# **Hyperbolic Fine-Tuning for Large Language Models**

# Menglin Yang<sup>1,2</sup>, Ram Samarth B B<sup>3</sup>, Aosong Feng<sup>4</sup>, Bo Xiong<sup>5</sup> Jiahong Liu<sup>6</sup>, Irwin King<sup>6</sup>, Rex Ying<sup>4</sup>

<sup>1</sup>HKUST(GZ); <sup>2</sup>HKUST; <sup>3</sup>Indian Institute of Science; <sup>4</sup>Yale University; <sup>5</sup>Stanford University; <sup>6</sup>The Chinese University of Hong Kong menglin.yang@outlook.com,ramsamarthbb@iisc.ac.in, aosong.feng@yale.edu, xiongbo@stanford.edu, {jhliu22, king}@cse.cuhk.edu.hk, rex.ying@yale.edu Code: https://github.com/marlin-codes/HypLLM

Project\* https://hyperboliclearning.github.io/work/hyplora

## **Abstract**

Large language models (LLMs) have demonstrated remarkable performance on various tasks. However, it remains an open question whether the default Euclidean space is the most suitable choice for embedding tokens in LLMs. In this study, we investigate the non-Euclidean characteristics of LLMs. Our findings reveal that token frequency follows a power-law distribution, with high-frequency tokens clustering near the origin and low-frequency tokens positioned farther away. Additionally, token embeddings exhibit a high degree of hyperbolicity, indicating a latent tree-like structure in the embedding space. Motivated by these observations, we propose to efficiently fine-tune LLMs in hyperbolic space to better exploit the underlying complex structures. However, we find that this hyperbolic fine-tuning cannot be achieved through the naive application of exponential and logarithmic maps when the embedding and weight matrices both reside in Euclidean space. To address this technical issue, we introduce hyperbolic low-rank efficient fine-tuning, HypLoRA, which performs low-rank adaptation directly on the hyperbolic manifold, preventing the cancellation effect produced by consecutive exponential and logarithmic maps and thereby preserving hyperbolic modeling capabilities. Extensive experiments across various base models and two different reasoning benchmarks, specifically arithmetic and commonsense reasoning tasks, demonstrate that HypLoRA substantially improves LLM performance.

## 1 Introduction

Large language models (LLMs) such as GPT-4 [1], LLaMA [2], Gemma [3], and Qwen [4] have demonstrated remarkable capabilities in understanding and generating human-like text [5, 6, 7]. Despite their impressive capabilities, these models often rely on Euclidean geometry for token representation, which may inadequately capture the inherently complex and hierarchical nature of real-world data structures [8, 9, 10, 11]. Consider how words naturally organize into nested categories with varying levels of abstraction: abstract concepts like "fruit" occupy higher positions in the semantic hierarchy, while specific instances such as "apple" or "banana" populate the lower levels. Representing such structures effectively is crucial for understanding the semantics of language in LLMs.

Recent advancements suggest that non-Euclidean geometries, particularly hyperbolic spaces [11, 12, 13, 14, 15, 16, 17, 18, 19, 20], offer promising alternatives for modeling hierarchical data. Hyperbolic space, distinguished by its negative curvature, is especially well-suited for representing tree-like

<sup>\*</sup>Corresponding author: Menglin Yang. Project lead: Rex Ying

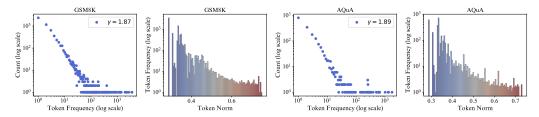


Figure 1: Token frequency distribution and token frequency vs. norm analysis for GSM8K and AQuA datasets in LLaMA3-8B. The left panels show the token frequency distributions (power-law distribution), while the right panels illustrate the relationship between token frequency and the corresponding norms. This visualization reveals the underlying geometric structure of the token embeddings. For additional data analysis and visualizations, please refer to Appendix A.

hierarchical data due to its exponential volume growth and geometric prior. This geometric property makes hyperbolic space especially useful for tasks involving complex, hierarchically structured information.

**Proposed Analysis Framework.** In this work, we first delve into how LLMs interact with token embeddings and explore the extent to which these embeddings exhibit non-Euclidean characteristics. We approach this from both a global and local perspective. At the global level, we analyze the overall distribution of tokens by frequency and investigate how these frequencies are arranged across the embedding space. At the local level, we measure the hyperbolicity of the metric space spanned by each input prompt, where the hyperbolicity of the Euclidean embedding serves as a proxy for evaluating the distance or dissimilarity between the underlying embedding structure and a tree-like hierarchy [21, 22, 14].

Our analysis in Section 4 reveals several key insights. First, token frequency follows a power-law distribution, which generally implies an underlying hierarchical structure similar to a branching tree [23, 24, 25, 26]. High-frequency tokens (e.g., abstract concepts, function words) tend to be located near the origin of the embedding space, while low-frequency tokens (e.g., specific terms) are farther away, as demonstrated in Table 1. Second, our investigation of hyperbolicity ( $\delta$  values) in Table 2 demonstrates that LLM token embeddings exhibit significant tree-like properties.

Based on our findings above, a natural consideration is to develop hyperbolic LLMs that explicitly incorporate a hyperbolic inductive bias<sup>2</sup>. However, training LLMs from scratch can be resource-intensive [27, 28, 29]. As a more resource-efficient alternative, we propose to build the first low-rank adaptation fine-tuning method in hyperbolic space. This approach is particularly advantageous given that existing LLMs are all Euclidean, and not all downstream tasks require hyperbolic geometry in their fine-tuning. By employing hyperbolic adapters for specific tasks on Euclidean LLMs, we can leverage the benefits of both geometries while maintaining computational efficiency.

**Challenges.** Adapting LLMs in non-Euclidean embedding spaces with classic techniques, *i.e.* applying exponential and logarithmic maps within tangent space [30, 31, 32, 33, 34] for weight adaptation is problematic in this case. This approach fails to fully capture the hyperbolic geometry, as the exponential and logarithmic maps are mutually inverse and can be canceled with consecutive operations<sup>3</sup>. Consequently, the inherent properties of the hyperbolic space are not effectively preserved, limiting the potential benefits of incorporating non-Euclidean geometries into the adaptation process.

**Proposed Method.** To address this limitation, we introduce HypLoRA to perform low-rank adaptation directly on the hyperbolic manifold without transformation to the tangent space, thus preserving hyperbolic modeling capabilities and counteracting the cancellation effect. HypLoRA integrates hyperbolic geometry into existing LLMs, implicitly introducing high-order interaction and accounting

<sup>&</sup>lt;sup>2</sup>The connection between power-law distribution and hyperbolic geometry is elaborated in Section 4.3.

 $<sup>^3</sup>$ The cancellation effect occurs because standard hyperbolic neural network [31, 30] approaches apply transformations in the tangent space at the origin, requiring the sequence: Euclidean embedding  $\rightarrow$  exponential map to hyperbolic space  $\rightarrow$  logarithmic map to tangent space  $\rightarrow$  linear transformation  $\rightarrow$  exponential map back to hyperbolic space  $\rightarrow$  projection to Euclidean space. When these operations are chained together, the maps are mutually inverse and effectively cancel out, reducing the entire sequence to approximately the original Euclidean transformation BAx without preserving the beneficial hyperbolic geometry.

for the token hierarchies, enabling them to benefit from hyperbolic characteristics while minimizing additional computational costs.

To summarize, our main contributions are twofold: (1) We conduct a comprehensive investigation into the hyperbolic characteristics of token embeddings in LLMs, revealing their inherent tree-like structure and strong hyperbolic properties. (2) We propose HypLoRA, a parameter-efficient fine-tuning method that integrates hyperbolic geometry into LLMs while preserving hyperbolic modeling capabilities. We conduct extensive experiments on various models and two different tasks, specifically arithmetic reasoning and commonsense reasoning, demonstrating clear advantages over competitive baselines.

