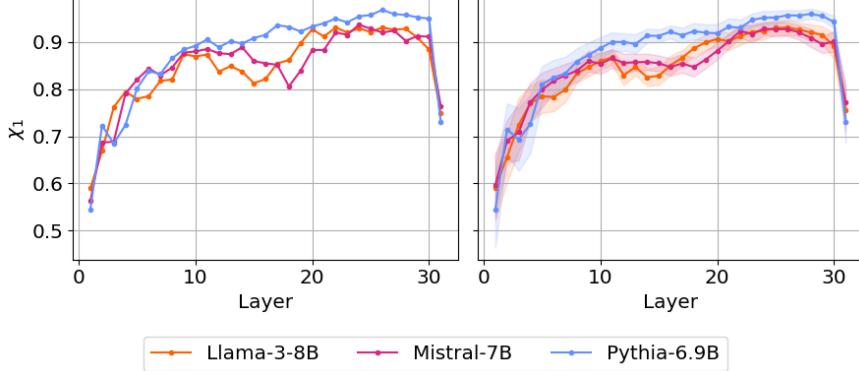


Figure 13: **Neighborhood Overlap.** Left Panel: neighborhood overlap for a single prompt as a function of layers. Right Panel: intrinsic dimension averaged over 2244 prompts as a function of layers. Shaded regions indicate standard deviation from the mean. The curves correspond to LLAMA (orange), MISTRAL (magenta) and PYTHIA (blue).

In this case, we observe that NO is generally higher for PYTHIA with respect to the other two models. The combined behavior of a lower ID peak and a higher NO in PYTHIA is similar to the structured case in the previous section. This might be a consequence of the fact that Pile is more in-distribution for PYTHIA than the other models. However, we note that a more comprehensive analysis would be required to confirm this statement, for instance by performing the analysis on PYTHIA using another dataset.

## C   Scale Analysis for GRIDE

In this section, we analyze the different choices of range scaling for the GRIDE algorithm discussed in Section 3. The prompts we analyze have $N = 1024$ tokens and in Fig. 14b, we check the dependence of ID estimate on range scaling $\in \{2, 4, 8, ..512\}$ for a single prompt on different models. This is to illustrate the scale dependence of a single prompt that we consider throughout the text.

In the main text, we focus on range scaling $= 2$ and here we extend the analysis to range scaling $= 4$ and $8$. In Figure 15, we find that PYTHIA's peak is more comparable to LLAMA and MISTRAL as the range scaling increases. In Figure 6, we notice that the correlation to loss becomes stronger for range scaling $= 4$ and $8$.



(a) Shuffled - GRIDE scale analysis for a shuffled prompt (prompt 3218) across layers.



(b) Unshuffled - GRIDE scale analysis for an unshuffled prompt (prompt number 3218) across layers.



