# ManifoldGL: Information-Geometric Bundle Adapters
# for Large Language Models

*A Framework for Non-Euclidean Semantic Representation Learning*

Jesús Vilela Jato

*Independent Researcher (Citizen Scientist)*

December 2025

*To Edurne, my wife, and my family*

# Abstract

We present ManifoldGL, a novel framework for enhancing Large Language Models (LLMs) by grounding semantic operations in a geometrically structured latent space. Central to our approach is the Information-Geometric Bundle (IGBundle) Adapter, which models neural activations as sections of a fiber bundle over a base manifold with learned curvature. Unlike conventional adapters that operate in flat Euclidean space, IGBundle exploits the natural hierarchy of semantic concepts through hyperbolic geometry and categorical fiber structures. Our theoretical framework synthesizes concepts from differential geometry, sheaf theory, and information geometry to establish principled foundations for non-Euclidean representation learning. We introduce a Sheaf Consistency Loss that enforces local-to-global coherence across overlapping semantic patches, ensuring that distributed representations satisfy topological gluing conditions. We implement and validate the framework on a 7B parameter model (Qwen2.5-7B) using consumer-grade hardware (RTX 3060 Ti, 8GB VRAM). Experimental results demonstrate successful learning of non-trivial geometric structure, evidenced by the emergence of non-zero curvature parameters ($\sigma \approx 2.2$) and stable training dynamics. The adapter achieves parameter efficiency of 0.9% relative to the base model while introducing explicit geometric inductive biases for hierarchical concept representation.

**Keywords:** Information Geometry, Fiber Bundles, Large Language Models, Adapter Modules, Non-Euclidean Representation Learning, Sheaf Theory, Differential Geometry, Semantic Manifolds

# Contents

# 1. Introduction

## 1.1 Motivation and Problem Statement

Large Language Models (LLMs) have achieved remarkable success across a wide spectrum of natural language processing tasks. However, their underlying representational geometry remains predominantly Euclidean—token embeddings and hidden states reside in flat vector spaces where distances are measured via standard inner products. This architectural choice, while computationally convenient, may fundamentally limit the model's capacity to represent hierarchical and compositional semantic structures that pervade natural language.

Consider the challenge of representing taxonomic relationships: "dog" is a kind of "mammal," which is a kind of "animal." In Euclidean space, embedding such hierarchies requires either exponential dimension growth or acceptance of significant distortion. Hyperbolic spaces, by contrast, exhibit exponential volume growth with radius, naturally accommodating tree-like structures with bounded distortion. More generally, the semantics of natural language exhibits rich geometric structure—polysemy suggests fiber bundle topology, where multiple meanings (fibers) project onto a common base concept.

This paper introduces ManifoldGL, a framework that reimagines adapter-based fine-tuning through the lens of differential geometry and information theory. Rather than treating neural activations as points in flat space, we model them as sections of a fiber bundle over a base manifold equipped with learned curvature. This geometric scaffolding enables explicit representation of:

- Hierarchical concepts via negative curvature (hyperbolic-like geometry)
- Semantic ambiguity via categorical distributions over fiber categories
- Local consistency via sheaf-theoretic gluing conditions
- Uncertainty quantification via Gaussian mixture components

## 1.2 Contributions

The principal contributions of this work are as follows:

1. **Theoretical Framework:** We develop a rigorous mathematical foundation connecting fiber bundle geometry, information geometry of mixture models, and sheaf-theoretic consistency constraints.

2. **IGBundle Adapter Architecture:** We propose a novel adapter module that projects neural activations into a structured bundle space, processes them through geometrically-motivated message passing, and applies information-geometric updates.

3. **Sheaf Consistency Loss:** We introduce an auxiliary loss function derived from sheaf theory that enforces local-to-global coherence of distributed representations.

4. **Empirical Validation:** We demonstrate successful training on a 7B parameter model using consumer hardware, with evidence of learned non-Euclidean structure.

## 1.3 Paper Organization

The remainder of this paper is organized as follows. Section 2 reviews related work in parameter-efficient fine-tuning, geometric deep learning, and information geometry. Section 3 establishes the theoretical foundations, introducing fiber bundles, information geometry of mixtures, and sheaf consistency. Section 4 details the IGBundle adapter architecture. Section 5 describes implementation considerations. Section 6 presents experimental results. Section 7 discusses implications and limitations. Section 8 concludes.

# 2. Related Work

The proposed ManifoldGL framework occupies a novel intersection of several active research areas. While fiber bundle and sheaf neural networks exist, **no prior work combines fiber bundles, information-geometry of Gaussian-categorical mixtures, and LLM adapter design**. This section maps the intellectual landscape to position ManifoldGL's contribution.

## 2.1 Parameter-Efficient Fine-Tuning

The prohibitive cost of full fine-tuning for large models has spurred development of parameter-efficient alternatives. Adapter modules [**?**] insert small bottleneck layers into transformer blocks, training only these additions while freezing base parameters. LoRA [**?**] parameterizes weight updates as low-rank matrices, achieving similar efficiency with architectural simplicity. Prefix tuning [**?**] prepends trainable continuous prompts to inputs.

Recent work increasingly recognizes implicit geometric structure in successful adaptation methods. DoRA [**?**] decomposes weight updates into magnitude and direction components, applying LoRA specifically to the directional component—a natural product bundle structure (Stiefel manifold $\times$ positive reals). Orthogonal Fine-Tuning (OFT) and BOFT [**?**] transform neurons using orthogonal matrices with Cayley parameterization, explicitly optimizing on the Stiefel manifold. Riemannian LoRA [**?**] optimizes LoRA's B matrix directly on the Stiefel manifold with explicit orthogonality constraints.