## 2 Related Work

Hyperbolic Representation Learning and Foundation Models. Hyperbolic geometry has been successfully applied to various neural network architectures and models [17, 19, 16], including shallow hyperbolic neural networks [13, 31, 35, 36, 37, 38], hyperbolic CNNs [39, 40, 41], hyperbolic GNNs [30, 42, 43], and hyperbolic attention networks or Transformers [44, 35, 36, 45]. These models leverage the inductive biases of hyperbolic geometry to achieve remarkable performance on various tasks and applications [30, 20, 46, 14, 15, 47, 48, 49, 50]. Recent efforts have focused on adapting LLMs and CLIP [51] to hyperbolic spaces. Key advancements include developing more expressive hyperbolic image-text representations [52], enabling compositional entailment learning for deeper vision-language understanding [53], designing safety-aware hyperbolic frameworks for content moderation [54], and creating core modules to facilitate the construction of novel hyperbolic foundation models [55]. However, training LLMs from scratch remains computationally expensive [56, 57]. The computational complexity increases further when considering Riemannian optimization [56, 57, 58] and additional hyperbolic operations, like Möbius addition.

Geometric Analysis of Language Model Embeddings. Prior work has made important observations about the geometry of embeddings that helped motivate our research. Reif et al. [59] demonstrated that BERT embeddings contain distinct syntactic and semantic subspaces and showed evidence of tree-like parse structures, while Gao et al. [60] revealed that token embeddings tend to cluster in a narrow cone during training, leading to representation degeneration. Building on these geometric insights, Rudman et al. [61] introduced IsoScore to formally quantify how uniformly embeddings utilize the ambient vector space. Additionally, Puccetti et al. [62] analyzed outlier dimensions in Transformers and showed their correlation with token frequencies. While these works provide crucial foundations for understanding embedding geometry, our work differs in that we specifically quantify and leverage the natural hyperbolicity of token embeddings.

Parameter-Efficient Fine-Tuning (PEFT) and LoRA. Fine-tuning LLMs [63, 64, 2] for downstream tasks poses significant challenges due to their massive number of parameters. To address this issue, PEFT methods have been proposed, which aim to train a small subset of parameters while achieving better performance compared to full fine-tuning. PEFT methods can be broadly categorized into prompt-based methods [65, 66, 67], adapter-based methods [68, 69], and reparameterization-based methods [29, 70, 71]. Among these, the reparameterization-based LoRA [29] has gained significant attention due to its simplicity, effectiveness, and compatibility with existing model architectures. Variants of LoRA, such as LoRA+[72], DoRA [73], and AdaLoRA [74], have been proposed to improve its performance and efficiency. Recent research has also investigated ensembles of multiple LoRAs [75, 76] and quantization techniques [77, 78, 79]. The proposed method is a foundational algorithm that is orthogonal to existing approaches and can potentially be combined with various LoRA variants to exploit their complementary strengths and achieve superior performance.

# 3 Preliminary

This section introduces the concepts used in our study, including the Lorentz model of hyperbolic geometry and the LoRA adapter.

**Hyperbolic Geometry.** Unlike the flat Euclidean geometry, hyperbolic geometry is characterized by a constant negative curvature. We utilize the Lorentz model, also known as the hyperboloid model, for our study due to its ability to effectively capture hierarchical structures and maintain numerical stability [12, 35, 80]. The Lorentz model in n dimensions with curvature -1/K(K > 0) is defined

as:

$$\mathcal{L}_K^n = \{ \mathbf{x} \in \mathbb{R}^{n+1} : \langle \mathbf{x}, \mathbf{x} \rangle_{\mathcal{L}} = -K, x_0 > 0 \}, \tag{1}$$

where  $\langle \cdot, \cdot \rangle_{\mathcal{L}}$  is the Lorentzian inner product, given by:  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{L}} = -x_0 y_0 + \sum_{i=1}^n \mathbf{x}_i \mathbf{y}_i$ .

**Tangent Space.** In the Lorentz model  $\mathcal{L}_K^n$ , the tangent space at a point  $\mathbf{x}$  is denoted as  $\mathcal{T}_{\mathbf{x}}\mathcal{L}_K^n$ . It is defined as the set of all vectors  $\mathbf{u}$  that are orthogonal to  $\mathbf{x}$  under the Lorentzian inner product:

$$\mathcal{T}_{\mathbf{x}}\mathcal{L}_K^n := \{ \mathbf{u} \in \mathbb{R}^{n+1} : \langle \mathbf{u}, \mathbf{x} \rangle_{\mathcal{L}} = 0 \}.$$
 (2)

To facilitate projection between the hyperboloid and its tangent spaces, one can utilize two critical mappings: the exponential and logarithmic maps. The *exponential map* at  $\mathbf{x}$ , denoted  $\exp_{\mathbf{x}}^K$ , projects a vector from the tangent space  $\mathcal{T}_{\mathbf{x}}\mathcal{L}_K^n$  back onto the hyperboloid. Conversely, the *logarithmic map*, denoted  $\log_{\mathbf{x}}^K$ , maps a point on the hyperboloid to the tangent space at  $\mathbf{x}$ . The relevant formulas are given in Appendix D.1.

**LoRA Adapter.** The LoRA adapter provides an efficient approach for modifying large LLMs with minimal computational overhead. Instead of retraining the entire model, LoRA focuses on adjusting specific components within the model's architecture to transform an input  $\mathbf{x}$  into an output  $\mathbf{z}$ . In practice, LoRA targets the weight matrices found in each Transformer layer of an LLM. Typically, the weight W of the Transformer, which resides in the dimensions  $\mathbb{R}^{d \times k}$ , is adapted through a low-rank approximation. This is achieved by introducing an additional term,  $\Delta W$ , to the original weight matrix:

$$\mathbf{z} = W_{\text{LoRA}}(\mathbf{x}) = W\mathbf{x} + \Delta W\mathbf{x} = W\mathbf{x} + BA\mathbf{x}.$$
 (3)

Here,  $B \in \mathbb{R}^{d \times r}$  and  $A \in \mathbb{R}^{r \times k}$  represent two smaller, learnable matrices where r—the rank of these matrices—is significantly less than either d or k. This design choice ensures that  $r \ll \min(d,k)$ , thereby reducing the complexity of the model adaptation. During the fine-tuning process, only the matrices A and B are adjusted, while the pre-existing weights W are kept frozen. This method significantly decreases the number of parameters that need to be trained, from dk to (d+k)r, enhancing the efficiency of the fine-tuning process. As a result, LoRA enables the targeted adaptation of LLMs, allowing them to transform an input  $\mathbf{x}$  into an output  $\mathbf{z}$  while maintaining high performance and adapting to new tasks or datasets with a fraction of the computational resources typically required.

# 4 Investigation

In this section, we present an in-depth investigation of token embeddings in LLMs from both global and local perspectives. Our goal is to uncover the geometric structures underlying pretrained token representations, specifically examining the global distribution of token frequencies and their spatial arrangement, as well as the local hyperbolicity of token embeddings across various datasets.

#### 4.1 Global Token Statistics

We begin by investigating the global distribution of token frequencies in the context of arithmetic reasoning datasets, focusing on datasets such as GSM8K [81], AQuA [82], MAWPS [83], and SVAMP [84]. We also provide a broader analysis across different types of datasets and LLMs in Appendix A. Figure 1 (left) presents the distribution of token frequencies, with a power-law exponent of approximately  $\gamma \approx 1.9$ , as estimated by the power-law package [85]. In such distributions, the exponent  $\gamma$  controls how quickly token frequencies decline: smaller values of  $\gamma$  (closer to 1) indicate a more gradual decay where frequent tokens dominate, while larger values signify a sharper decline, with most tokens being rare.

This power-law behavior aligns with the tree-like hierarchical nature of language [11, 14, 86, 49, 20]. High-frequency tokens often correspond to more abstract or general concepts, while low-frequency tokens represent specific or rare terms. This distribution naturally suggests a hierarchical organization of the token space, where general concepts serve as the "roots" and specific terms "branch out" as we move through the hierarchy.