(c) Scale analysis for GRIDE estimation across models averaged among prompts for different layers.
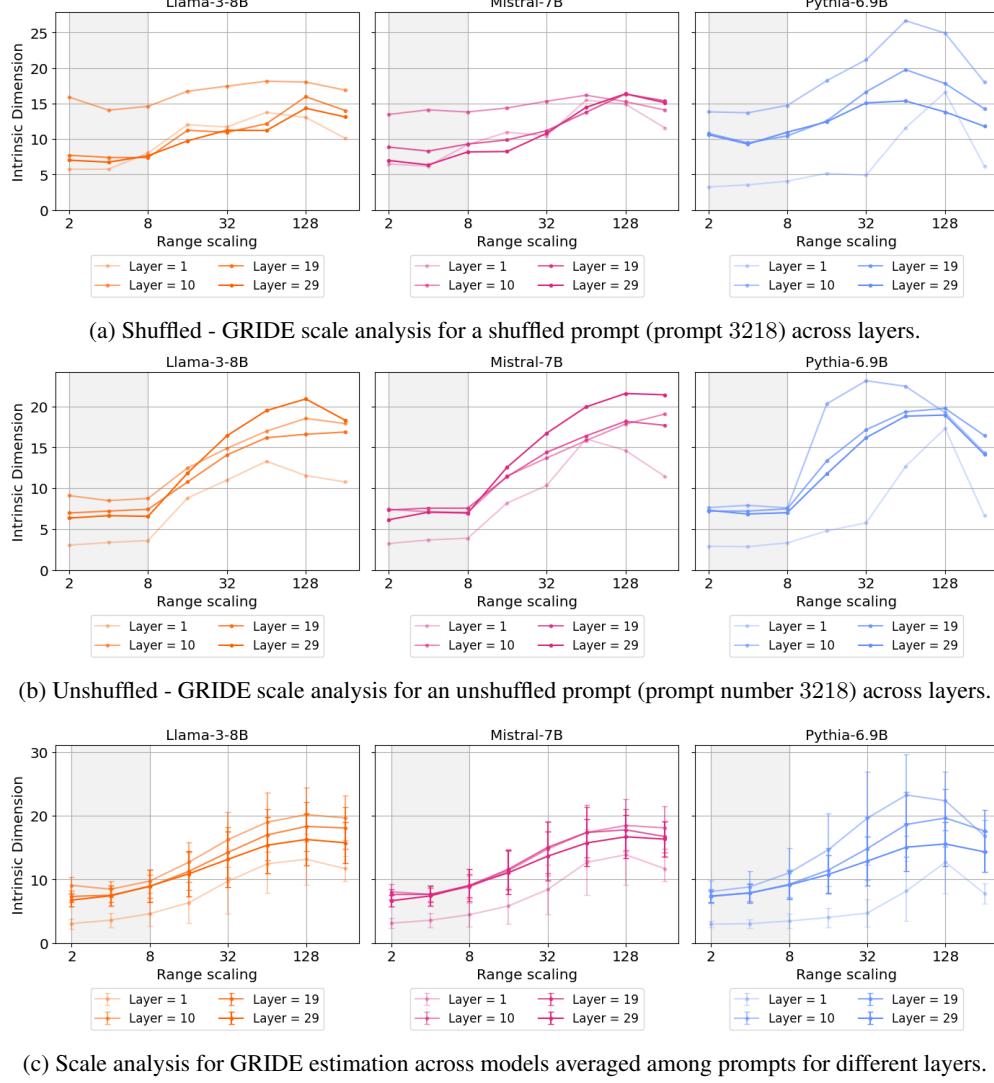
Figure 14: **Scale analysis for GRIDE estimation across models for shuffled and unshuffled prompts.** (a) Results for a single shuffled prompt (prompt number 3218), (b) Results for a single prompt (prompt number 3218), and (c) averaged results across unshuffled prompts, both showing different layers with early layers in lighter colors and late layers in darker colors.
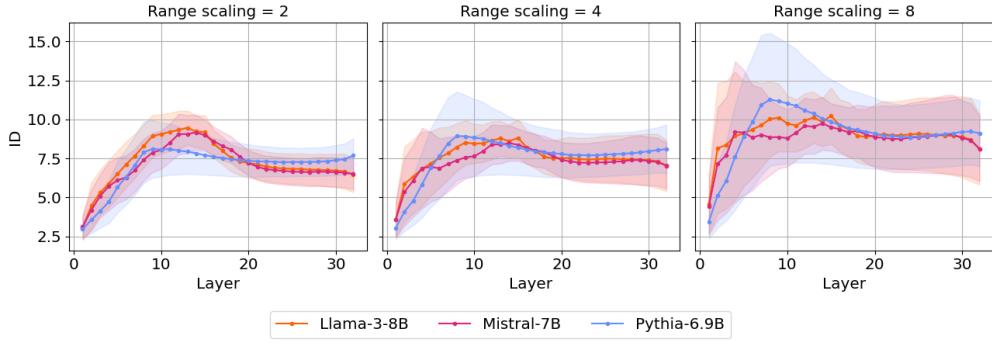
Figure 15: **Scale analysis for average ID profile.** The ID profile averaged over 2244 prompts for range scaling = $2, 4, 8$, with shaded regions indicating the standard deviation from the mean.
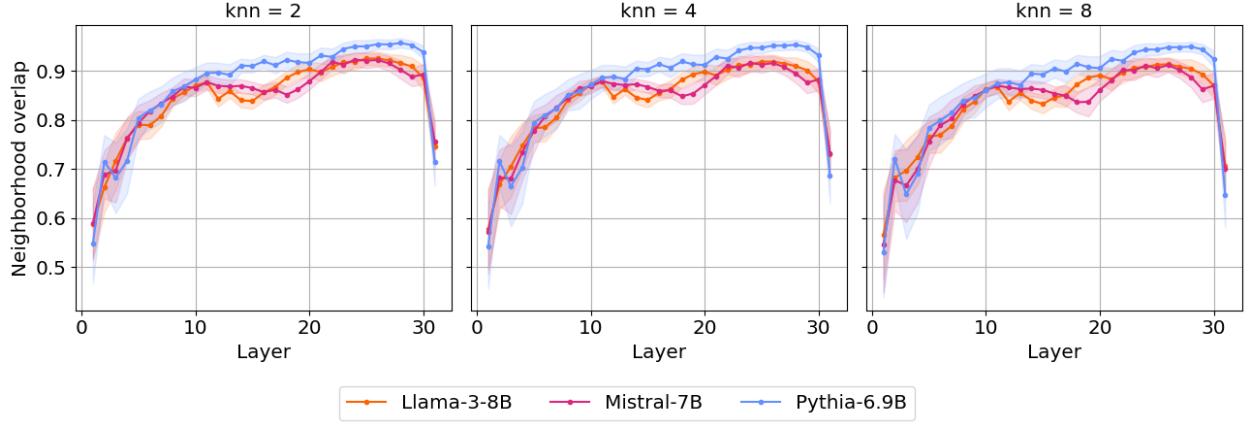


