

# A Markov Categorical Framework for Language Modeling

Yifan Zhang

## Abstract

Auto-regressive language models factorize sequence probabilities as  $P_\theta(\mathbf{w}) = \prod_t P_\theta(w_t|\mathbf{w}_{<t})$  and are trained by minimizing the negative log-likelihood (NLL) objective. While empirically powerful, a deep theoretical understanding of why this simple objective yields such versatile representations remains elusive. This work introduces a unifying analytical framework using Markov Categories (MCs) to deconstruct the AR generation process and the NLL objective. We model the single-step generation map  $\mathbf{w}_{<t} \mapsto P_\theta(\cdot|\mathbf{w}_{<t})$  as a composition of Markov kernels,  $k_{\text{gen},\theta} = k_{\text{head}} \circ k_{\text{bb}} \circ k_{\text{emb}}$ , in the category **Stoch**. This compositional view, when enriched with statistical divergences, allows us to dissect information flow and learned geometry. Our framework makes three main contributions. First, we provide a formal, information-theoretic rationale for the success of modern speculative decoding methods like EAGLE (Li et al., 2024), quantifying the information surplus in hidden states that these methods exploit. Second, we formalize how NLL minimization forces the model to learn not just the next token, but the data’s intrinsic conditional stochasticity, a process we analyze using categorical entropy. Third, and most centrally, we prove that NLL training acts as an implicit form of spectral contrastive learning. By analyzing the information geometry of the model’s prediction head, we show that NLL implicitly forces the learned representation space to align with the eigenspectrum of a predictive similarity operator, thereby learning a geometrically structured space without explicit contrastive pairs. This compositional and information-geometric perspective reveals the deep structural principles underlying the effectiveness of modern LMs.

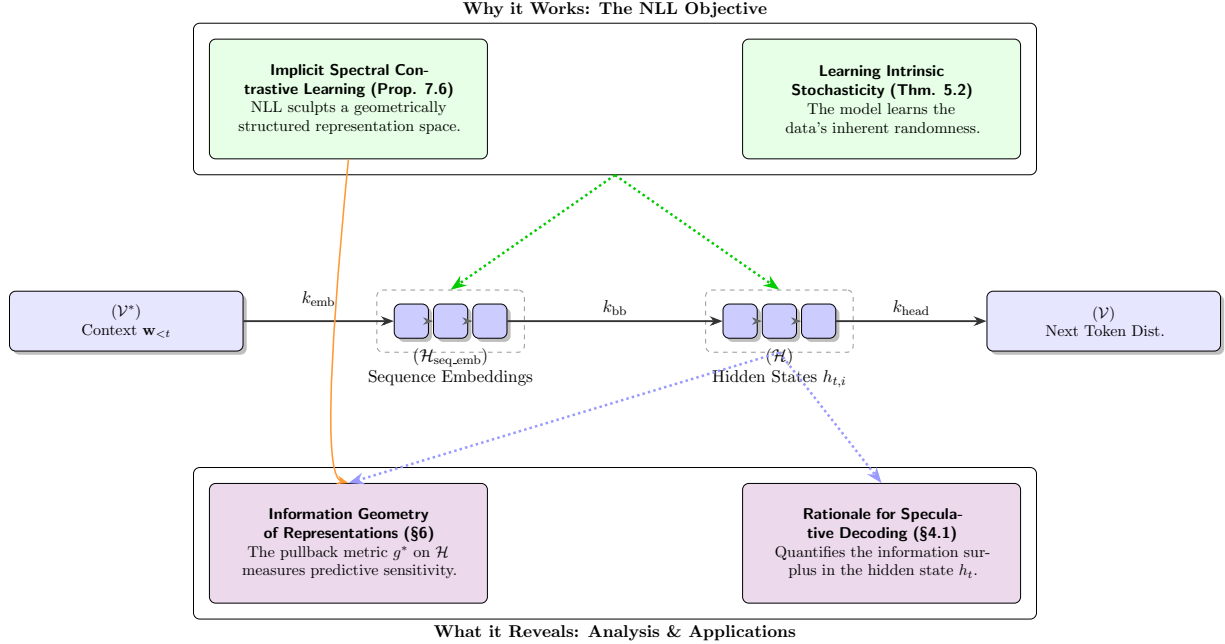
Date: July 25, 2025

Project Page: <https://github.com/asiresearch/lm-theory>

## 1 Introduction

Auto-regressive language models (AR LMs), particularly those based on the Transformer architecture (Vaswani et al., 2017; Radford et al., 2019; Brown et al., 2020), have achieved remarkable success, defining the state-of-the-art in natural language generation and demonstrating impressive few-shot learning capabilities. These models operate by sequentially predicting the next token in a sequence based on the preceding context. Formally, given a sequence  $\mathbf{w} = w_1 \dots w_L$  with tokens  $w_i$  from a finite vocabulary  $\mathbb{V}$ , the model learns a parameterized probability distribution  $P_\theta$  that factorizes as:

$$P_\theta(\mathbf{w}) = \prod_{t=1}^L P_\theta(w_t|\mathbf{w}_{<t}), \quad (1.1)$$



**Figure 1** A conceptual overview of our framework. **Center:** The core thesis models the auto-regressive generation step as a composition of Markov kernels  $k_{\text{gen}} = k_{\text{head}} \circ k_{\text{bb}} \circ k_{\text{emb}}$  in the category **Stoch**. This separates the deterministic context encoding ( $k_{\text{emb}}, k_{\text{bb}}$ ) from the stochastic prediction ( $k_{\text{head}}$ ). **Top:** This compositional lens reveals the deeper mechanisms of the NLL objective, which we re-frame as minimizing the average KL divergence between the model and true data kernels. We prove that this single objective implicitly forces the model to learn two key properties of the data: its intrinsic conditional stochasticity (measured via categorical entropy) and its underlying geometric structure, which we show is equivalent to performing spectral contrastive learning on a predictive similarity graph. **Bottom:** The framework yields new analytical tools. By pulling back the Fisher-Rao metric, we endow the representation space  $\mathcal{H}$  with an information geometry that quantifies predictive sensitivity. This provides a formal basis for analyzing the information surplus in hidden states, giving a rigorous theoretical foundation for modern speculative decoding methods.

where  $\mathbf{w}_{<t} := w_1 \dots w_{t-1}$  is the context sequence, and  $\theta$  denotes the model parameters, typically optimized by minimizing the negative log-likelihood (NLL) on vast text corpora. The core computational step is the mapping from a context  $\mathbf{w}_{<t}$  to the conditional probability distribution  $P_\theta(\cdot | \mathbf{w}_{<t})$  over  $\mathbb{V}$  for the next token  $w_t$ .