**Empirical Observation.** To better understand the relationship between token frequency and their spatial arrangement within the embedding space, we calculate the average token frequency as a function of their distance from the origin. The results are shown in Figure 1 (right), indicating that more frequent tokens tend to have smaller norms and vice versa. Table 1 provides representative

Table 1: Mean, Minimum, and Maximum Frequency and Norm Values for Different Models and Groups. Group 1: to, in, have, that, and, is, for, Group 2: how, much, many, time, cost, Group 3: animal, fruit, number, color, size, Group 4: dog, cow, apple, banana, 380, 480, purple, red, medium, small, large.

Model	Group	Frequency (Mean [Min~Max])	Norm (Mean [Min~Max])
	Group 1	$4934.4 [1838 \sim 8539]$	$3.160 [3.060 \sim 3.299]$
Gemma-7B	Group 2	$2709.4 [474 \sim 6681]$	$3.561 [3.488 \sim 3.627]$
Genna-/B	Group 3	$292.0 [34 \sim 1191]$	$3.765 [3.623 \sim 3.887]$
	Group 4	$114.3 [25 \sim 284]$	$3.998 [3.660 \sim 4.520]$
	Group 1	$4993.9 [1838 \sim 8547]$	$0.951 [0.793 \sim 1.060]$
LLaMA-7B	Group 2	$2712.6 [474 \sim 6683]$	$1.222 [1.118 \sim 1.299]$
LLaWIA-/D	Group 3	$299.8 [34 \sim 1200]$	$1.325 [1.274 \sim 1.428]$
	Group 4	$139.1 [26 \sim 286]$	$1.364 [1.326 \sim 1.417]$
	Group 1	$4937.4 [1838 \sim 8547]$	$0.353 [0.330 \sim 0.396]$
LLaMA3-8B	Group 2	$2710.0 [474 \sim 6683]$	$0.456 [0.394 \sim 0.499]$
LLaWA3-0D	Group 3	$292.6 [34 \sim 1191]$	$0.499 [0.452 \sim 0.549]$
	Group 4	$97.1 [13 \sim 284]$	$0.569 [0.499 \sim 0.675]$
LLaMA-13B	Group 1	$4993.9 [1838 \sim 8547]$	$1.027 [0.833 \sim 1.255]$
	Group 2	$2712.6 [474 \sim 6683]$	$1.429 [1.346 \sim 1.489]$
LLawiA-13D	Group 3	$299.8 [34 \sim 1200]$	$1.494 [1.453 \sim 1.532]$
	Group 4	$139.1 [26 \sim 286]$	$1.501 [1.470 \sim 1.526]$

tokens with different norm ranges within the embedding space of different base models. The results presented in Table 1 demonstrate several critical findings. *First*, we observe a statistically significant separation between functional/abstract words (Group 1) and specific terms (Group 4) across all models, with Group 1 consistently exhibiting the smallest embedding norms and highest frequencies, while Group 4 shows the largest norms and lowest frequencies. *Second*, the relative ordering of groups remains consistent across all examined models, with Group 1 < Group 2 < Group 3 < Group 4 in terms of embedding norms, despite absolute magnitude variations. Most notably, even across different architectural families (LLaMA vs. Gemma), the hierarchical organization principle remains preserved, though with different absolute scales, where Gemma-7B exhibits systematically larger embedding norms (mean Group 1 norm: 3.160) compared to LLaMA models (mean Group 1 norm:  $0.353 \sim 1.027$ ), yet maintains the same relative hierarchical structure.

Conclusion. These findings suggest that the spatial organization of token embeddings reflects the inherent hierarchical relationships in language, supporting the hypothesis that token embedding in LLMs exhibits a tree-like structure, with spatial positioning aligned with token frequency and specificity. It is worth noting, however, that a power-law distribution of token frequency alone does not guarantee the emergence of a hierarchical token embedding, as it also depends on the training objectives. Our analysis demonstrates that the hierarchy is strongly correlated with token frequencies, which can be understood through the lens of LLMs tokenization and co-occurrence pattern learning during training. While the exact mechanisms underlying this relationship require further investigation in future work, the spatial distribution of token embeddings remains crucial as it provides the primary motivation for our methodological approach.

## 4.2 $\delta$ -Hyperbolicity of Local Token Embeddings

To rigorously quantify the hierarchical nature of token embeddings, we examine the  $\delta$ -hyperbolicity of the space spanned by the token embedding.  $\delta$ -Hyperbolicity, introduced by Gromov [87], is a measure that captures the degree to which a metric space deviates from an exact tree structure. Lower values of  $\delta$  imply a space more similar to a perfect tree, while higher values indicate deviation from a tree-like structure. A brief explanation of  $\delta$ -hyperbolicity can be found on Wikipedia<sup>4</sup>.

We compute  $\delta$ -hyperbolicity using the four-point condition, which compares the Gromov products between any four points a, b, c, and w in the metric space. Specifically, the hyperbolicity is defined as:

$$[a, c]_w \ge \min([a, b]_w, [b, c]_w) - \delta,$$
 (4)

<sup>4</sup>https://en.wikipedia.org/wiki/Hyperbolic\_metric\_space

Table 2: Comparison of  $\delta$ -Hyperbolicity across various metric spaces and datasets. The left table provides reference values for baseline metric spaces, allowing for a clearer interpretation of hyperbolicity in the analyzed datasets in the right table.

Metric Space	$\textbf{Hyperbolicity}(\delta)$
Sphere Space	$0.99 \pm 0.01$
Random Graph	$0.62 \pm 0.34$
PubMed Graph	$0.40 \pm 0.45$
Scale-free Graph	0.00
Tree Graph	0.00

$\overline{\textbf{Hyperbolicity}(\delta)}$	MAWPS	SVAMP	GSM8K	AQuA
LLaMA-7B	$0.08 \pm 0.02$	$0.09 \pm 0.01$	$0.10 \pm 0.01$	$0.10 \pm 0.01$
LLaMA-13B	$0.08 \pm 0.01$	$0.09 \pm 0.01$	$0.09 \pm 0.01$	$0.10 \pm 0.01$
Gemma-7B	$0.11 \pm 0.01$	$0.11 \pm 0.01$	$0.11 \pm 0.01$	$0.12 \pm 0.01$
LLaMA3-8B	$0.06 \pm 0.01$	$0.07 \pm 0.01$	$0.07 \pm 0.01$	$0.08 \pm 0.01$
Average	$0.08 \pm 0.01$	$0.09 \pm 0.01$	$0.09 \pm 0.01$	$0.10 \pm 0.01$

where the Gromov product  $[a, b]_w$  is:

$$[a,b]_w = \frac{1}{2}(d(a,w) + d(b,w) - d(a,b)).$$
(5)

Quantitative Analysis. To measure the hyperbolicity of token embeddings, we apply this algorithm to various open-source LLMs. Following the methodologies proposed by Khrulkov et al. [14] and Cetin et al. [15], we estimate  $\delta$ -hyperbolicity using the efficient algorithm introduced by Fournier et al. [88]. To ensure scale invariance, we normalize  $\delta$  by the diameter of the embedding space, diam(X), yielding a relative measure:  $\delta_{rel} = \frac{2\delta}{\operatorname{diam}(X)}$ . This relative measure ranges from 0 to 1, with values closer to 0 indicating a highly hyperbolic (tree-like) structure, and values near 1 indicating a non-hyperbolic, flat structure. Following previous works [14], we employ Euclidean distance as a measure of the shortest distance. To further validate the correctness of this approach, we generate a series of random graphs with predefined hyperbolicity, embed them using a graph neural network (GNN), and then compute the hyperbolicity in Euclidean space. Details of this process are provided in Appendix B. Our experiments reveal a positive correlation between the hyperbolicity of the embeddings and the original graphs. Consequently, we utilize this method as a proxy for estimating the hyperbolicity of token embeddings. In our analysis, we calculate hyperbolicity at the prompt level, treating each token within a prompt as a point in the metric space spanned by the embeddings. By averaging the hyperbolicity across all prompts, we assess the overall hyperbolic structure of token embeddings in each dataset.