Figure 16: **Scale analysis for average NO profile.** The neighborhood overlap profile averaged over 2244 prompts for range scaling = $2, 4, 8$, with shaded regions indicating the standard deviation from the mean.
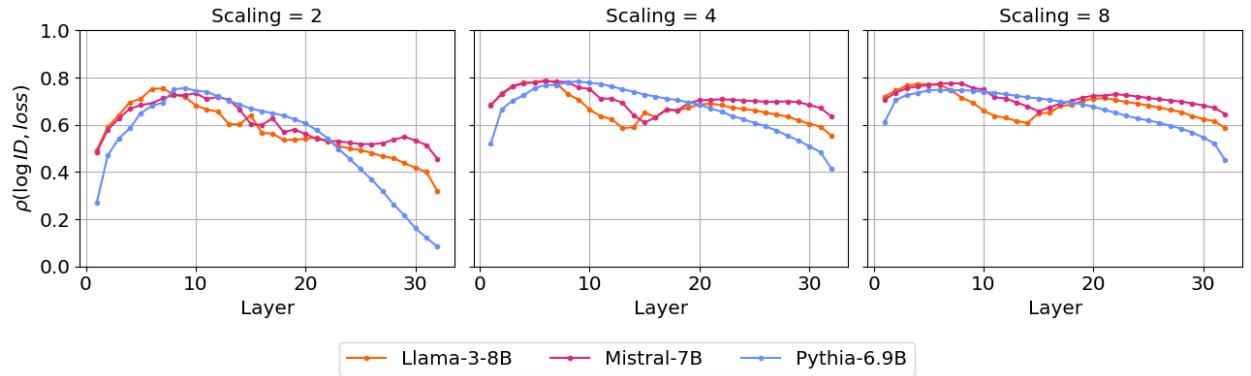


Figure 17: **Scale analysis for the correlation between intrinsic dimension and loss.** Pearson coefficient between the logarithm of intrinsic dimension and model loss at scalings = $2, 4, 8$ for different models.
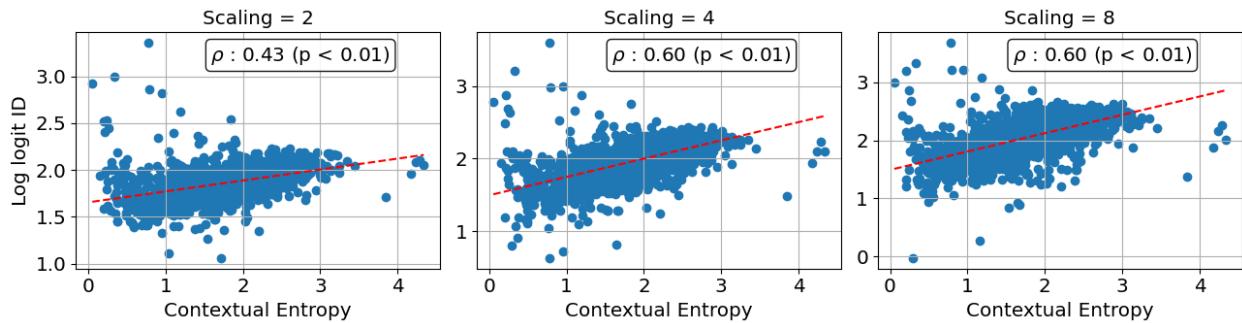
Figure 18: **Scale analysis for the correlation between intrinsic dimension of logits and contextual entropy.** Pearson coefficient between the logarithm of the intrinsic dimension of the logits and model contextual entropy for scalings $= 2, 4, 8$ for LLAMA.