Despite their empirical triumphs, a deep theoretical understanding of their internal mechanisms remains incomplete (Manning et al., 2020; Elhage et al., 2021; Yuan, 2023). Current analysis often relies on empirical probes (Hewitt and Manning, 2019) or studies of specific components like attention heads (Olsson et al., 2022). While insightful, these methods can be fragmented and often lack a unified mathematical language to describe the model’s compositional and stochastic nature as a whole. A fundamental open question is why the simple NLL objective is so effective, leading to representations that capture complex linguistic and world knowledge. Another critical challenge is improving the slow, sequential nature of AR generation. Recent advances in speculative decoding, such as EAGLE (Li et al., 2024), have achieved significant speedups by predicting multiple tokens in parallel, suggesting that the final hidden state  $h_t$  contains far more information than is needed for predicting only the single next token  $w_t$ . However, a formal understanding of this information

surplus is lacking.

This paper addresses this gap by introducing a unifying analytical framework for AR LMs. Our central thesis is that the language of *Markov Categories (MCs)* (Cho and Jacobs, 2019; Fritz, 2020b) provides the natural mathematical setting to provide a single, unified mathematical language to formally connect several concepts that are usually discussed separately: information flow through the model’s components, the geometry of the learned representation space, and the structural effects of the NLL training objective.

While many individual mathematical tools we employ—such as the pullback of the Fisher-Rao metric or the connection between NLL and KL divergence—are well-established, our primary contribution is their novel synthesis and application to dissect AR LMs. The originality of this work lies in using the category **Stoch** to formally model compositional information flow, leading to new insights uniquely enabled by this perspective. Unlike information-theoretic analyses that treat models as monolithic black boxes analyzing external behavior (e.g., the entropy of the output sequence), our framework uses categorical information theory to analyze the internal transformations and the learned geometry of the representation space at each stage of processing.

This paper introduces an analytical framework focused on the internal mechanics of the AR generation step  $\mathbf{w}_{<t} \mapsto P_\theta(\cdot|\mathbf{w}_{<t})$ . We leverage the category **Stoch**, a canonical MC whose objects are standard Borel spaces (like the continuous representation space  $\mathcal{H} \cong \mathbb{R}^{d_{\text{model}}}$ ) and whose morphisms are Markov kernels (Kallenberg and Kallenberg, 1997; Fritz, 2020b). We provide the first formal model of the AR generation step as a composite kernel in **Stoch**:

$$k_{\text{gen},\theta} := k_{\text{head}} \circ k_{\text{bb}} \circ k_{\text{emb}} : (\mathbb{V}^*, \mathcal{B}(\mathbb{V}^*)) \rightarrow (\mathbb{V}, \mathcal{P}(\mathbb{V})). \quad (1.2)$$

Here,  $k_{\text{emb}}$  and  $k_{\text{bb}}$  represent the typically deterministic context embedding and backbone transformations that produce the final hidden state  $h_t \in \mathcal{H}$ , while  $k_{\text{head}}$  is the generally stochastic kernel mapping  $h_t$  to the predictive distribution  $P_\theta(\cdot|\mathbf{w}_{<t})$ .

A crucial aspect of our framework is enriching **Stoch** with a statistical divergence  $D$  (e.g.,  $D_{\text{KL}}$ ) (Baez et al., 2016; Perrone, 2023a,b). This allows for defining intrinsic, categorical information measures like entropy  $\mathcal{H}_D$  and mutual information  $I_D$  (Perrone, 2023a), which automatically satisfy the Data Processing Inequality (DPI). Leveraging this unified framework, this paper makes the following contributions:

1. **A Formal Compositional Model for Information Flow (Section 3):** We formally model the AR generation step as a composite kernel,  $k_{\text{gen},\theta} = k_{\text{head}} \circ k_{\text{bb}} \circ k_{\text{emb}}$ . This compositional structure is a powerful tool for reasoning about how information is transformed, preserved, or lost at each distinct stage of processing.
2. **Information-Theoretic Rationale for Speculative Decoding (Section 4.1):** We leverage the DPI to provide a formal, quantitative explanation for the success of methods like EAGLE (Li et al., 2024). We show that the information surplus these methods exploit—the information in a hidden state  $H_t$  about multiple future tokens—can be rigorously quantified within our framework.
3. **A Unified View of the NLL Objective (Sections 5 and 7):** We show how the MC framework unifies three critical interpretations of NLL training under one theoretical roof:

- **Compression:** NLL as KL minimization, equivalent to optimal source coding.
- **Learning Stochasticity:** NLL forces the model to learn the data’s inherent randomness, a process we formalize using categorical entropy. We show that optimizing NLL implies that the model’s learned stochasticity converges to that of the data (Theorem 5.2).
- **Implicit Spectral Contrastive Learning:** We prove that NLL acts as an implicit spectral contrastive objective. By analyzing the information geometry of the prediction head, we demonstrate that NLL forces the model to learn a geometrically structured representation space by aligning representations with the eigenspectrum of a predictive similarity operator, effectively separating representations of predictively dissimilar contexts (Proposition 7.6).

The primary goal of this paper is to establish the theoretical framework itself, which is a substantial contribution that provides the essential groundwork for future empirical investigation. While estimating some of our proposed information-theoretic quantities in high dimensions is a known challenge (discussed in Section 4), the framework’s immediate value lies in providing a new, more powerful lens for theoretical analysis. By formalizing how information is transformed (Section 4), how predictive sensitivity is encoded in representation geometry (Section 6), and how the NLL objective implicitly structures representations (Section 7), we can move towards more principled approaches to model design, interpretation, and control.

The paper is organized as follows. Sections 2 and 3 introduce the MC framework and our compositional model of LMs. Section 4 uses the framework to analyze information flow, providing a rationale for speculative decoding. Section 5 connects the NLL objective to learning the data’s intrinsic stochasticity. Sections 6 and 7 present our main theoretical result, showing how NLL performs implicit spectral learning by shaping the geometry of the representation space. Sections 8 and 9 discuss related work and conclude.