**Conclusion.** Our results, as shown in Table 2, reveal that token embeddings exhibit significant hyperbolicity, suggesting that the embedding space has a strong tree-like structure. This observation further corroborates our findings from the global token statistics, where the arrangement of tokens in the embedding space mirrors hierarchical relationships seen in language data. We also provide the hyperbolicity analysis of the final hidden layer in Appendix A.3.

# 4.3 Connection between Power-law Distribution and Hyperbolic Geometry

The observation of a power-law distribution in token frequencies, as discussed in Section 4, is not merely a statistical curiosity. It has deep connections to the underlying geometry of the data, particularly to hyperbolic spaces, which are well-suited for representing hierarchical structures [11, 23, 89, 90]. For instance, Nickel and Kiela [11] highlighted that the existence of power-law degree distributions can often be traced back to hierarchical structures. Similarly, Ravasz and Barabási [86] established that the scaling law  $P(k) \sim k^{-\gamma}$  can signify the co-existence of a hierarchy of nodes with varying degrees of clustering. Krioukov et al. [23] further strengthened this connection by showing that the exponent of the power-law degree distribution is a function of the hyperbolic space curvature. Building on this geometric understanding, Papadopoulos et al. [90] demonstrated that complex (scale-free) network topologies naturally emerge when networks grow within an underlying hyperbolic metric space, and importantly, that the resulting hyperbolic embedding of these dynamic scale-free networks facilitates highly efficient greedy forwarding.

To formalize this connection with hyperbolic geometry, we can consider embedding tokens in a hyperbolic space. A common model for hyperbolic space is the Poincaré disk model ( $\mathbb{H}^2$ ) with curvature K=-1. In such a space, both the circumference C(r) and area A(r) of a circle of radius r exhibit exponential growth:

$$C(r) = 2\pi \sinh(r) \sim e^r \quad \text{as } r \to \infty,$$
 (6)

$$A(r) = 2\pi(\cosh(r) - 1) \sim e^r \quad \text{as } r \to \infty.$$
 (7)

Table 3: Accuracy comparison of various LLMs on arithmetic reasoning tasks. The percentage following each dataset indicates the proportion of prompts relative to the total number of inference prompts. M.AVG represents the micro-average accuracy. For a full comparison table, please see Appendix F.

Base Model	PEFT Method	# Params (%)	MAWPS(8.5%)	SVAMP(35.6%)	GSM8K(46.9%)	AQuA(9.0%)	M.AVG
GPT-3.5	None	None	87.4	69.9	56.4	38.9	62.3
LLaMA-7B	LoRA	0.83	<b>81.9</b>	48.2	38.3	18.5	43.7
	HypLoRA	0.83	79.2	<b>49.3</b>	<b>39.2</b>	<b>20.7</b>	<b>44.6</b>
LLaMA-13B	LoRA	0.67	<b>83.5</b>	54.7	48.5	18.5	51.0
	HypLoRA	0.67	83.4	<b>55.8</b>	<b>49.7</b>	<b>21.5</b>	<b>52.2</b>
Gemma-7B	LoRA	0.79	89.5	74.2	61.4	34.3	65.9
	HypLoRA	0.79	<b>89.5</b>	<b>75.6</b>	<b>68.8</b>	<b>46.5</b>	<b>71.0</b>
LLaMA3-8B	LoRA	0.70	<b>91.2</b>	<b>81.2</b>	69.8	42.1	73.2
	HypLoRA	0.70	89.1	80.6	<b>72.6</b>	<b>46.9</b>	<b>74.5</b>
Gemma3-4B	LoRA HypLoRA	1.04 1.04	89.0 <b>90.8</b>	<b>80.3</b> 78.9	67.7 <b>67.8</b>	43.5 <b>48.4</b>	71.8 <b>72.0</b>

If we consider token embeddings in a hyperbolic space with polar coordinates  $(r,\theta)$ , where  $r\in\mathbb{R}^+$  is the radial coordinate (correlating with token frequency) and  $\theta\in[0,2\pi]$  is the angular coordinate (encoding semantic similarity), the radial distribution of tokens follows  $p(r)\sim e^{-\zeta r}$ , where  $\zeta>0$  relates to the hyperbolic curvature K. The frequency function k(r) for tokens at radius r is then given by  $k(r)\sim e^{-r}$ . Through coordinate transformation, we can derive the power-law frequency distribution:

$$P(k) \sim P(r) \left| \frac{dr}{dk} \right| \sim k^{-\gamma}.$$

The relationship between the hyperbolic curvature and the power-law exponent  $\gamma$  can be given by  $\gamma=2+\frac{1}{\zeta}$  (as concluded by Krioukov et al. [23] in the context of complex networks, where  $\zeta$  relates to the effective temperature or network structure). This relationship underscores the theoretical connection between the power-law behavior observed in token frequencies and the inherent hyperbolic geometry of the embedding space.

Hyperbolic space offers distinct advantages for modeling language hierarchies, especially when addressing the structural and spatial constraints of token co-occurrence: (1) Separation of Low-Frequency Tokens. Tokens with low frequencies, which typically represent more specific or granular concepts, require clear separation from each other to maintain semantic clarity. (2) Proximity to High-Frequency Hypernyms. Simultaneously, these low-frequency tokens should remain close to their corresponding high-frequency hypernyms or function words. Hyperbolic space is uniquely suited for capturing these dual constraints due to its exponential volume growth, which inherently supports hierarchical structure and allows for ample separation of specific entities while keeping them close to their parent categories. This contrasts with Euclidean space, where such arrangements can lead to crowding or distortion of distances.

**Overall Conclusion.** Through these analyses, we demonstrate that token embeddings in LLMs exhibit hierarchical organization and significant hyperbolicity. This understanding not only sheds light on the geometric nature of token embeddings but also motivates the development of methods that can better capture and preserve these underlying geometric properties.

# 5 Hyperbolic Fine-tuning for LLMs

The core technique in the LoRA adapter involves linear transformations. The conventional approach to implementing linear transformations in the Lorentz model of hyperbolic geometry is through operations in the tangent space, while maintaining the learnable weights in Euclidean space [31, 30]. However, this approach presents a significant challenge for our application. Since the hidden states of LLMs exist in Euclidean space, we would need to project these states to hyperbolic space and subsequently map them back to the tangent space. This process results in consecutive logarithmic and exponential mappings  $(\log_{\mathbf{o}}^K(\exp_{\mathbf{o}}^K(\mathbf{x})))$ , which effectively cancel each other out, reducing the method to the original LoRA approach and nullifying any benefits from hyperbolic geometry.

**Direct Lorentz Low-rank Transformation (LLR)**. To overcome this limitation, we propose a Direct Lorentz Low-rank Transformation (LLR) that operates directly on the hyperbolic manifold without

Table 4: Comparison on Commonsense Reasoning Tasks. These datasets contain relatively similar amounts of data, so we use AVG to represent the average accuracy.