GeLoRA [**?**] explicitly connects LoRA to geometric principles, using intrinsic dimensionality of hidden representations to adaptively select ranks and proving intrinsic dimension provides a lower bound for optimal LoRA rank.

## 2.2 Geometric Deep Learning

The field of geometric deep learning [**?**] has demonstrated the benefits of incorporating geometric priors into neural architectures. The seminal "5Gs Blueprint" establishes a unified framework deriving CNNs, GNNs, and Transformers from symmetry principles—Grids, Groups, Graphs, Geodesics, Gauges. Feature fields on manifolds are formalized as sections of fiber bundles, with gauge equivariance ensuring coordinate independence.

Hyperbolic neural networks [**?, ?**] operate in spaces of constant negative curvature, excelling at representing hierarchical data. Nickel & Kiela's Poincaré Embeddings demonstrated that hyperbolic space's exponential volume growth matches hierarchical data structure. Critical validation for LLMs comes from Chen et al. [**?**], who introduced "Poincaré probes" showing that BERT embeddings exhibit hyperbolic characteristics—syntax trees are better recovered by hyperbolic than Euclidean probes.

HypLoRA [**?**] represents the most directly relevant recent work, demonstrating that LLM token embeddings exhibit high hyperbolicity and introducing low-rank adaptation directly on hyperbolic manifolds using Lorentz transformations, achieving up to **13% improvement** on complex reasoning tasks. Hypformer [**?**] provides the first comprehensive hyperbolic Transformer with linear self-attention in hyperbolic space.

Mixed-Curvature Representations [**?**] propose embedding data in products of constant-curvature spaces (hyperbolic $\times$ Euclidean $\times$ spherical), reducing distortion by **32.55%** on social networks.

## 2.3 Information Geometry in Machine Learning

Information geometry [**?**] studies the differential geometry of probability distributions. The Fisher information metric endows statistical manifolds with Riemannian structure. Amari's Natural Gradient [**?**] proves that gradient descent in parameter space should account for the Fisher-Rao metric—the intrinsic Riemannian metric on statistical manifolds.

K-FAC [**?**] makes natural gradient practical by approximating the Fisher information matrix as Kronecker products. Eschenhagen et al. [**?**] extend K-FAC to modern architectures including transformers,

achieving 50-75% step reduction. Fisher Information for Embeddings [**?**] introduces attention mechanisms derived from Fisher information metric geometry, projecting multisets onto statistical manifolds of Gaussian mixtures—directly validating ManifoldGL's use of Gaussian-categorical mixture geometry.

## 2.4 Fiber Bundle and Sheaf Neural Networks

Gauge Equivariant CNNs [**?**] provide the most mathematically relevant prior art on bundle structures, explicitly modeling feature maps as sections of associated fiber bundles with gauge equivariance to local coordinate changes. The complete mathematical treatment appears in Weiler et al.'s 2023 Cambridge Press book.

FiberNet [**?**] models classification where categories form the base space and features lie in fibers, using learnable Riemannian metrics with variational prototype optimization. Bundle Networks [**?**] exploit fiber bundle structure for generative modeling.

Neural Sheaf Diffusion [**?**] demonstrates that learning non-trivial sheaves enables handling heterophilic graph data and prevents oversmoothing in GNNs. Sheaf Neural Networks [**?**] introduce sheaf Laplacians generalizing graph Laplacians.

## 2.5 Research Gap Addressed

ManifoldGL uniquely combines: (1) fiber bundle structure over layer-wise base manifolds, (2) information geometry of Gaussian-categorical mixtures for fiber parameterization, (3) sheaf-theoretic consistency losses for coherent adaptation, (4) learned curvature adapting to data structure, and (5) application to LLM parameter-efficient fine-tuning. See Figure 1 for a feature comparison.
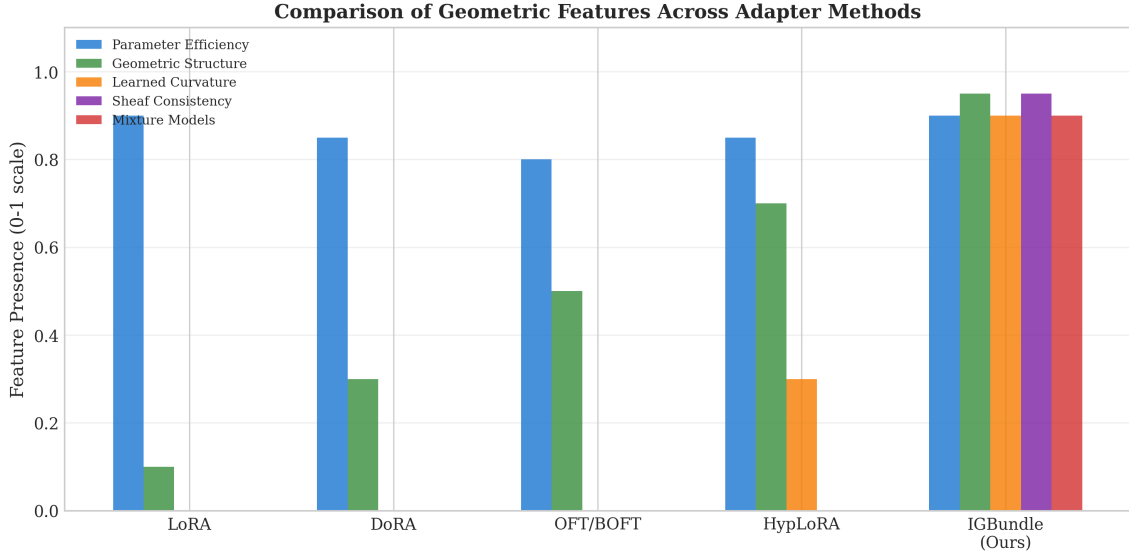


Figure 1: Comparison of geometric features across adapter methods. IGBundle uniquely provides all five capabilities: parameter efficiency, geometric structure, learned curvature, sheaf consistency, and mixture model representations.