# D   Qualitative Comparison of Token-level Geometry to Previous Prompt-level Studies

Previous work [15, 16, 17, 18, 19, 20, 57] have studied internal representations from a geometric point of view by considering point clouds of last token representations. While the approach is similar in spirit, token-level and prompt-level measures of intrinsic dimension probe different manifolds and thus different features of LLMs. We expect this because the relationship between the last tokens is different from that of tokens within the same prompt. In the upcoming analysis, we understand this difference intuitively by looking at the geometry of the shuffled and the unshuffled prompts at the prompt and token-level around the peak layers.

While prompt-level and token-level ID profiles exhibit similar behavior qualitatively, e.g. they peak in early-middle layers, there is a notable difference in the shuffled and unshuffled prompts. At the prompt level, we see that the unshuffled ID has a more prominent peak than the shuffled ID, whereas it is the other way around at the token level. In the shuffled case, the last token representations are less likely to share semantic content, leading to a lower intrinsic dimension at the prompt level. At the token level, the lesser prominence of the peak of the unshuffled case can be explained using the ID loss correlation. Since the loss is expected to be lower for the unshuffled prompts, we can expect their ID peak to be less prominent than that of the shuffled prompts.
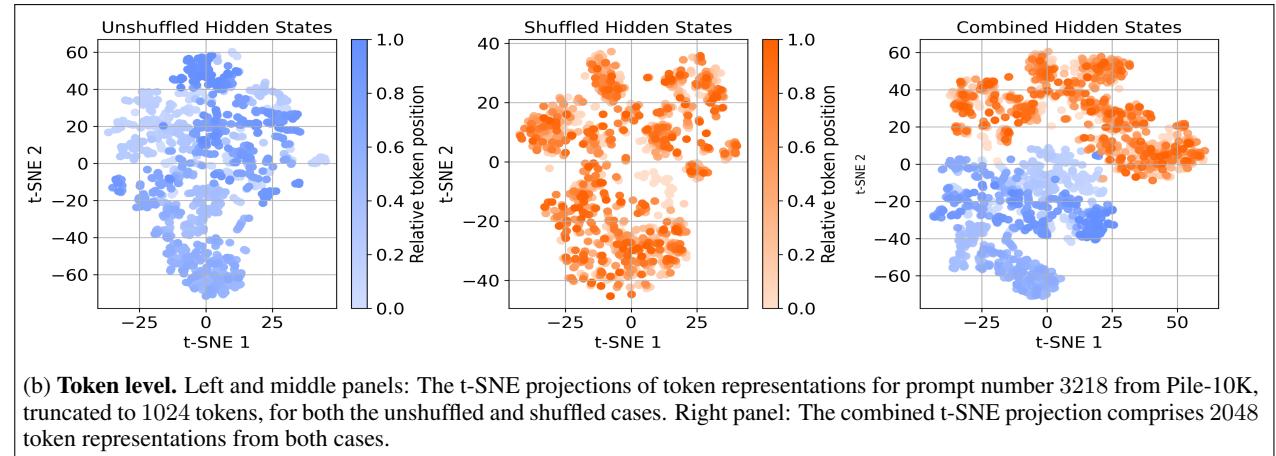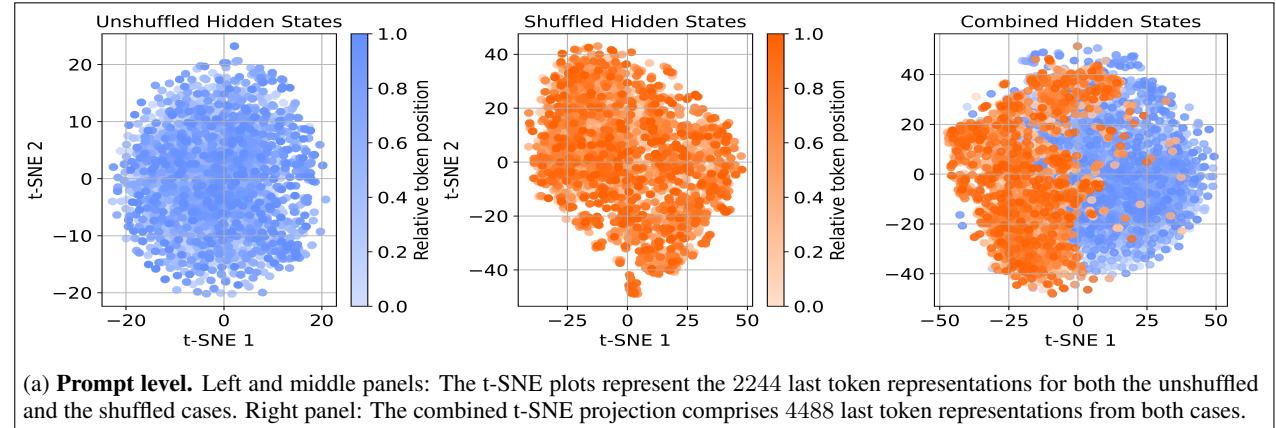