Base Model	PEFT Method	# Params (%)	BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	AVG
GPT-3.5	None	None	73.1	85.4	68.5	78.5	66.1	89.8	79.9	74.8	77.0
LLaMA3-8B	LoRA HypLoRA (Ours)	0.70 0.70	70.8 <b>74.1</b>	85.2 <b>87.6</b>	79.9 <b>80.6</b>	91.7 <b>94.5</b>	84.3 <b>84.7</b>	84.2 <b>90.4</b>	71.2 <b>81.2</b>	79.0 <b>85.2</b>	80.8 <b>84.8</b>
Gemma3-4B	LoRA HypLoRA (Ours)	1.04 1.04	68.1 <b>70.0</b>	83.2 <b>84.3</b>	77.2 <b>79.2</b>	88.9 <b>91.5</b>	<b>80.5</b> 80.3	84.5 <b>89.1</b>	69.9 <b>75.9</b>	83.6 <b>86.4</b>	79.5 <b>82.5</b>
Qwen2.5-7B	LoRA HypLoRA (Ours)	0.71 0.71	<b>73.4</b> 72.8	<b>89.5</b> 89.3	79.5 <b>79.8</b>	93.6 <b>94.8</b>	84.1 <b>84.4</b>	92.8 <b>95.5</b>	82.0 <b>87.5</b>	87.0 <b>90.8</b>	85.2 <b>87.0</b>

relying on tangent space mappings. This approach allows us to perform low-rank adaptation while preserving the advantages of hyperbolic geometry:

$$\mathbf{z}^{E} = W_{\text{LoRA}}(\mathbf{x}^{E}) = W\mathbf{x}^{E} + \Delta W\mathbf{x}^{E}$$

$$= W\mathbf{x}^{E} + \Pi_{H \to E}^{K}(\mathbf{LLR}(BA, \Pi_{E \to H}^{K}(\mathbf{x}^{E}))),$$
(8)

where LLR represents our Direct Lorentz Low-Rank Transformation that operates directly on the hyperbolic representation  $\mathbf{x}^H = \Pi_{E \to H}^K(\mathbf{x}^E)$ :

$$\mathbf{LLR}(BA, \mathbf{x}^H) = (\sqrt{\|BA\mathbf{x}_s^H\|_2^2 + K}, BA\mathbf{x}_s^H),\tag{9}$$

where  $\mathbf{x}_s^H$  is the space-like component of  $\mathbf{x}^H$ , i.e.,  $\mathbf{x}_s^H = \mathbf{x}_{[1:n]}^H$  without the first time-like dimension  $\mathbf{x}_{[0:1]}^H$ . The operators  $\Pi_{E \to H}^K$  and  $\Pi_{H \to E}^K$  represent projections from Euclidean to hyperbolic space and vice versa, respectively. While exponential maps and logarithmic maps are commonly used for these projections, we simplify the process by utilizing stereographic projection and its inverse [91] as a more computationally efficient alternative for LLM adaptation, where the detailed formula is given in Appendix D.1. It can be verified that  $\mathbf{LLR}(BA,\mathbf{x}^H) \in \mathcal{L}^n$ , ensuring that our transformation remains within the Lorentz model of hyperbolic space. This transformation primarily affects the space-like dimensions, functioning similarly to a Lorentz rotation. The linear transformation is inspired by hyperbolic neural networks [35, 45, 92]. For efficient integration with LLMs, the transformation removes normalization and non-linear activation terms in [35], varying curvatures in [45], and orthogonal constraints in [92]. Our main contribution lies in applying hyperbolic low-rank adaptation for LLMs, while the specific linear transformation itself is flexible, and other transformations on the manifold could also be compatible with our approach.

By adapting in the hyperbolic domain, HypLoRA captures more complex hierarchical relationships than traditional Euclidean-based methods, as detailed in Proposition 1. Additionally, the low-rank nature of the adaptation matrices A and B promotes parameter efficiency, making HypLoRA well-suited for LLMs.

**Time Complexity.** HypLoRA has similar theoretical time complexity as the Euclidean LoRA, which is  $\mathcal{O}(r \cdot (d+k))$ , where d and k represent the input and output dimensions, respectively. However, in practical implementation, HypLoRA introduces additional computations due to the logarithmic and exponential maps. These additional operations, nevertheless, can be completed within  $\mathcal{O}(N)$  where N is the number of input tokens.

**Proposition 1.** Let  $\mathbf{x} \in \mathbb{R}^d$  denote the input token embeddings. The HypLoRA adaptation, applied to  $\mathbf{x}$ , involves a sequence of projection into hyperbolic space, a Direct Lorentz Low-rank Transformation (LLR), and projection back to Euclidean space. Due to the non-linear nature of these hyperbolic operations, the effective transformation applied by HypLoRA introduces higher-order terms with respect to  $\mathbf{x}$ . As detailed in Appendix  $\mathbf{E}$ , these terms exhibit explicit dependency on the L2 norm,  $\|\mathbf{x}\|_2$ , of the input embeddings. This norm-dependent, higher-order modification enables HypLoRA to capture hierarchical relationships in the embedding space, thereby achieving natural alignment with the underlying hyperbolic geometry of the token representations.

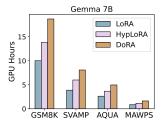
## 5.1 Experimental Settings

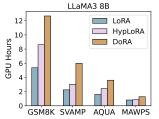
**Datasets.** Following the experiment setup outlined in [93], we utilize two high-quality datasets, Math10K and Commonsense170K, tailored for mathematical and commonsense reasoning, respectively. Math10K consists of training data from GSM8K [81], MAWPS, MAWPS-single [83], and

Table 5: Results for varying curvature K on the Gemma3-4B model

Dataset	K=0.5	K=1.0
MAWPS	90.7	91.9
SVAMP	78.9	78.0
GSM8K	67.8	68.8
AQuA	48.4	41.7
M.AVG	71.9	71.6

Figure 2: GPU (A100) usage during inference





1,000 samples from AQuA [82], augmented with ChatGPT-generated step-by-step rationales to reinforce reasoning capabilities. The test set includes GSM8K, AQuA, MAWPS, and SVAMP [84], ensuring no overlap with the training data. Commonsense170K is constructed by reformatting samples from BoolQ, PIQA, SIQA, HellaSwag, WinoGrande, ARC-e, ARC-c, and OBQA using standardized templates that outline the task, content, and answer, resulting in 170K training samples. The test datasets are drawn from the same sources, with strict separation from training samples. **Model Comparison.** We evaluate the performance of LLaMA-7B, LLaMA-13B, LLaMA3-8B, Gemma-7B, Qwen2.5-7B, and Gemma3-4B base models. Except for Gemma-7B, we use the *instruct* variants of all models to ensure consistency with instruction-tuned behavior. For fine-tuning methods, we compare with LoRA [29].

Implementation Details. To ensure consistency and comparability, our experimental setup closely followed the training configurations outlined in Hu et al. [93]. Across all fine-tuning tasks, we employed the AdamW optimizer with a learning rate of  $3 \times 10^{-4}$  and trained for a total of three epochs. LoRA modules (and consequently, HypLoRA adapters) were integrated into both the Multi-Head Attention (MHA) and MLP layers of the foundation models. A key hyperparameter for HypLoRA is the curvature K (defining the hyperbolic curvature as -1/K), which was determined by searching the set  $\{0.5, 1.0\}$ . For evaluation, final scores were micro-averaged for arithmetic reasoning and averaged for commonsense reasoning across the datasets, thereby giving equal weight to each individual prompt, regardless of the varying number of questions per dataset (e.g., 1,319 in GSM8K versus 238 in MAWPS). All experiments were executed on a single NVIDIA A100 GPU with 80GB of memory.

# 5.2 Experimental Results

Table 3 summarizes our key experimental outcomes on arithmetic reasoning tasks, while Table 4 presents results for commonsense reasoning benchmarks. Our primary comparison contrasts LoRA and HypLoRA to demonstrate the effectiveness of the proposed approach, with additional baselines provided in Appendix F. Our findings consistently demonstrate that HypLoRA achieves superior performance compared to the standard LoRA across a diverse range of LLMs and reasoning tasks. This empirical evidence supports the theoretical advantages outlined in Proposition 1, suggesting that adapting LLMs within the hyperbolic domain allows for a more effective capture of the complex, often hierarchical, relationships.

**Arithmetic Reasoning Performance.** On arithmetic reasoning tasks, as indicated by results in Table 3, HypLoRA shows notable efficacy, especially on datasets recognized for their complexity, such as GSM8K and AQuA. These datasets demand robust multi-step reasoning and a nuanced understanding of numerical and textual relationships. The enhanced performance of HypLoRA in these areas aligns with its design; by operating in hyperbolic space, it can better model the hierarchical structure of problems and distinguish subtle yet critical differences in input embeddings. This is further corroborated by the theoretical analysis (Appendix E), which posits that HypLoRA introduces higher-order, norm-dependent terms. These terms allow the model to develop a more refined sensitivity to token importance and inter-token relationships.