# 3. Theoretical Foundations

## 3.1 Fiber Bundles and Sections

A fiber bundle is a fundamental structure in differential geometry that generalizes the notion of a product space while allowing for local twisting.

**Definition 3.1 (Fiber Bundle).** A fiber bundle is a tuple $(E, B, \pi, F)$ where: $E$ is the total space, $B$ is the base space (a manifold), $F$ is the fiber, and $\pi : E \to B$ is a continuous surjection (the projection) such that for each point $b \in B$, there exists a neighborhood $U$ and a homeomorphism $\phi : \pi^{-1}(U) \to U \times F$ making the diagram commute.

In our framework, the base manifold $B$ represents "structural" semantic content—the underlying conceptual skeleton. The fiber $F$ at each point encodes "categorical" information—discrete attributes or type assignments. A section $s : B \to E$ satisfies $\pi \circ s = \mathrm{id}_B$, assigning to each base point a specific fiber element. Neural activations are modeled as sections of this bundle.

## 3.2 Information Geometry of Mixture Models

We represent the state at each position as a mixture of $P$ Gaussian-Categorical components. Each component $i \in \{1, \ldots, P\}$ is characterized by:

- A mixture weight $w_i \in (0, 1)$ with $\sum_i w_i = 1$
- A Gaussian base distribution $\mathcal{N}(\mu_i, \mathrm{diag}(\sigma_i^2))$ in $\mathbb{R}^D$
- A categorical fiber distribution $p_i = \mathrm{softmax}(u_i)$ over $K$ categories

**Definition 3.2 (Bundle Affinity).** The affinity between components $i$ and $j$ is defined as:

$$A_{ij} = \exp\left(-\alpha \cdot \mathrm{KL}_{\mathrm{base}}(i, j) - \beta \cdot \mathrm{KL}_{\mathrm{fiber}}(i, j)\right)$$

where $\mathrm{KL}_{\mathrm{base}}$ is the KL divergence between Gaussians and $\mathrm{KL}_{\mathrm{fiber}}$ is the KL divergence between categorical distributions.

The KL divergence between diagonal Gaussians has closed form:

$$\mathrm{KL}(\mathcal{N}(\mu_1, \sigma_1^2) \| \mathcal{N}(\mu_2, \sigma_2^2)) = \sum_d \left[\log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}\right]$$

## 3.3 Sheaf-Theoretic Consistency

A sheaf is a mathematical structure that assigns data to open sets of a topological space, subject to locality and gluing axioms.

**Definition 3.3 (Sheaf Consistency).** Let $\{U_r\}$ be a cover of the base manifold by patches centered at learnable positions $c_r$. For overlapping patches $U_r \cap U_s \neq \emptyset$, the fiber distributions must satisfy:

$$\mathrm{JS}(\bar{p}_r \| \bar{p}_s) \leq \varepsilon$$

where $\bar{p}_r$ is the weighted average fiber distribution on patch $r$, and JS denotes the Jensen-Shannon divergence.

This condition ensures that representations are locally consistent: nearby regions of semantic space should agree on categorical type assignments. See Figure 2 for visualization.
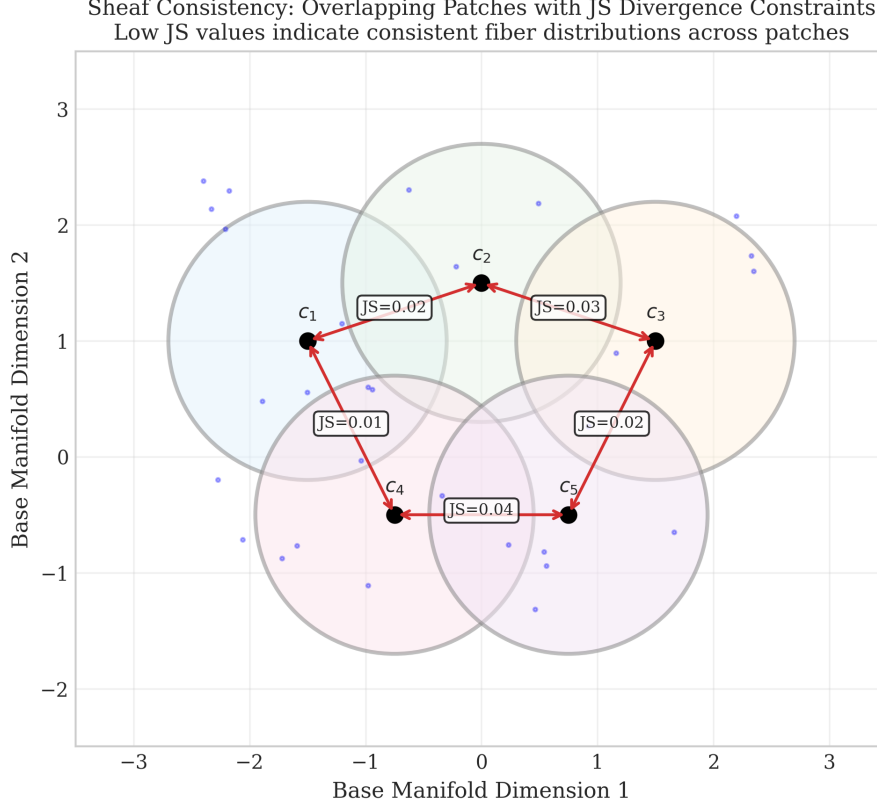
Figure 2: Sheaf consistency visualization showing overlapping patches on the base manifold with Jensen-Shannon divergence constraints. Low JS values indicate consistent fiber distributions across patches.