## 2 Background

This section reviews the essential mathematical concepts forming the foundation of our framework: the definition of Markov Categories and the specific category **Stoch**, followed by the enrichment of **Stoch** with statistical divergences leading to categorical information measures.

### 2.1 Markov Categories and Stoch

Markov Categories provide an axiomatic framework for probability and stochastic processes using category theory (Fritz, 2020b).

**Definition 2.1** (Markov Category (Fritz, 2020b)). A Markov category  $(\mathcal{C}, \otimes, I)$  is a symmetric monoidal category where each object  $X$  is equipped with a commutative comonoid structure  $(\Delta_X : X \rightarrow X \otimes X, !_X : X \rightarrow I)$  that is natural in  $X$ , and the monoidal unit  $I$  is a terminal object (the *causality* axiom:  $!_X$  is the unique map  $X \rightarrow I$ ).

Morphisms  $k : X \rightarrow Y$  are interpreted as stochastic processes. Composition  $h \circ k$  is sequential processing, while  $k \otimes h$  is parallel processing. The comonoid maps  $\Delta_X$  (copy) and  $!_X$  (discard) abstractly model the duplication and deletion of information. The causality axiom enforces that

discarding information is deterministic and ultimately reflects probability normalization ( $\int k(x, dy) = 1$ ) in concrete examples like **Stoch**. States (probability distributions) on an object  $X$  are represented as morphisms  $p : I \rightarrow X$ .

The key example for our purposes is the category **Stoch**.

**Definition 2.2** (Category **Stoch** (Fritz, 2020b; Perrone, 2023a)). The Markov category **Stoch** is defined by:

- **Objects:** Standard Borel spaces  $(X, \mathcal{B}(X))$ . These are general measure spaces that include finite sets (like a vocabulary  $\mathbb{V}$ ), countable sets, and continuous spaces like Euclidean space  $\mathbb{R}^d$  or other Polish spaces. This ensures the framework can handle both discrete tokens and continuous representations. The monoidal unit  $I$  is a singleton space  $(\{\star\}, \{\emptyset, \{\star\}\})$ .
- **Morphisms:** Markov kernels  $k : X \rightarrow Y$ . A map  $k : X \times \mathcal{B}(Y) \rightarrow [0, 1]$  where  $k(x, \cdot)$  is a probability measure on  $Y$  for each  $x \in X$ , and  $k(\cdot, A)$  is a measurable function on  $X$  for each  $A \in \mathcal{B}(Y)$ .
- **Composition:** Given  $k : X \rightarrow Y$  and  $h : Y \rightarrow Z$ , the composite  $h \circ k : X \rightarrow Z$  is  $(h \circ k)(x, C) := \int_Y h(y, C) k(x, dy)$  (Chapman-Kolmogorov). Identity  $\text{id}_X(x, A) = \delta_x(A)$ .
- **Monoidal Product ( $\otimes$ ):** Product space  $(X \times Y, \mathcal{B}(X) \otimes \mathcal{B}(Y))$  with the product  $\sigma$ -algebra. Product kernel  $(k \otimes h)((x, y), \cdot) := k(x, \cdot) \otimes h(y, \cdot)$  (product measure).
- **Symmetry:** Swap map  $\sigma_{X,Y} : X \otimes Y \rightarrow Y \otimes X$  is  $\sigma_{X,Y}((x, y), \cdot) = \delta_{(y,x)}$ .
- **Comonoid Structure:** Copy  $\Delta_X : X \rightarrow X \otimes X$  is  $\Delta_X(x, \cdot) = \delta_{(x,x)}$ . Discard  $!_X : X \rightarrow I$  maps to the unique point measure on  $I$ ,  $!_X(x, \{\star\}) = 1$ .
- **Causality:**  $I$  is terminal,  $!_Y \circ k = !_X$  holds, reflecting probability normalization.

**Remark 2.3** (Interpretation). In **Stoch**, objects represent the types of random outcomes (e.g., sequences, vectors, tokens). Morphisms represent stochastic processes or channels mapping inputs to probability distributions over outputs. Deterministic functions  $f : X \rightarrow Y$  correspond to deterministic kernels  $k_f(x, \cdot) = \delta_{f(x)}$ . States  $p : I \rightarrow X$  correspond bijectively to probability measures  $\mu_p \in \mathcal{P}(X)$  via  $\mu_p(A) = p(\star, A)$ . Marginalization arises from discarding information, e.g., for a joint state  $p : I \rightarrow X \otimes Y$ , the  $X$ -marginal is  $p_X = (\text{id}_X \otimes !_Y) \circ p$ .

## 2.2 Divergence Enrichment and Categorical Information Measures

The structure of **Stoch** is particularly powerful when enriched with a statistical divergence  $D$ , quantifying the dissimilarity between probability measures (states)  $p, q : I \rightarrow X$ , written  $D_X(p||q)$  (Perrone, 2023a). Examples include KL divergence ( $D_{\text{KL}}$ ), Total Variation ( $d_{\text{TV}}$ ), Rényi divergences ( $D_\alpha$ ), and the broad class of  $f$ -divergences ( $D_f$ ) (Amari and Nagaoka, 2000; Nowozin et al., 2016).

A fundamental property linking divergences and Markov kernels is the Data Processing Inequality (DPI), which holds for most standard divergences (e.g.,  $f$ -divergences, Rényi  $\alpha \in [0, \infty]$ ).

**Theorem 2.4** (Data Processing Inequality (DPI)). Let  $D$  be a statistical divergence satisfying the DPI. For any Markov kernel  $k : X \rightarrow Y$  in **Stoch** and any pair of states  $p, q : I \rightarrow X$ :

$$D_Y(k \circ p || k \circ q) \leq D_X(p || q) \quad (2.1)$$

Processing through the channel  $k$  cannot increase the  $D$ -divergence between the distributions.

Based on this, Perrone (Perrone, 2023a) introduced categorical definitions of entropy and mutual information intrinsically tied to the divergence  $D$  and the MC structure.

**Definition 2.5** (Categorical Entropy (Perrone, 2023a)). Let  $(\mathbf{Stoch}, D)$  be enriched with a DPI-satisfying divergence  $D$ .