Commonsense Reasoning Performance. The robust performance of HypLoRA extends to commonsense reasoning, as detailed in Table 4. For the Gemma3-4B model, HypLoRA achieved an average accuracy of 82.5% across all datasets, surpassing LoRA's 79.5%. Similarly, on the Qwen2.5-7B model, HypLoRA obtained an average of 87.0% compared to LoRA's 85.2%. These improvements are distributed across various commonsense benchmarks, including notable gains on datasets like ARC-c and OBQA for Gemma3-4B, and ARC-c, ARC-e, and OBQA for Qwen2.5-7B. Commonsense

reasoning often relies on understanding implicit relationships and contextual nuances, which may not always be explicitly hierarchical but still benefit from the richer representational capacity offered by hyperbolic geometry. The ability of HypLoRA to better discern these subtleties, likely due to the mechanisms described in Proposition 1, contributes to these observed performance gains, showcasing the broad applicability of hyperbolic fine-tuning.

The Impact of Curvature on Performance. Curvature in hyperbolic space is a key hyperparameter in HypLoRA, directly affecting its capacity to model underlying structures and geometries. To evaluate its impact, we experiment with a learnable curvature initialised with different curvature values on the Gemma3-4B model, as shown in Table 5, where the curvature is defined as -1/K. Our results demonstrate that curvature does influence model performance. For Gemma-7B and Gemma-3-4B, a curvature value of 0.5 consistently yields the best overall performance across both arithmetic and commonsense reasoning benchmarks. Similarly, for LLaMa3-8B, 0.5 proves optimal. In commonsense reasoning benchmarks, a curvature of 1.0 performs best for LLaMa3-8B and Qwen2.5-7B.

**Efficiency.** In Section 5, we analyze the time complexity of our approach, which remains consistent with that of LoRA. However, during actual inference, HypLoRA incurs additional computational overhead due to operations such as the exponential and logarithmic mappings, or the inverse stereographic and stereographic projections when using the stereographic method. These operations introduce some additional runtime, particularly for larger models. The GPU hours for inference on four datasets are presented in Figure 2. Despite this overhead, our method demonstrates improved efficiency when compared to the previous competitive model, DoRA. Notably, HypLoRA still outperforms DoRA in terms of both runtime and overall efficiency. Besides, all models can be fine-tuned in approximately one hour for optimal training efficiency.

#### 6 Conclusion

In this work, we investigated the non-Euclidean geometric properties inherent in LLM token embeddings, confirming their strong hyperbolic characteristics, which suggest underlying hierarchical structures. Building on these insights, we introduced HypLoRA, a novel hyperbolic low-rank adaptation technique. HypLoRA performs fine-tuning directly on the hyperbolic manifold. Extensive experiments show that HypLoRA significantly improves LLM performance on arithmetic reasoning and commonsense tasks. By leveraging the hyperbolic structure of the data, HypLoRA enhances the model's ability to capture and utilize intricate relationships, leading to better reasoning capabilities. While the theoretical time complexity aligns with LoRA, the practical application of hyperbolic projections introduces a slight computational overhead. This is a manageable aspect that can be readily addressed through optimized numerical libraries or by exploring alternative, more computationally efficient projection techniques without sacrificing the geometric benefits.

**Broader Impact.** Enhancing reasoning-oriented LLMs can help education, scientific assistance, and safer decision-support systems, but the same improvements may also accelerate misuse (e.g., automating complex disinformation or amplifying biased advice) and increase energy consumption because of added hyperbolic projections. We therefore advocate releasing checkpoints and code with usage guidelines (as in our public repo), tracking compute budgets when scaling HypLoRA further, and pairing capability gains with evaluations focused on safety, robustness, and fairness.

# Acknowledgment

The authors would like to express their sincere gratitude to the anonymous reviewers and the Area Chair for their valuable comments and insightful suggestions, which have greatly improved this work. We thank our colleagues and lab mates at HKUST(GZ), HKUST, IISc, Yale, Stanford, and CUHK for motivating discussions on geometric learning for LLMs, and we acknowledge the administrators of the institutional GPU clusters that enabled the large-scale experiments. Part of this research was conducted while Menglin Yang was at Yale University, whose support and computing resource were instrumental in shaping the initial ideas of this project.

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 Technical Report. arXiv preprint arXiv:2303.08774, 2023.
- [2] 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.
- [3] Google Deepmind Gemma Team. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- [4] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [5] Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. Is chatgpt a general-purpose natural language processing task solver? *arXiv* preprint *arXiv*:2302.06476, 2023.
- [6] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. HuggingGPT: Solving AI tasks with ChatGPT and its friends in Hugging Face. In *Advances in Neural Information Processing Systems*, volume 36, 2024.
- [7] Jiahong Liu, Zexuan Qiu, Zhongyang Li, Quanyu Dai, Wenhao Yu, Jieming Zhu, Minda Hu, Menglin Yang, Tat-Seng Chua, and Irwin King. A survey of personalized large language models: Progress and future directions. *arXiv preprint arXiv:2502.11528*, 2025.
- [8] Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.
- [9] Gregor Bachmann, Gary Bécigneul, and Octavian Ganea. Constant curvature graph convolutional networks. In *ICML*, pages 486–496. PMLR, 2020.
- [10] Atsushi Suzuki, Atsushi Nitanda, Jing Wang, Linchuan Xu, Kenji Yamanishi, and Marc Cavazza. Generalization error bound for hyperbolic ordinal embedding. In *ICML*, pages 10011–10021. PMLR, 2021.
- [11] Maximillian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. In *Advances in Neural Information Processing Systems*, pages 6338–6347, 2017.
- [12] Maximillian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz model of hyperbolic geometry. In *ICML*, pages 3779–3788, 2018.
- [13] Octavian Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic entailment cones for learning hierarchical embeddings. In *ICML*, pages 1646–1655. PMLR, 2018.
- [14] Valentin Khrulkov, Leyla Mirvakhabova, Evgeniya Ustinova, Ivan Oseledets, and Victor Lempitsky. Hyperbolic image embeddings. In CVPR, pages 6418–6428, 2020.
- [15] Edoardo Cetin, Benjamin Chamberlain, Michael Bronstein, and Jonathan J Hunt. Hyperbolic deep reinforcement learning. *arXiv preprint arXiv:2210.01542*, 2022.
- [16] Wei Peng, Tuomas Varanka, Abdelrahman Mostafa, Henglin Shi, and Guoying Zhao. Hyperbolic deep neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [17] Menglin Yang, Min Zhou, Zhihao Li, Jiahong Liu, Lujia Pan, Hui Xiong, and Irwin King. Hyperbolic graph neural networks: A review of methods and applications. *arXiv preprint arXiv:2202.13852*, 2022.