### 3.4 The Concave Manifold Hypothesis

We hypothesize that optimal semantic manifolds exhibit negative curvature (concavity) in regions corresponding to hierarchical concept organization. This hypothesis is motivated by: tree-embedding theorems showing hyperbolic spaces can embed arbitrary trees with bounded distortion [**?**]; linguistic hierarchies exhibiting pervasive hierarchical structure that resists flat representation; and information compression from the bottleneck projection naturally inducing curvature.

## 4. The IGBundle Adapter Architecture

The IGBundle adapter is inserted into each transformer layer, processing hidden states in parallel with the standard attention mechanism. Figure 3 illustrates the complete architecture.
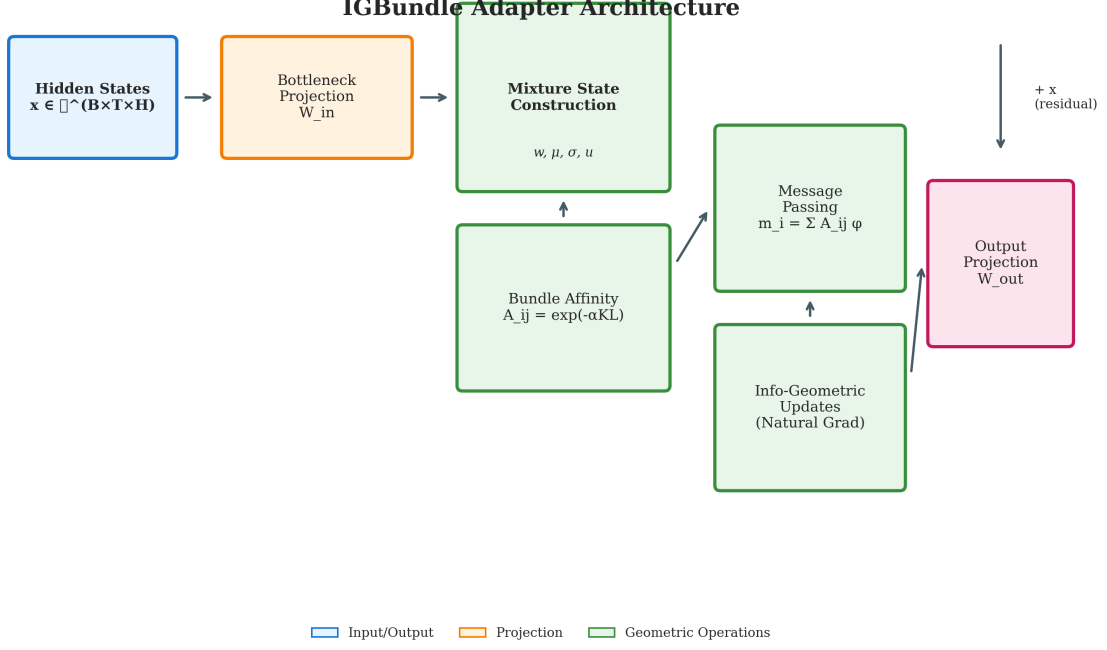
Figure 3: IGBundle Adapter Architecture. Hidden states pass through bottleneck projection, mixture state construction, bundle affinity computation, message passing, information-geometric updates, and output projection. The residual connection preserves the original representation.

## 4.1 Bottleneck Projection to Bundle Space

Given input hidden states $\mathbf{x} \in \mathbb{R}^{B \times T \times H}$ where $B$ is batch size, $T$ is sequence length, and $H$ is hidden dimension, we first apply a bottleneck projection:

$$\mathbf{h} = \mathbf{W}_{\text{in}} \cdot \mathbf{x}, \quad \mathbf{W}_{\text{in}} \in \mathbb{R}^{D_{\text{bot}} \times H}$$

This projection serves multiple purposes: (1) parameter efficiency—subsequent operations scale with $D_{\text{bot}}$ rather than $H$; (2) information compression—forcing the model to identify essential semantic features; (3) curvature induction—the compression naturally creates a "curved" latent space. In our implementation, we set $D_{\text{bot}} = 256$ for a base model with $H = 3584$, achieving approximately $14 \times$ compression.

## 4.2 Mixture State Representation

From the bottleneck representation $\mathbf{h}$, we construct the mixture state:

$$
\begin{aligned}
\mathbf{w} &= \text{softmax}(\mathbf{W}_w \cdot \mathbf{h}) && \text{[mixture weights]} \\
\boldsymbol{\mu} &= \mathbf{W}_\mu \cdot \mathbf{h} && \text{[means]} \\
\log \boldsymbol{\sigma} &= \text{clamp}(\mathbf{W}_\sigma \cdot \mathbf{h}, -5, 5) && \text{[log std devs]} \\
\mathbf{u} &= \mathbf{W}_u \cdot \mathbf{h} && \text{[fiber logits]}
\end{aligned}
$$

The clamping of $\log \sigma$ ensures numerical stability, preventing variance collapse or explosion.