1. The **Categorical Entropy** of a kernel  $k : X \rightarrow Y$  measures its intrinsic stochasticity:

$$\mathcal{H}_D(k) := D_{Y \otimes Y}(\Delta_Y \circ k \parallel (k \otimes k) \circ \Delta_X) \quad (2.2)$$

Intuitively, this compares two processes: (1) applying  $k$  to an input  $x$  and then copying the output to get  $(y, y)$ ; versus (2) copying the input  $x$  and applying  $k$  independently to each copy to get  $(y_1, y_2)$ . The more stochastic  $k$  is, the greater the divergence. If  $k$  is deterministic, the entropy is zero.

2. The **Categorical Mutual Information** for a joint state  $p : I \rightarrow X \otimes Y$  is defined as  $I_D(p) := D_{X \otimes Y}(p \parallel p_X \otimes p_Y)$ , where  $p_X$  and  $p_Y$  are the marginal states. It measures the dependence between  $X$  and  $Y$ .

**Remark 2.6** (Properties and Connections). When  $D = D_{\text{KL}}$ ,  $I_{D_{\text{KL}}}(p)$  recovers the standard Shannon mutual information  $I(X; Y)$ .  $\mathcal{H}_{D_{\text{KL}}}(k)$  provides an intrinsic measure of the kernel’s stochasticity, related to but distinct from average conditional Shannon entropy (Perrone, 2023a). Crucially, these definitions automatically satisfy the DPI. For instance, for a state  $p : I \rightarrow X \otimes Y$  and a kernel  $h : Y \rightarrow Z$ , the DPI implies  $I_D(X; Y) \geq I_D(X; Z)$ , reflecting that processing  $(Y \rightarrow Z)$  cannot create information about  $X$ . Information geometry (Amari and Nagaoka, 2000) arises naturally, as the Fisher-Rao metric is induced by the local quadratic approximation of KL divergence.

### 3 Auto-Regressive Language Models as Composed Kernels

We now apply the Markov Category framework established in Section 2 to model auto-regressive language models. Specifically, we model the single-step generation mapping  $\mathbf{w}_{<t} \mapsto P_\theta(\cdot | \mathbf{w}_{<t})$  as a composition of Markov kernels within the category  $\mathbf{Stoch}$ .

The relevant measurable spaces (objects in  $\mathbf{Stoch}$ ) are:

- Input context space:  $(\mathbb{V}^*, \mathcal{B}(\mathbb{V}^*)) = (\mathbb{V}^*, \mathcal{B}(\mathbb{V}^*))$ , where  $\mathbb{V}^*$  is the set of finite sequences over the vocabulary  $\mathbb{V}$ , equipped with a suitable  $\sigma$ -algebra making it standard Borel (e.g., considering it as a disjoint union of finite products  $\mathbb{V}^n$ ).
- Initial sequence representation space:  $(\mathcal{H}_{\text{seq.emb}}, \mathcal{B}(\mathcal{H}_{\text{seq.emb}})) = (\mathcal{H}_{\text{seq.emb}}, \mathcal{B}(\mathcal{H}_{\text{seq.emb}}))$ , the space of initial vector sequences (e.g.,  $\bigcup_n (\mathbb{R}^{d_{\text{model}}})^n$ ), also equipped with a standard Borel structure.
- Final hidden state space:  $(\mathcal{H}, \mathcal{B}(\mathcal{H})) = (\mathcal{H}, \mathcal{B}(\mathcal{H}))$ , typically  $(\mathbb{R}^{d_{\text{model}}}, \mathcal{B}(\mathbb{R}^{d_{\text{model}}}))$ .
- Output vocabulary space:  $(\mathbb{V}, \mathcal{P}(\mathbb{V})) = (\mathbb{V}, \mathcal{P}(\mathbb{V}))$ , a finite measurable space.

Standard Borel spaces are chosen because they form a well-behaved class of measurable spaces (isomorphic to Borel subsets of Polish spaces) closed under countable products, sums, and containing



standard examples like  $\mathbb{R}^d$  and finite sets, ensuring measure-theoretic regularity (Kallenberg and Kallenberg, 1997).

The generation process decomposes into three kernels (morphisms in **Stoch**):

1, **Embedding Layer Kernel** ( $k_{\text{emb}} : (\mathbb{V}^*, \mathcal{B}(\mathbb{V}^*)) \rightarrow (\mathcal{H}_{\text{seq\_emb}}, \mathcal{B}(\mathcal{H}_{\text{seq\_emb}}))$ ): This kernel encapsulates the initial processing of the discrete input sequence  $\mathbf{w}_{<t} \in \mathbb{V}^*$ . It typically involves applying a token embedding function  $\mathcal{E} : \mathbb{V} \rightarrow \mathbb{R}^{d_{\text{model}}}$  to each token  $w_i$  and potentially incorporating absolute positional encodings. Let  $f_{\text{emb}} : \mathbb{V}^* \rightarrow \mathcal{H}_{\text{seq\_emb}}$  denote the overall deterministic function computing the initial sequence representation  $E_{<t}$ . Since this mapping is deterministic, the kernel  $k_{\text{emb}}$  is defined via the Dirac measure  $\delta$ :

$$k_{\text{emb}}(\mathbf{w}_{<t}, A) := \delta_{f_{\text{emb}}(\mathbf{w}_{<t})}(A) = \mathbf{1}_A(f_{\text{emb}}(\mathbf{w}_{<t})), \quad \text{for } A \in \mathcal{B}(\mathcal{H}_{\text{seq\_emb}}). \quad (3.1)$$

This is a valid morphism in **Stoch**.