- [18] Jiahong Liu, Xinyu Fu, Menglin Yang, Weixi Zhang, Rex Ying, and Irwin King. Client-specific hyperbolic federated learning. In KDD, 2024.
- [19] Pascal Mettes, Mina Ghadimi Atigh, Martin Keller-Ressel, Jeffrey Gu, and Serena Yeung. Hyperbolic deep learning in computer vision: A survey. *arXiv preprint arXiv:2305.06611*, 2023.
- [20] Menglin Yang, Zhihao Li, Min Zhou, Jiahong Liu, and Irwin King. Hicf: Hyperbolic informative collaborative filtering. In *KDD*, pages 2212–2221, 2022.
- [21] Michele Borassi, Alessandro Chessa, and Guido Caldarelli. Hyperbolicity measures democracy in real-world networks. *Physical Review E*, 92(3):032812, 2015.
- [22] W Sean Kennedy, Onuttom Narayan, and Iraj Saniee. On the hyperbolicity of large-scale networks. *arXiv preprint arXiv:1307.0031*, 2013.
- [23] Dmitri Krioukov, Fragkiskos Papadopoulos, Maksim Kitsak, Amin Vahdat, and Marián Boguná. Hyperbolic geometry of complex networks. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 82(3):036106, 2010.
- [24] Albert-László Barabási, Zoltán Dezső, Erzsébet Ravasz, Soon-Hyung Yook, and Zoltán Oltvai. Scale-free and hierarchical structures in complex networks. In *AIP Conference Proceedings*, volume 661, pages 1–16. American Institute of Physics, 2003.
- [25] Kai Nakaishi, Ryo Yoshida, Kohei Kajikawa, Koji Hukushima, and Yohei Oseki. Rethinking the relationship between the power law and hierarchical structures. *arXiv preprint arXiv:2505.04984*, 2025.
- [26] Enrique Alvarez-Lacalle, Beate Dorow, J-P Eckmann, and Elisha Moses. Hierarchical structures induce long-range dynamical correlations in written texts. *Proceedings of the National Academy of Sciences*, 103(21):7956–7961, 2006.
- [27] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint* arXiv:1711.05101, 2017.
- [28] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In SC20: International Conference for High Performance Computing, Networking, Storage and Analysis, pages 1–16. IEEE, 2020.
- [29] 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. *arXiv* preprint arXiv:2106.09685, 2021.
- [30] Ines Chami, Zhitao Ying, Christopher Ré, and Jure Leskovec. Hyperbolic graph convolutional neural networks. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [31] Octavian Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic neural networks. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [32] Menglin Yang, Min Zhou, Hui Xiong, and Irwin King. Hyperbolic temporal network embedding. *TKDE*, 2022.
- [33] Sungjun Cho, Seunghyuk Cho, Sungwoo Park, Hankook Lee, Honglak Lee, and Moontae Lee. Curve your attention: Mixed-curvature transformers for graph representation learning. *arXiv* preprint arXiv:2309.04082, 2023.
- [34] Xingcheng Fu, Yisen Gao, Yuecen Wei, Qingyun Sun, Hao Peng, Jianxin Li, and Xianxian Li. Hyperbolic geometric latent diffusion model for graph generation. *arXiv preprint* arXiv:2405.03188, 2024.
- [35] Weize Chen, Xu Han, Yankai Lin, Hexu Zhao, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. Fully hyperbolic neural networks. *arXiv preprint arXiv:2105.14686*, 2021.
- [36] Ryohei Shimizu, Yusuke Mukuta, and Tatsuya Harada. Hyperbolic neural networks++. *arXiv* preprint arXiv:2006.08210, 2020.

- [37] Xiaomeng Fan, Yuwei Wu, Zhi Gao, Mehrtash Harandi, and Yunde Jia. Curvature learning for generalization of hyperbolic neural networks: X. fan et al. *IJCV*, pages 1–37, 2025.
- [38] Yidan Mao, Jing Gu, Marcus C Werner, and Dongmian Zou. Klein model for hyperbolic neural networks. *arXiv preprint arXiv:2410.16813*, 2024.
- [39] Ahmad Bdeir, Kristian Schwethelm, and Niels Landwehr. Hyperbolic geometry in computer vision: A novel framework for convolutional neural networks. arXiv preprint arXiv:2303.15919, 2023.
- [40] Max van Spengler, Erwin Berkhout, and Pascal Mettes. Poincaré resnet. In *ICCV*, pages 5419–5428, 2023.
- [41] Raiyan R Khan, Philippe Chlenski, and Itsik Pe'er. Hyperbolic genome embeddings. *arXiv* preprint arXiv:2507.21648, 2025.
- [42] Qi Liu, Maximilian Nickel, and Douwe Kiela. Hyperbolic graph neural networks. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [43] Menglin Yang, Min Zhou, Lujia Pan, and Irwin King.  $\kappa$ hgcn: Tree-likeness modeling via continuous and discrete curvature learning. In *KDD*, pages 2965–2977, 2023.
- [44] Caglar Gulcehre, Misha Denil, Mateusz Malinowski, Ali Razavi, Razvan Pascanu, Karl Moritz Hermann, Peter Battaglia, Victor Bapst, David Raposo, Adam Santoro, et al. Hyperbolic attention networks. *arXiv preprint arXiv:1805.09786*, 2018.
- [45] Menglin Yang, Harshit Verma, Delvin Ce Zhang, Jiahong Liu, Irwin King, and Rex Ying. Hypformer: Exploring efficient transformer fully in hyperbolic space. *arXiv preprint arXiv:2407.01290*, 2024.
- [46] Jianing Sun, Zhaoyue Cheng, Saba Zuberi, Felipe Pérez, and Maksims Volkovs. HGCF: Hyperbolic graph convolution networks for collaborative filtering. In *WWW*, pages 593–601, 2021.
- [47] Zhenzhen Weng, Mehmet Giray Ogut, Shai Limonchik, and Serena Yeung. Unsupervised discovery of the long-tail in instance segmentation using hierarchical self-supervision. In *CVPR*, pages 2603–2612, 2021.
- [48] Bo Xiong, Michael Cochez, Mojtaba Nayyeri, and Steffen Staab. Hyperbolic embedding inference for structured multi-label prediction. In *Advances in Neural Information Processing Systems*, volume 35, pages 33016–33028, 2022.
- [49] Menglin Yang, Min Zhou, Marcus Kalander, Zengfeng Huang, and Irwin King. Discrete-time temporal network embedding via implicit hierarchical learning in hyperbolic space. In *KDD*, pages 1975–1985, 2021.
- [50] Zhi Gao, Yuwei Wu, Yunde Jia, and Mehrtash Harandi. Curvature generation in curved spaces for few-shot learning. In *ICCV*, pages 8691–8700, 2021.
- [51] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- [52] Karan Desai, Maximilian Nickel, Tanmay Rajpurohit, Justin Johnson, and Ramakrishna Vedantam. Hyperbolic image-text representations. In *ICML*, 2023.
- [53] Avik Pal, Max van Spengler, Guido Maria D'Amely di Melendugno, Alessandro Flaborea, Fabio Galasso, and Pascal Mettes. Compositional entailment learning for hyperbolic vision-language models. *arXiv preprint arXiv:2410.06912*, 2024.
- [54] Tobia Poppi, Tejaswi Kasarla, Pascal Mettes, Lorenzo Baraldi, and Rita Cucchiara. Hyperbolic safety-aware vision-language models. In *CVPR*, 2025.

- [55] Neil He, Menglin Yang, and Rex Ying. Hypercore: The core framework for building hyperbolic foundation models with comprehensive modules. In *ICLR*, 2025.
- [56] Max Kochurov, Rasul Karimov, and Serge Kozlukov. Geoopt: Riemannian optimization in pytorch. *arXiv preprint arXiv:2005.02819*, 2020.
- [57] Steven Thomas Smith. Optimization techniques on riemannian manifolds. arXiv preprint arXiv:1407.5965, 2014.
- [58] Gary Bécigneul and Octavian-Eugen Ganea. Riemannian adaptive optimization methods. *arXiv preprint arXiv:1810.00760*, 2018.
- [59] Emily Reif, Ann Yuan, Martin Wattenberg, Fernanda B Viegas, Andy Coenen, Adam Pearce, and Been Kim. Visualizing and measuring the geometry of bert. In Advances in Neural Information Processing Systems, volume 32, 2019.
- [60] Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tie-Yan Liu. Representation degeneration problem in training natural language generation models. arXiv preprint arXiv:1907.12009, 2019.
- [61] William Rudman, Nate Gillman, Taylor Rayne, and Carsten Eickhoff. Isoscore: Measuring the uniformity of embedding space utilization. *arXiv preprint arXiv:2108.07344*, 2021.
- [62] Giovanni Puccetti, Anna Rogers, Aleksandr Drozd, and Felice Dell'Orletta. Outliers dimensions that disrupt transformers are driven by frequency. arXiv preprint arXiv:2205.11380, 2022.
- [63] OpenAI Foundation. Introducing chatgpt. https://openai.com/index/chatgpt, November 2022.
- [64] OpenAI Foundation. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- [65] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv* preprint arXiv:2104.08691, 2021.
- [66] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv* preprint arXiv:2101.00190, 2021.
- [67] Yujia Qin, Xiaozhi Wang, Yusheng Su, Yankai Lin, Ning Ding, Jing Yi, Weize Chen, Zhiyuan Liu, Juanzi Li, Lei Hou, et al. Exploring universal intrinsic task subspace via prompt tuning. arXiv preprint arXiv:2110.07867, 2021.
- [68] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *ICML*, pages 2790–2799. PMLR, 2019.
- [69] Yaoming Zhu, Jiangtao Feng, Chengqi Zhao, Mingxuan Wang, and Lei Li. Counter-interference adapter for multilingual machine translation. arXiv preprint arXiv:2104.08154, 2021.
- [70] Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. *arXiv preprint arXiv:2012.13255*, 2020.
- [71] Ali Edalati, Marzieh Tahaei, Ivan Kobyzev, Vahid Partovi Nia, James J Clark, and Mehdi Rezagholizadeh. Krona: Parameter efficient tuning with kronecker adapter. arXiv preprint arXiv:2212.10650, 2022.
- [72] Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. *arXiv preprint arXiv:2402.12354*, 2024.
- [73] Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv* preprint arXiv:2402.09353, 2024.