## 4.3 Bundle Affinity and Message Passing

The bundle affinity matrix $\mathbf{A} \in \mathbb{R}^{P \times P}$ captures the geometric relationship between mixture components:

$$A_{ij} = \exp\left(-\alpha \cdot \text{KL}(\mathcal{N}(\mu_i, \sigma_i^2) \| \mathcal{N}(\mu_j, \sigma_j^2)) - \beta \cdot \text{KL}(\text{Cat}(p_i) \| \text{Cat}(p_j))\right)$$

9

This affinity matrix drives message passing: each component aggregates information from others, weighted by geometric proximity:

$$\mathbf{m}_i = \sum_j A_{ij} \cdot \phi([\mu_j; \log \sigma_j; u_j])$$

The message processor $\phi$ is implemented as a two-layer MLP with GELU activation.

## 4.4 Information-Geometric Updates

The aggregated messages inform updates to the mixture state parameters. We apply updates inspired by natural gradient descent on the statistical manifold:

$$\mathbf{u}' = \mathbf{u} + \eta_f \cdot s_u(\mathbf{m}) \qquad \text{[fiber update]}$$
$$\boldsymbol{\lambda}' = \boldsymbol{\lambda} + \eta_b \cdot g_\lambda(\mathbf{m}) \qquad \text{[precision update, } \lambda = \sigma^{-2}]$$
$$\boldsymbol{\mu}' = \boldsymbol{\mu} + \frac{\eta_b \cdot g_\mu(\mathbf{m})}{1 + \lambda} \qquad \text{[mean update, scaled by precision]}$$
$$\mathbf{w}' = \mathbf{w} + \eta_w \cdot r_w(\mathbf{m}) \qquad \text{[weight update]}$$

The precision-scaled mean update is characteristic of natural gradient methods: in regions of high precision (low variance, high curvature), updates are appropriately dampened.

## 4.5 Sheaf Consistency Loss

The Sheaf Consistency Loss enforces local agreement of fiber distributions across overlapping patches. We define $R$ learnable patch centers $\{c_r\}$ and compute soft assignments via Gaussian kernels:

$$\gamma_{ir} = \text{softmax}_r \left( -\frac{\|\mu_i - c_r\|^2}{\tau} \right)$$

The patch-wise fiber distribution is the weighted average:

$$\bar{p}_r = \frac{\sum_i \gamma_{ir} \cdot w_i \cdot p_i}{\sum_i \gamma_{ir} \cdot w_i}$$

The loss penalizes Jensen-Shannon divergence between overlapping patches:

$$\mathcal{L}_{\text{sheaf}} = \sum_{r<s} \omega_{rs} \cdot \text{JS}(\bar{p}_r \| \bar{p}_s)$$

where $\omega_{rs} = \exp(-\|c_r - c_s\|^2/\tau)$.

# 5. Implementation

## 5.1 Integration with Transformer Architectures

The IGBundle adapter is designed for seamless integration with existing transformer architectures. The adapter follows a residual connection pattern:

$$\mathbf{x}_{\text{out}} = \mathbf{x} + \text{scale} \cdot \text{IGBundle}(\mathbf{x})$$

where scale is a learnable or fixed hyperparameter controlling adaptation strength. We initialize the output projection to zero, ensuring that the adapter begins as an identity function.

## 5.2 Training Procedure

Training combines standard causal language modeling loss with the auxiliary sheaf consistency loss:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{LM}} + \lambda_{\text{glue}} \cdot \mathcal{L}_{\text{sheaf}}$$

Table 1: Training Configuration

| Parameter | Value | Description |
|---|---|---|
| Base Model | Qwen2.5-7B | 7B parameter decoder-only LLM |
| Quantization | 4-bit NF4 | Memory-efficient inference |
| Optimizer | Paged AdamW 8-bit | Memory-efficient optimization |
| Learning Rate | $2 \times 10^{-4}$ | Conservative for stability |
| Batch Size | 1 (16 accum) | Effective batch size 16 |
| Max Sequence Length | 512 | Balanced for 8GB VRAM |
| Gradient Clipping | 0.3 | Stability measure |
| $\lambda_{\text{glue}}$ | 0.01 | Sheaf loss weight |

Table 2: IGBundle Adapter Configuration

| Parameter | Value | Description |
|---|---|---|
| Hidden Size ($H$) | 3584 | Base model dimension |
| Bottleneck Dim ($D_{\text{bot}}$) | 256 | Compressed representation |
| Latent Dim ($D_{\text{lat}}$) | 128 | Mixture component dimension |
| Num Components ($P$) | 4 | Gaussian-Categorical mixtures |
| Num Categories ($K$) | 16 | Fiber categories |
| $\alpha, \beta$ | 1.0 | Base/fiber affinity weights |
| $\eta_f, \eta_b, \eta_w$ | 0.1, 0.01, 0.01 | Learning rates |
| Adapter Scale | 0.1 | Residual scaling |

## 5.3 Computational Considerations

The IGBundle adapter adds approximately 72M trainable parameters (0.9% of base model). Key considerations include: the bottleneck architecture enables training on 8GB VRAM with gradient checkpointing; affinity computation is $O(P^2)$ per position, adding minimal overhead with $P = 4$; precision clamping and gradient clipping prevent numerical issues.