2, **Backbone Transformation Kernel** ( $k_{\text{bb}} : (\mathcal{H}_{\text{seq\_emb}}, \mathcal{B}(\mathcal{H}_{\text{seq\_emb}})) \rightarrow (\mathcal{H}, \mathcal{B}(\mathcal{H}))$ ): This kernel represents the core computation, usually a deep neural network like a Transformer stack. Let  $f_{\text{bb}} : \mathcal{H}_{\text{seq\_emb}} \rightarrow \mathcal{H}$  be the function mapping the initial sequence representation  $E_{<t}$  to the final hidden state  $h_t \in \mathcal{H}$  (often the output vector at the last sequence position). This function incorporates complex operations like multi-head self-attention and feed-forward layers. Relative positional information, such as Rotary Position Embeddings (RoPE) (Su et al., 2024), is implemented within the function  $f_{\text{bb}}$  by modifying attention computations based on token positions. Assuming the backbone computation is deterministic for a given  $E_{<t}$  and parameters  $\theta$ , the kernel  $k_{\text{bb}}$  is also deterministic:

$$k_{\text{bb}}(E_{<t}, B) := \delta_{f_{\text{bb}}(E_{<t})}(B) = \mathbf{1}_B(f_{\text{bb}}(E_{<t})), \quad \text{for } B \in \mathcal{B}(\mathcal{H}). \quad (3.2)$$

This is also a morphism in **Stoch**.

3, **LM Head Kernel** ( $k_{\text{head}} : (\mathcal{H}, \mathcal{B}(\mathcal{H})) \rightarrow (\mathbb{V}, \mathcal{P}(\mathbb{V}))$ ): This final kernel maps the summary hidden state  $h_t \in \mathcal{H}$  to a probability distribution over the finite vocabulary  $\mathbb{V}$ . Typically,  $h_t$  is passed through a linear layer ( $f_{\text{head}} : \mathcal{H} \rightarrow \mathbb{R}^{|\mathbb{V}|}$ ) producing logits  $\mathbf{z} = f_{\text{head}}(h_t)$ , followed by the softmax function:  $P(w|h_t) = [\text{softmax}(\mathbf{z})]_w$ . This defines a genuinely stochastic Markov kernel:

$$k_{\text{head}}(h, A) := \sum_{w \in A} [\text{softmax}(f_{\text{head}}(h))]_w \quad \text{for } h \in \mathcal{H}, A \subseteq \mathbb{V}. \quad (3.3)$$

This kernel maps each point  $h$  in the representation space to a probability measure on the discrete space  $\mathbb{V}$ , satisfying the required measurability conditions. It is a morphism in **Stoch**.

The overall single-step generation kernel  $k_{\text{gen}, \theta} : (\mathbb{V}^*, \mathcal{B}(\mathbb{V}^*)) \rightarrow (\mathbb{V}, \mathcal{P}(\mathbb{V}))$  is the composition  $k_{\text{head}} \circ k_{\text{bb}} \circ k_{\text{emb}}$  in the category **Stoch**. This composition precisely represents the model’s learned conditional probability map  $P_\theta(\cdot | \mathbf{w}_{<t})$ . It is crucial to note that this formalism applies to **any** AR model, including Transformers. The attention mechanism provides a powerful, history-dependent parameterization of this single-step kernel. The subsequent sections use this representation to analyze the model’s behavior.

## 4 Information-Theoretic Analysis via Categorical Metrics

The MC framework allows us to define principled metrics for internal analysis and to formally reason about information flow. We focus on two key applications that are central to the paper’s main

arguments: quantifying the information surplus exploited by speculative decoding and measuring the intrinsic stochasticity of the prediction head.

We operate within the probabilistic setting induced by a distribution  $P_{\text{ctx}}$  over input contexts, corresponding to an initial state  $p_{W_{<t}} : I \rightarrow (\mathbb{V}^*, \mathcal{B}(\mathbb{V}^*))$ . Processing this state through the composed kernels induces distributions over the hidden state  $H_t$  (state  $p_{H_t} : I \rightarrow (\mathcal{H}, \mathcal{B}(\mathcal{H}))$ ) and the next token  $W_t$  (state  $p_{W_t} : I \rightarrow (\mathbb{V}, \mathcal{P}(\mathbb{V}))$ ).

#### 4.1 Information Flow Bounds and Rationale for Speculative Decoding

The Data Processing Inequality (DPI), a cornerstone of information theory, is a native property of our framework. It provides more than just a theoretical constraint; it offers a formal, quantitative explanation for the success of modern efficiency-boosting techniques like multi-token prediction (MTP) (Gloeckle et al., 2024) and speculative decoding (e.g., EAGLE (Li et al., 2024)). These methods are predicated on the hypothesis that the hidden state  $H_t$  contains sufficient information to predict a sequence of future tokens, not just the single next token  $W_t$ . Our framework allows us to formalize and dissect this hypothesis.

Let  $W_{t:t+K-1} = (W_t, W_{t+1}, \dots, W_{t+K-1})$  be the sequence of  $K$  tokens following the context  $W_{<t}$ . The generation of this sequence from the hidden state  $H_t$  forms a Markov chain:  $H_t \rightarrow W_{t:t+K-1} \rightarrow W_t$ . The second step,  $W_{t:t+K-1} \rightarrow W_t$ , is simply a deterministic projection that extracts the first element of the sequence. Applying the Data Processing Inequality for categorical mutual information  $I_D$  (see Remark 2.6) to this processing step directly yields a chain of inequalities:

$$I_D(H_t; W_{t:t+K-1}) \geq I_D(H_t; W_t). \quad (4.1)$$

The empirical success of speculative decoding is explained by the existence of a significant information surplus, which we define as:

$$\text{Information Surplus} := I_D(H_t; W_{t:t+K-1}) - I_D(H_t; W_t) \geq 0. \quad (4.2)$$

This surplus represents the information about future tokens ( $W_{t+1}, \dots, W_{t+K-1}$ ) that is present in the state  $H_t$  but ignored by the standard single-token prediction head  $k_{\text{head}}$ . This is not merely an artifact; it is a necessity for coherent generation. The state  $H_t$  must encode thematic, syntactic, and semantic context that will influence a long span of future text, far beyond the immediate next token. The information surplus is the formal measure of this latent, long-range predictive capacity.