(a) **Prompt level.** Left and middle panels: The t-SNE plots represent the 2244 last token representations for both the unshuffled and the shuffled cases. Right panel: The combined t-SNE projection comprises 4488 last token representations from both cases.



(b) **Token level.** Left and middle panels: The t-SNE projections of token representations for prompt number 3218 from Pile-10K, truncated to 1024 tokens, for both the unshuffled and shuffled cases. Right panel: The combined t-SNE projection comprises 2048 token representations from both cases.

Figure 19: **Prompt geometry and token geometry.** A qualitative comparison of last-token representations at the prompt level (top panel) and the token level (bottom panel) geometry at layer 11 using t-SNE projections. All the plots are obtained using the representations from LLAMA.

For the prompt-level analysis, we use a corpus of 2244 prompts (the same corpus used for the token level analysis), drawn from Pile-10K and consisting of prompts with at least 1024 tokens. The last token representations are extracted from these prompts as follows - we choose tokens at positions 512 through 532 that result in a 20-token sequence for the unshuffled case[11]. We randomly permute aforementioned the 20-token sequences in the shuffled case and obtain

---

[11]This is a simplified setup of the experiments in [20].

the last token representations. The token-level analysis is done on prompt number 3218 from the Pile-10K dataset. In Figure 19, we plot the t-SNE projections of the shuffled and unshuffled along with ID for different scalings at both the prompt and token levels. We notice that in both levels, the shuffled and unshuffled representations lie on separate manifolds [58].

For the sake of completeness, we compare the results of the ID-loss correlation at the prompt and the token level in the next section.

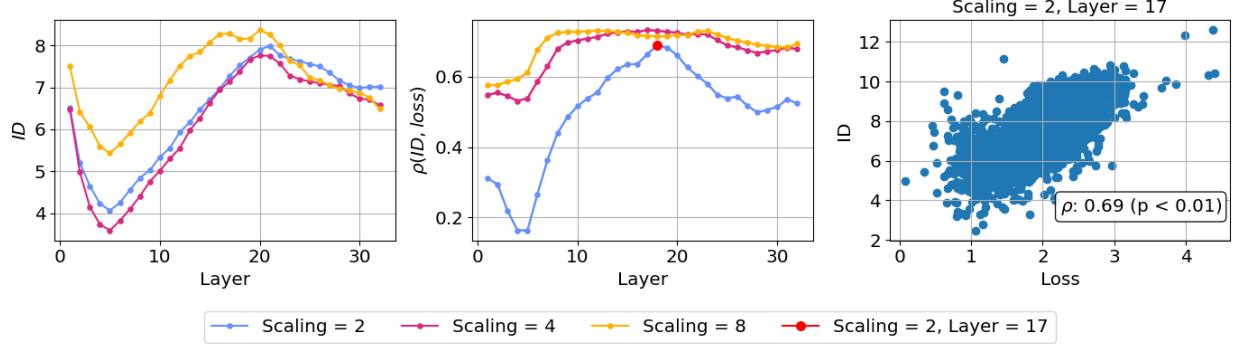### D.1 Token level ID is more strongly correlated to cross-entropy loss