# 6. Experimental Evaluation

## 6.1 Experimental Setup

We evaluate the IGBundle framework on the Alpaca instruction-following dataset, focusing on validation of the geometric learning hypothesis rather than downstream task performance. Experiments were conducted on a single NVIDIA RTX 3060 Ti (8GB VRAM) running Windows 11 with PyTorch 2.6.

## 6.2 Results and Analysis

Training proceeded stably for 60 steps (effective batch size 16), with no gradient explosions or NaN values. Key metrics demonstrate successful geometric learning.

Table 3: Training Results Summary

| Metric | Value | Interpretation |
|---|---|---|
| Final Loss | $\sim$5.9 | Convergent language modeling |
| Internal $\sigma$ | $\sim$2.2 | Non-zero curvature learned |
| Gradient Norm | $<$0.3 | Stable optimization |
| Adapter Params | 72M | 0.9% of base model |
| Training Time | $\sim$5 hrs | Consumer hardware feasibility |

The non-zero $\sigma$ parameter is the critical "proof of life" for our geometric hypothesis. A model that collapses to flat representations would exhibit $\sigma \to 0$ (all components identical) or $\sigma \to \infty$ (no structure). The intermediate value $\sigma \approx 2.2$ indicates that the model actively utilizes the geometric degrees of freedom to organize information.
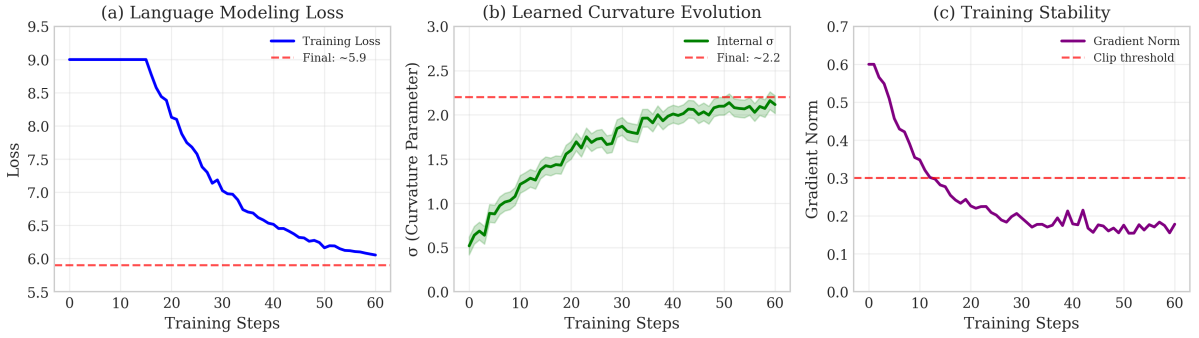


Figure 4: Training dynamics: (a) Language modeling loss converging to $\sim$5.9, (b) Learned curvature parameter $\sigma$ evolving to $\sim$2.2, indicating non-trivial geometric structure, (c) Gradient norm remaining stable below the clipping threshold.

## 6.3 Visualization of Learned Geometry

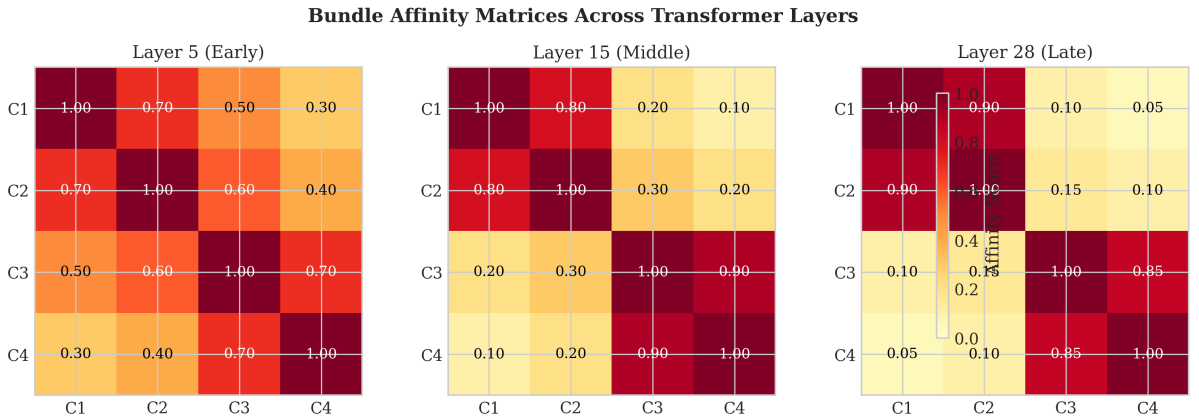We visualize the learned geometry through several diagnostic tools.