Our framework can model the mechanism of speculative decoding explicitly. These models introduce a second, lightweight “drafting” kernel,  $k_{\text{draft}} : (\mathcal{H}, \mathcal{B}(\mathcal{H})) \rightarrow (\mathbb{V}^K, \mathcal{P}(\mathbb{V}^K))$ , which operates in parallel with the original “verification” head, which we now denote  $k_{\text{verify}} \equiv k_{\text{head}}$ . Both kernels take the same hidden state  $H_t$  as input, creating two parallel information pathways:

$$\begin{array}{ccc} \mathbb{V}^* & \xrightarrow{k_{\text{enc}} := k_{\text{bb}} \circ k_{\text{emb}}} & \mathcal{H} & \xrightarrow{k_{\text{verify}}} & \mathbb{V} \\ & & \downarrow k_{\text{draft}} & & \\ & & \mathbb{V}^K & & \end{array}$$

This compositional diagram makes the structure clear. The framework allows for a principled analysis of these pathways. For instance, one could compare the average categorical entropy (equation (4.4))



of the two heads. It is plausible that an effective drafting head  $\bar{\mathcal{H}}_D(k_{\text{draft}})$  would be higher than that of the verification head  $\bar{\mathcal{H}}_D(k_{\text{verify}})$ , reflecting that its role is to propose a diverse set of candidate sequences, while the verifier’s role is to make a high-confidence final selection.

The key insight is that our framework moves the understanding of speculative decoding from an empirical heuristic to a principled consequence of information theory. The information surplus is not just a concept but a formal, quantifiable resource that can be measured and, in principle, optimized, transforming the abstract bound of the DPI into a concrete tool for designing more efficient language models.

## 4.2 Metric: LM Head Categorical Entropy (Prediction Stochasticity)

To quantify the intrinsic stochasticity or uncertainty associated with the final prediction step, embodied by the LM head kernel  $k_{\text{head}} : (\mathcal{H}, \mathcal{B}(\mathcal{H})) \rightarrow (\mathbb{V}, \mathcal{P}(\mathbb{V}))$ . This metric is crucial for understanding how NLL training forces the model to learn the data’s inherent randomness (Section 5). For a given input  $h \in \mathcal{H}$ , the categorical entropy quantifies the stochasticity of the output distribution  $k_{\text{head}}(h, \cdot)$ . We obtain a single summary statistic by averaging this value over the distribution of hidden states  $p_{H_t}$ .

**Definition 4.1** (Categorical Entropy of  $k_{\text{head}}$ ). The Categorical Entropy of  $k_{\text{head}}$  is defined using equation (2.2) with  $X = \mathcal{H}$ ,  $Y = \mathbb{V}$ , and  $k = k_{\text{head}}$ :

$$\mathcal{H}_D(k_{\text{head}}) : \mathcal{H} \rightarrow \mathbb{R}_{\geq 0} := D_{\mathbb{V} \otimes \mathbb{V}}(\Delta_{\mathbb{V}} \circ k_{\text{head}} \parallel (k_{\text{head}} \otimes k_{\text{head}}) \circ \Delta_{\mathcal{H}}). \quad (4.3)$$

The categorical entropy  $\mathcal{H}_D(k_{\text{head}})(h)$  measures the divergence between generating a correlated pair  $(W, W)$  versus an independent pair  $(W_1, W_2)$ , where  $W, W_1, W_2 \sim k_{\text{head}}(h, \cdot)$ . This quantifies how far the output distribution for a given  $h$  is from a deterministic point mass. To obtain a single metric for the LM head’s overall stochasticity, we compute its expectation with respect to the hidden state distribution  $p_{H_t}$ :

$$\bar{\mathcal{H}}_D(k_{\text{head}}; p_{H_t}) := \mathbb{E}_{h \sim p_{H_t}} \left[ D_{\mathbb{V} \otimes \mathbb{V}} \left( \sum_{w \in \mathbb{V}} k_{\text{head}}(h, \{w\}) \delta_{(w,w)} \parallel k_{\text{head}}(h, \cdot) \otimes k_{\text{head}}(h, \cdot) \right) \right]. \quad (4.4)$$

**Interpretation.** This metric measures the intrinsic conditional stochasticity of the LM head mapping. If  $k_{\text{head}}$  were deterministic (i.e., for each  $h$ , it mapped to a single specific  $w_h$ , so  $p_h = \delta_{w_h}$ ), then both measures inside the divergence would be  $\delta_{(w_h, w_h)}$ , and the entropy would be  $D(\delta_{(w_h, w_h)} \parallel \delta_{(w_h, w_h)}) = 0$ . A higher value of  $\bar{\mathcal{H}}_D(k_{\text{head}}; p_{H_t})$  indicates greater average uncertainty or “spread” in the output distribution  $p_h = k_{\text{head}}(h, \cdot)$ , meaning the kernel is inherently more stochastic. It quantifies how far the prediction process is from a deterministic assignment, measured in the geometry of  $\mathbb{V} \otimes \mathbb{V}$  induced by  $D$ .

For the specific case  $D = D_{\text{KL}}$ , the average categorical entropy  $\bar{\mathcal{H}}_{D_{\text{KL}}}(k_{\text{head}}; p_{H_t})$  is closely related to the average conditional Shannon entropy  $H(W_t | H_t)$ :

$$\mathbb{E}_{h \sim p_{H_t}} [H(k_{\text{head}}(h, \cdot))] = \mathbb{E}_{h \sim p_{H_t}} \left[ - \sum_{w \in \mathbb{V}} p_h(w) \log p_h(w) \right]. \quad (4.5)$$

Both measures capture the average uncertainty in the next-token prediction given the hidden state. As discussed in [Section 5](#), minimizing NLL implicitly drives the model to match this intrinsic stochasticity present in the data.

## 5 Pretraining Objective, Compression, and Learning Intrinsic Stochasticity

A central question surrounding large language models is why the seemingly simple auto-regressive objective of next-token prediction, trained via minimizing cross-entropy loss (equivalently, negative log-likelihood or NLL), yields such powerful and versatile capabilities, often exhibiting behaviors associated with understanding and reasoning. The framework of Markov Categories and categorical entropy provides a lens through which to interpret this phenomenon, connecting it to fundamental ideas about compression and learning the inherent stochasticity of the data generating process.

Let  $k_{\text{data}} : (\mathbb{V}^*, \mathcal{B}(\mathbb{V}^*)) \rightarrow (\mathbb{V}, \mathcal{P}(\mathbb{V}))$  be the (potentially unknown) Markov kernel representing the true data-generating process, such that  $k_{\text{data}}(\mathbf{w}_{<t}, \cdot)$  corresponds to the true conditional probability measure  $P_{\text{data}}(\cdot | \mathbf{w}_{<t})$  on the vocabulary  $\mathbb{V}$ . Let  $p_{W_{<t}}$  denote the marginal probability measure on the context space  $(\mathbb{V}^*, \mathcal{B}(\mathbb{V}^*))$ , derived from the underlying joint distribution  $P_{\text{data}}$  over sequences observed in the training corpus.