Figure 20: **Summary of results for Opt-6.7B at the token-level.** Left panel: The ID curve for Opt-6.7B for scaling $= 2, 4, 8$ for prompt number $3218$ from Pile-10K. We observe a peak around layer $20$ as in the prompt level [20]. Middle panel: Spearman correlation between ID and loss for Opt-6.7B for different range scalings at the token level as a function of layers. Right panel: Scatter plot with the $ID$ (y-axis) and the average cross-entropy loss (x-axis) at scaling = 2, layer 17 for the $2244$ prompts we consider in this text.

Since there is an extensive amount of work done for the case of Opt-6.7B at the prompt level regarding the ID-cross entropy correlation, we compare the token level results to the prompt level for Opt-6.7B. Before proceeding here is a summary of the prompt level results from [19] and [20] that are relevant for our comparison.

- In [19], the authors show a positive Spearman correlation of $0.51$ for Opt-6.7B (Figure 2a in [19]) using the ID estimator Expected Simplex Skewness (ESS) [59] between ID at the peak and the cross-entropy loss.
- An analysis at a higher range scaling is done in [20] where they show a **negative correlation** with cross-entropy (Figure 6 in [20], there mentioned as surprisal) among a population consisting of different models and datasets with a relatively less statistical significance since it has a high $p$-value $= 0.09$.

On the other hand, using the token-level approach, we measure a **higher layerwise positive correlation** with cross-entropy. We summarize the results in Table 1.

| | Prompt level (ESS) | Prompt level (2NN) | Prompt level (high scaling) (many models $\times$ corpus) | Token level (2NN) | Token level (scaling = 8) |
|---|---|---|---|---|---|
| Spearman $\rho$ | 0.51 | 0.13 | -0.46 | 0.69 | 0.73 |
| $p$-value | 0.01 | 0.5 | 0.09 | $< 0.01$ | $< 0.01$ |

Table 1: Summary of Spearman correlations between ID and loss from prompt and token level analysis for Opt-6.7B. The results for token level are from Figure 20 and the prompt level are from [19] and [20].
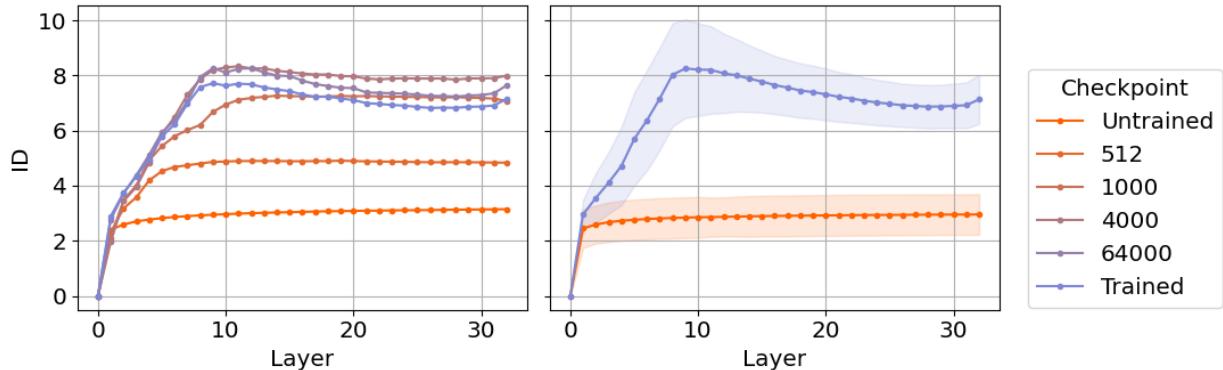
# E  Token-level ID during Training



Figure 21: **Intrinsic Dimension profile over training for PYTHIA.** Left Panel: intrinsic dimension profile for a single random prompt as a function of layers for different levels of training. Right Panel: intrinsic dimension averaged over 50 prompts as a function of layers for the untrained (orange) and trained (blue) model. The shaded regions indicate the standard deviation from the mean.