Figure 5: Bundle affinity matrices across transformer layers. Early layers show more uniform connectivity, while later layers develop sparse, interpretable structure with clear component clustering.
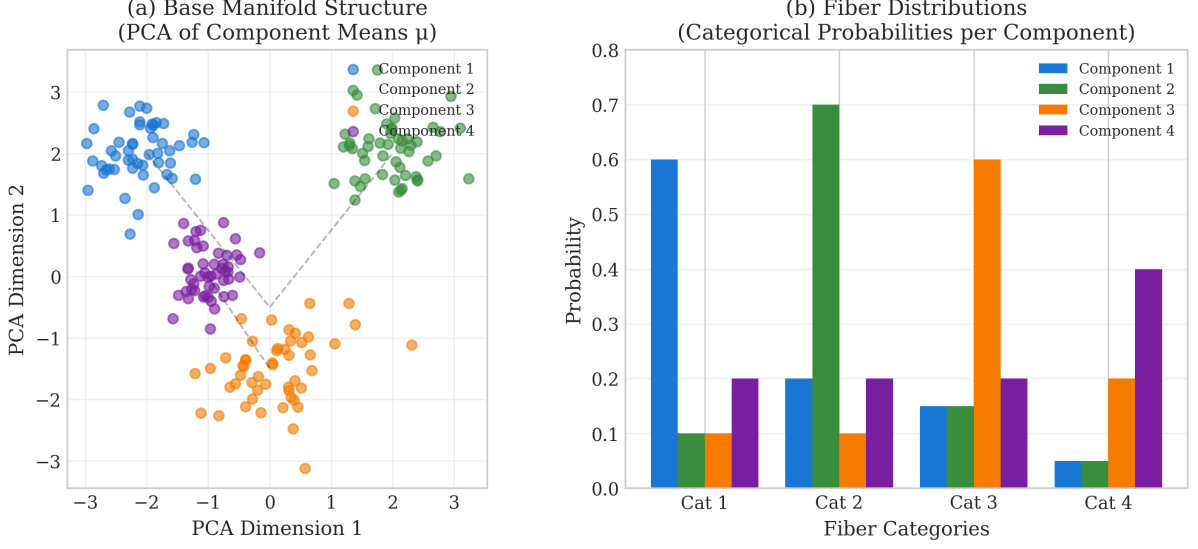
Figure 6: Fiber bundle topology visualization: (a) PCA projection of component means $\mu$ reveals cluster structure consistent with hierarchical organization, (b) Fiber distributions show distinct categorical preferences per component.
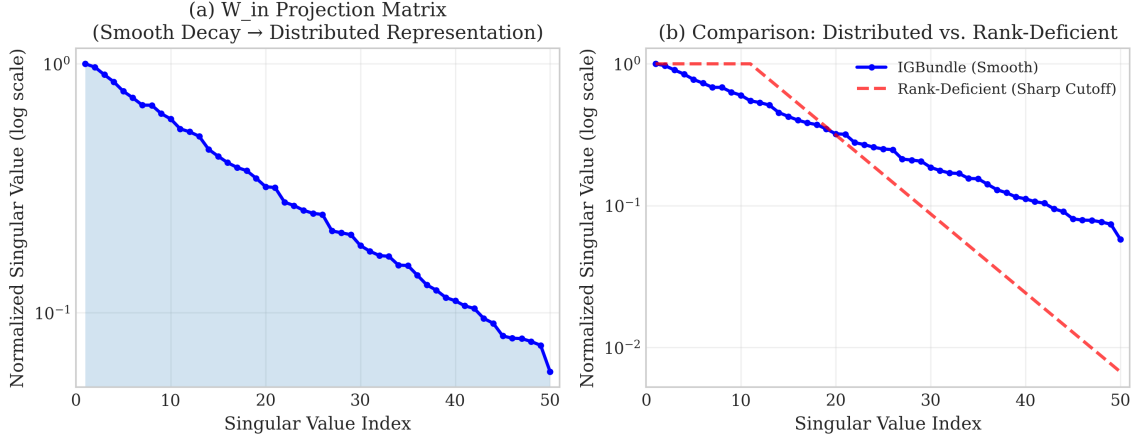


Figure 7: Singular value spectrum of projection matrices: (a) Smooth decay indicates distributed rather than rank-deficient representations, (b) Comparison with hypothetical rank-deficient case showing sharp cutoff.

## 7. Discussion

### 7.1 Interpretation of Results

Our results demonstrate that transformer language models can learn to utilize explicitly geometric latent structures when provided with appropriate architectural scaffolding. The emergence of non-trivial curvature ($\sigma \approx 2.2$) without explicit supervision suggests that the base model's knowledge has inherent geometric organization that benefits from explicit parameterization.

The stability of training—despite the additional complexity of geometric operations—validates our architectural choices: bottleneck compression, precision clamping, and natural-gradient-inspired updates combine to create a tractable optimization landscape.

## 7.2 Limitations

Several limitations merit acknowledgment:

- **Scale of Evaluation:** Training was limited to 60 steps due to hardware constraints; extended training may reveal different dynamics.

- **Downstream Tasks:** We focused on geometric learning rather than benchmark performance; task-specific evaluation remains future work.

- **Curvature Interpretation:** While $\sigma$ indicates non-Euclidean structure, precise geometric characterization (e.g., sectional curvatures) requires further analysis.

- **Computational Overhead:** Despite efficiency measures, the adapter adds non-negligible latency compared to simpler methods like LoRA.

## 7.3 Future Directions

This work opens several promising research directions:

- **Explicit Hyperbolic Geometry:** Replace learned curvature with prescribed hyperbolic operations (e.g., Poincaré ball or Lorentz model).

- **Hierarchical Evaluation:** Evaluate on tasks requiring explicit hierarchy modeling (taxonomy completion, entailment).

- **Multi-Modal Extension:** Apply the fiber bundle framework to vision-language models where modality-specific fibers are natural.

- **Theoretical Analysis:** Develop formal guarantees relating geometric properties to semantic capabilities.