The standard pretraining objective for an AR LM parameterized by  $\theta$  is to minimize the negative log-likelihood (NLL) of the next token  $w_t$  given the preceding context  $\mathbf{w}_{<t}$ , averaged over the training data distribution  $P_{\text{data}}$ . This is equivalent to minimizing the cross-entropy:

$$L_{\text{CE}}(\theta) = -\mathbb{E}_{(\mathbf{w}_{<t}, w_t) \sim P_{\text{data}}} [\log P_{\theta}(w_t | \mathbf{w}_{<t})] \quad (5.1)$$

where  $P_{\theta}(w_t | \mathbf{w}_{<t})$  is the model's predicted probability. Let  $k_{\text{gen}, \theta} : (\mathbb{V}^*, \mathcal{B}(\mathbb{V}^*)) \rightarrow (\mathbb{V}, \mathcal{P}(\mathbb{V}))$  be the model's overall generation kernel,  $k_{\text{gen}, \theta} = k_{\text{head}} \circ k_{\text{bb}} \circ k_{\text{emb}}$ , such that  $P_{\theta}(\cdot | \mathbf{w}_{<t}) = k_{\text{gen}, \theta}(\mathbf{w}_{<t}, \cdot)$ . The objective can be precisely stated in terms of Kullback-Leibler (KL) divergence between the data kernel and the model kernel.

**Theorem 5.1** (NLL Minimization as Average KL Minimization). It is a well-known result in information theory that minimizing the cross-entropy loss  $L_{\text{CE}}(\theta)$  ([equation \(5.1\)](#)) is equivalent to minimizing the average KL divergence between the true and model conditional distributions. We state it here in the language of our framework to ground the subsequent analysis.

$$\operatorname{argmin}_{\theta} L_{\text{CE}}(\theta) = \operatorname{argmin}_{\theta} \mathcal{L}_{\text{KL}}(\theta) := \operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{w}_{<t} \sim p_{W_{<t}}} [D_{\text{KL}}(k_{\text{data}}(\mathbf{w}_{<t}, \cdot) \| k_{\text{gen}, \theta}(\mathbf{w}_{<t}, \cdot))] \quad (5.2)$$

where the expectation is taken over contexts  $\mathbf{w}_{<t}$  drawn according to the data's marginal context distribution  $p_{W_{<t}}$ . The minimum value of  $\mathcal{L}_{\text{KL}}(\theta)$  is non-negative. If the model class  $\{k_{\text{gen}, \theta} \mid \theta \in \Theta\}$  is sufficiently expressive to contain  $k_{\text{data}}$  (i.e.,  $k_{\text{data}} = k_{\text{gen}, \theta_{\text{true}}}$  for some  $\theta_{\text{true}} \in \Theta$ ), then the minimum value is 0, achieved if and only if  $k_{\text{gen}, \theta^*}(\mathbf{w}_{<t}, \cdot) = k_{\text{data}}(\mathbf{w}_{<t}, \cdot)$  for  $p_{W_{<t}}$ -almost every context  $\mathbf{w}_{<t}$ .

*Proof Sketch.* The equivalence arises from rewriting cross-entropy as

$$L_{\text{CE}}(\theta) = \mathcal{L}_{\text{KL}}(\theta) + H(W_t | W_{<t})_{\text{data}},$$

where the conditional entropy  $H(W_t|W_{<t})_{\text{data}} = \mathbb{E}_{\mathbf{w}_{<t}}[-\sum_{w_t} P_{\text{data}}(w_t|\mathbf{w}_{<t}) \log P_{\text{data}}(w_t|\mathbf{w}_{<t})]$  is independent of  $\theta$ . Since  $H(W_t|W_{<t})_{\text{data}}$  is constant during optimization, minimizing  $L_{\text{CE}}(\theta)$  is identical to minimizing  $\mathcal{L}_{\text{KL}}(\theta)$ . The non-negativity and condition for achieving zero follow from the fundamental properties of KL divergence. (Full proof in Appendix A.1).  $\square$

This theorem frames NLL training as driving the model kernel  $k_{\text{gen},\theta}$  to match the data kernel  $k_{\text{data}}$ . The connection to compression arises from Shannon’s source coding theorem. The minimal average code length required to losslessly encode the next token  $w_t$ , given the context  $\mathbf{w}_{<t}$  and using an optimal code based on the true distribution  $P_{\text{data}}(\cdot|\mathbf{w}_{<t})$ , is the conditional Shannon entropy  $H(W_t|W_{<t})_{\text{data}}$ . The cross-entropy loss  $L_{\text{CE}}(\theta)$  achieved by the model represents the average code length when using a code based on the model’s distribution  $P_\theta(\cdot|\mathbf{w}_{<t})$ . Therefore, minimizing NLL (theorem 5.1) is equivalent to finding a model that provides the most efficient compression of the training data sequences, achieving an average code length that approaches the theoretical minimum  $H(W_t|W_{<t})_{\text{data}}$ . The widely discussed hypothesis that “compression implies understanding” posits that achieving high compression rates on complex data like natural language necessitates learning the underlying structure, rules, and statistical regularities, which may manifest as emergent capabilities.

Beyond matching the predictive distributions point-wise on average, successful NLL training implies that the model also learns to replicate the intrinsic stochasticity or uncertainty inherent in the data generation process at the prediction step. Within our framework, this intrinsic conditional stochasticity can be quantified using the concept of average categorical entropy (equation (4.4)). Recall that for a kernel  $k : X \rightarrow Y$  and an input distribution  $p_X$ , the average categorical entropy with respect to divergence  $D$  is:

$$\bar{\mathcal{H}}_D(k; p_X) := \mathbb{E}_{x \sim p_X} \left[ D_{Y \otimes Y} \left( \sum_{y \in Y} k(x, \{y\}) \delta_{(y,y)} \parallel k(x, \cdot) \otimes k(x, \cdot) \right) \right]. \quad (5.3)$$

This measures the average divergence between deterministically copying the output  $y \sim k(x, \cdot)$  versus generating two independent outputs  $y_1, y_2 \sim k(x, \cdot)$ . It quantifies the average “spread” or non-determinism of the kernel  $k$ .