- [74] Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *ICLR*, 2023.
- [75] Xi Wang, Laurence Aitchison, and Maja Rudolph. Lora ensembles for large language model fine-tuning. *arXiv preprint arXiv:2310.00035*, 2023.
- [76] Pengjie Ren, Chengshun Shi, Shiguang Wu, Mengqi Zhang, Zhaochun Ren, Maarten de Rijke, Zhumin Chen, and Jiahuan Pei. Mini-ensemble low-rank adapters for parameter-efficient fine-tuning. *arXiv preprint arXiv:2402.17263*, 2024.
- [77] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. In Advances in Neural Information Processing Systems, volume 36, 2024.
- [78] Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhensu Chen, Xiaopeng Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language models. *arXiv preprint arXiv:2309.14717*, 2023.
- [79] Yixiao Li, Yifan Yu, Chen Liang, Pengcheng He, Nikos Karampatziakis, Weizhu Chen, and Tuo Zhao. Loftq: Lora-fine-tuning-aware quantization for large language models. *arXiv* preprint arXiv:2310.08659, 2023.
- [80] Gal Mishne, Zhengchao Wan, Yusu Wang, and Sheng Yang. The numerical stability of hyperbolic representation learning. In *ICML*, pages 24925–24949. PMLR, 2023.
- [81] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [82] Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. *arXiv preprint arXiv:1705.04146*, 2017.
- [83] Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. Mawps: A math word problem repository. In *NAACL*, pages 1152–1157, 2016.
- [84] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math word problems? *arXiv preprint arXiv:2103.07191*, 2021.
- [85] Jeff Alstott, Ed Bullmore, and Dietmar Plenz. Powerlaw: a python package for analysis of heavy-tailed distributions. *PloS one*, 9(1):e85777, 2014.
- [86] Erzsébet Ravasz and Albert-László Barabási. Hierarchical organization in complex networks. *Physical review E*, 67(2):026112, 2003.
- [87] Mikhael Gromov. Hyperbolic groups. In *Essays in group theory*, pages 75–263. Springer, 1987.
- [88] Hervé Fournier, Anas Ismail, and Antoine Vigneron. Computing the gromov hyperbolicity of a discrete metric space. *Information Processing Letters*, 115(6-8):576–579, 2015.
- [89] Dmitri Krioukov, Fragkiskos Papadopoulos, Amin Vahdat, and Marián Boguná. Curvature and temperature of complex networks. *Physical Review E*, 80(3):035101, 2009.
- [90] Fragkiskos Papadopoulos, Dmitri Krioukov, Marián Boguná, and Amin Vahdat. Greedy forwarding in dynamic scale-free networks embedded in hyperbolic metric spaces. In 2010 Proceedings IEEE Infocom, pages 1–9. IEEE, 2010.
- [91] Ondrej Skopek, Octavian-Eugen Ganea, and Gary Bécigneul. Mixed-curvature variational autoencoders. *arXiv preprint arXiv:1911.08411*, 2019.
- [92] Jindou Dai, Yuwei Wu, Zhi Gao, and Yunde Jia. A hyperbolic-to-hyperbolic graph convolutional network. In *CVPR*, pages 154–163, 2021.

- [93] Zhiqiang Hu, Yihuai Lan, Lei Wang, Wanyu Xu, Ee-Peng Lim, Roy Ka-Wei Lee, Lidong Bing, and Soujanya Poria. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. *arXiv preprint arXiv:2304.01933*, 2023.
- [94] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- [95] Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics, and function using networkx. In Gäel Varoquaux, Travis Vaught, and Jarrod Millman, editors, Proceedings of the 7th Python in Science Conference (SciPy2008), pages 11–15. Pasadena, CA USA, 2008.
- [96] Philip McCord Morse and Herman Feshbach. Methods of theoretical physics. Technology Press, 1946.
- [97] Valter Moretti. The interplay of the polar decomposition theorem and the lorentz group. *arXiv* preprint math-ph/0211047, 2002.
- [98] Ines Chami, Albert Gu, Dat Nguyen, and Christopher Ré. HoroPCA: Hyperbolic dimensionality reduction via horospherical projections. In *ICML*, volume 139 of *Proceedings of Machine Learning Research*, pages 1419–1429, 2021.
- [99] P. Thomas Fletcher, Conglin Lu, Stephen M. Pizer, and Sarang Joshi. Principal geodesic analysis for the study of nonlinear statistics of shape. *IEEE Transactions on Medical Imaging*, 23(8):995–1005, 2004.
- [100] Stefan Sommer, François Lauze, and Mads Nielsen. Optimization over geodesics for exact principal geodesic analysis. *Advances in Computational Mathematics*, 40:283–313, 2014.
- [101] Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. *arXiv* preprint arXiv:2410.05229, 2024.

# **NeurIPS Paper Checklist**

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The paper's abstract and introduction explicitly state the main claims, including the contribution of the proposed method and analysis, and its scope.

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The paper discusses a limitation in the Conclusion (Section 6).

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

# 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: The paper presents Proposition 1 regarding HypLoRA, introducing higherorder terms. It states that the details and derivation of these terms are provided in Appendix F.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

## 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Section 5.1 "Experimental Settings" describes the datasets, models, and implementation details, including optimizer (AdamW),etc.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

# 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The datasets used (e.g., GSM8K, AQUA, various commonsense reasoning benchmarks) are publicly available and cited.

#### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
  to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

# 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Section 5.1 details the experimental setup.

# Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

#### 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Table 2, which presents hyperbolicity across various metric spaces and datasets, includes error bars reported as mean and standard deviation. Tables 3 and 4 report mean accuracies over three runs.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

# 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The paper states in Section 5.1 that "All experiments were executed on a single NVIDIA A100 GPU with 80GB of memory."

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <a href="https://neurips.cc/public/EthicsGuidelines">https://neurips.cc/public/EthicsGuidelines</a>?

Answer: [Yes]

Justification: Based on the content of the paper, the research focuses on investigating geometric properties of LLMs and proposing a new fine-tuning method. The datasets used are standard benchmarks in the field. There is no indication of activities that would violate the NeurIPS Code of Ethics, such as plagiarism, falsification of data, or unethical use of human subjects.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

### 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Conclusion part.

#### Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: It does not introduce new large-scale pretrained models or novel scraped datasets that would pose a high risk for misuse requiring specific safeguards beyond those applicable to the original models it builds upon.

## Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
  necessary safeguards to allow for controlled use of the model, for example by requiring
  that users adhere to usage guidelines or restrictions to access the model or implementing
  safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The creators/original owners of assets (datasets like GSM8K, MAWPS, AQUA; models like LLaMA, Gemma; and software like the Powerlaw Package ) are credited via citations.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.

- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No new assets.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

# 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The datasets are existing benchmarks.

# Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The research does not appear to involve human subjects in a way that would necessitate IRB approval, as per the justification for question 14.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

## 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [No]

Justification: No, except for writing, editing, or formatting purposes.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.