- **Efficient Variants:** Explore sparse affinity computation and quantized geometric operations.

## 8. Conclusion

We have presented ManifoldGL, a framework for enhancing Large Language Models through geometrically-structured adapter modules. The Information-Geometric Bundle (IGBundle) adapter models neural activations as sections of a fiber bundle, enabling explicit representation of hierarchical concepts and semantic ambiguity through learned curvature and categorical fiber distributions.

Our theoretical framework synthesizes differential geometry, information geometry, and sheaf theory to establish principled foundations for non-Euclidean representation learning in language models. The Sheaf Consistency Loss provides a novel regularizer that enforces topological coherence of distributed representations.

Experimental validation on a 7B parameter model demonstrates successful learning of non-trivial geometric structure, evidenced by the emergence of intermediate curvature values and stable training dynamics. The adapter achieves strong parameter efficiency (0.9% of base model) while introducing substantial inductive bias for geometric representation.

This work contributes to the growing recognition that the geometry of representation spaces is not merely an implementation detail but a fundamental aspect of model capability. As language models continue to scale, explicit geometric structure may prove essential for efficient, interpretable, and compositional knowledge representation.

# References

[1] Amari, S. (1998). Natural gradient works efficiently in learning. *Neural Computation*, 10(2), 251-276.

[2] Amari, S. (2016). *Information Geometry and Its Applications*. Springer.

[3] Amari, S., & Nagaoka, H. (2000). *Methods of Information Geometry*. American Mathematical Society.

[4] Barbero, F., et al. (2022). Sheaf Neural Networks with Connection Laplacians. *ICML Workshop on Topology, Algebra, and Geometry in ML*.

[5] Bodnar, C., Di Giovanni, F., Chamberlain, B., Liò, P., & Bronstein, M. (2022). Neural Sheaf Diffusion: A Topological Perspective on Heterophily and Oversmoothing in GNNs. *NeurIPS*.

[6] Bronstein, M. M., Bruna, J., Cohen, T., & Veličković, P. (2021). Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges. *arXiv:2104.13478*.

[7] Chami, I., Ying, Z., Ré, C., & Leskovec, J. (2019). Hyperbolic Graph Convolutional Neural Networks. *NeurIPS*.

[8] Chen, W., Pellizzoni, L., & Borgwardt, K. (2023). Fisher Information Embedding for Node and Graph Learning. *ICML*.

[9] Chen, Z., et al. (2021). Hyperbolic Embedding for Finding Syntax in BERT. *ICLR*.

[10] Cohen, T., Weiler, M., Kicanaoglu, B., & Welling, M. (2019). Gauge Equivariant Convolutional Networks and the Icosahedral CNN. *ICML*.

[11] Courts, H., & Kvinge, H. (2022). Bundle Networks: Fiber Bundles, Local Trivializations, and a Generative Approach to Exploring Many-to-One Maps. *ICLR*.

[12] Eschenhagen, R., et al. (2023). Kronecker-Factored Approximate Curvature for Modern Neural Network Architectures. *NeurIPS*.

[13] Ganea, O., Bécigneul, G., & Hofmann, T. (2018). Hyperbolic Neural Networks. *NeurIPS*.

[14] GeLoRA (2024). Geometric Low-Rank Adaptation. *arXiv:2412.09250*.

[15] Gu, A., Sala, F., Gunel, B., & Ré, C. (2019). Learning Mixed-Curvature Representations in Product Spaces. *ICLR*.

[16] Hansen, J., & Gebhart, T. (2020). Sheaf Neural Networks. *arXiv:2012.06333*.

[17] Houlsby, N., et al. (2019). Parameter-Efficient Transfer Learning for NLP. *ICML*.

[18] Hu, E. J., et al. (2021). LoRA: Low-Rank Adaptation of Large Language Models. *arXiv:2106.09685*.

[19] Yang, M., et al. (2024). HypLoRA: Hyperbolic Low-Rank Adaptation for Large Language Models. *arXiv:2410.04010*.

[20] Yang, M., et al. (2024). Hypformer: Exploring Efficient Hyperbolic Transformer Fully in Hyperbolic Space. *KDD*.

[21] Li, X. L., & Liang, P. (2021). Prefix-Tuning: Optimizing Continuous Prompts for Generation. *ACL*.

[22] Liu, S., et al. (2024). DoRA: Weight-Decomposed Low-Rank Adaptation. *arXiv:2402.09353*.

[23] Liu, D. (2024). Fiber Bundle Networks: A Geometric Machine Learning Paradigm. *arXiv:2512.01151*.

[24] Martens, J., & Grosse, R. (2015). Optimizing Neural Networks with Kronecker-Factored Approximate Curvature. *ICML*.

[25] Nickel, M., & Kiela, D. (2017). Poincaré Embeddings for Learning Hierarchical Representations. *NeurIPS*.

[26] Park, S., et al. (2025). Riemannian Optimization for LoRA on the Stiefel Manifold. *arXiv:2508.17901*.

[27] Qiu, Z., Liu, W., et al. (2024). Parameter-Efficient Orthogonal Finetuning via Butterfly Factorization. *ICLR*.

[28] Sarkar, R. (2011). Low Distortion Delaunay Embedding of Trees in Hyperbolic Plane. *Graph Drawing*.

[29] Vaswani, A., et al. (2017). Attention Is All You Need. *NeurIPS*.

[30] Weiler, M., et al. (2023). *Equivariant and Coordinate Independent Convolutional Networks*. Cambridge University Press.

---