Let  $k_{\text{head},\theta} : \mathcal{H} \rightarrow (\mathbb{V}, \mathcal{P}(\mathbb{V}))$  be the LM head kernel corresponding to parameters  $\theta$ . Let  $p_{H_t,\theta}$  be the distribution over hidden states  $h_t \in \mathcal{H}$  induced by processing contexts  $\mathbf{w}_{<t} \sim p_{W_{<t}}$  through the model’s encoder  $k_{\text{bb}} \circ k_{\text{emb}}$  (parameterized by  $\theta$ ). The following theorem formalizes the idea that convergence in average KL divergence implies convergence in the learned average intrinsic stochasticity.

**Theorem 5.2** (Convergence of Average Categorical Entropy via NLL Minimization). Let  $D$  be a statistical divergence (e.g., an  $f$ -divergence with  $f$  differentiable at 1,  $f(1) = 0$ , or  $d_{\text{TV}}$ ) defined on  $\mathcal{P}(\mathbb{V} \times \mathbb{V})$ . Assume a sequence of model parameters  $\theta_n$  achieves convergence in the training objective, such that the average KL divergence  $\mathcal{L}_{\text{KL}}(\theta_n) \rightarrow \inf_{\theta} \mathcal{L}_{\text{KL}}(\theta)$  as  $n \rightarrow \infty$ . Let  $k_n = k_{\text{gen},\theta_n}$  be the corresponding sequence of model kernels, and  $k_{\text{head},n} = k_{\text{head},\theta_n}$  be the LM head kernels. Assume this convergence implies:

- (i) Pointwise convergence of the head kernel’s output distributions:  $k_{\text{head},n}(h, \cdot) \rightarrow k_{\text{head},\theta^*}(h, \cdot)$  in a suitable topology on  $\mathcal{P}(\mathbb{V})$  (e.g., in total variation, which is strong for finite  $\mathbb{V}$ ) for  $p_{H_t,\theta^*}$ -almost every  $h$ , where  $\theta^*$  minimizes  $\mathcal{L}_{\text{KL}}(\theta)$ .

- (ii) Weak convergence of the induced hidden state distributions:  $p_{H_t, \theta_n} \Rightarrow p_{H_t, \theta^*}$ .
- (iii) The function  $\Psi_D(h, p) := D_{\mathbb{V} \otimes \mathbb{V}}(\sum_w p(w) \delta_{(w, w)} \| p \otimes p)$  is continuous and bounded as a function of  $p = k_{\text{head}}(h, \cdot)$  for  $h$  in the support of  $p_{H_t, \theta^*}$ , with respect to the convergence topology in (i). (This holds for standard divergences like KL and TV on finite  $\mathbb{V}$ ).

Then, the average categorical entropy of the learned LM head converges to that of the optimal LM head:

$$\lim_{n \rightarrow \infty} \bar{\mathcal{H}}_D(k_{\text{head}, n}; p_{H_t, \theta_n}) = \bar{\mathcal{H}}_D(k_{\text{head}, \theta^*}; p_{H_t, \theta^*}). \quad (5.4)$$

Furthermore, if the model class is sufficiently expressive such that the optimal model kernel  $k_{\text{gen}, \theta^*}$  matches the data kernel  $k_{\text{data}}$  (i.e.,  $\mathcal{L}_{\text{KL}}(\theta^*) = 0$ ), then the learned average categorical entropy approximates that of the implicit final stochastic step of the true data generating process. Assuming  $k_{\text{data}}$  admits a similar factorization with a final stochastic kernel  $k_{\text{head}, \text{data}}$  acting on some “true” state  $h_{\text{data}}$ , then:

$$\bar{\mathcal{H}}_D(k_{\text{head}, \theta^*}; p_{H_t, \theta^*}) \approx \bar{\mathcal{H}}_D(k_{\text{head}, \text{data}}; p_{H_t, \text{data}}). \quad (5.5)$$

*Proof Sketch.* The average categorical entropy is an expectation:

$$\bar{\mathcal{H}}_D(k_{\text{head}, n}; p_{H_t, \theta_n}) = \mathbb{E}_{h \sim p_{H_t, \theta_n}} [\Psi_D(h, k_{\text{head}, n}(h, \cdot))].$$

The assumptions ensure that the random variable inside the expectation converges in distribution. Specifically, weak convergence of  $p_{H_t, \theta_n}$  (ii) combined with the pointwise convergence of  $k_{\text{head}, n}(h, \cdot)$  (i) and the continuity of  $\Psi_D$  (iii) allow us to apply variants of the continuous mapping theorem or dominated convergence theorem (leveraging the assumed boundedness in (iii) and the fact that probability measures are bounded). This yields the convergence of the expectation to  $\mathbb{E}_{h \sim p_{H_t, \theta^*}} [\Psi_D(h, k_{\text{head}, \theta^*}(h, \cdot))] = \bar{\mathcal{H}}_D(k_{\text{head}, \theta^*}; p_{H_t, \theta^*})$ . The final approximation holds if  $k_{\text{gen}, \theta^*} = k_{\text{data}}$ , which implies  $k_{\text{head}, \theta^*} \approx k_{\text{head}, \text{data}}$  and  $p_{H_t, \theta^*} \approx p_{H_t, \text{data}}$  under reasonable assumptions about the factorization. (Full proof in Appendix A.2).  $\square$

Theorem 5.2 provides a formal basis for the claim that NLL training compels the model to learn not just the most likely next token, but also the degree of uncertainty or stochasticity associated with that prediction, as dictated by the data. By minimizing the average KL divergence  $\mathcal{L}_{\text{KL}}(\theta)$ , the model  $k_{\text{gen}, \theta}$  must align its output distributions  $k_{\text{gen}, \theta}(\mathbf{w}_{< t}, \cdot)$  with the data distributions  $k_{\text{data}}(\mathbf{w}_{< t}, \cdot)$ . This alignment necessarily includes matching the “shape” or “spread” of these distributions, which is precisely what is quantified by the average categorical entropy  $\bar{\mathcal{H}}_D$ . The parameters  $\theta$  and the compositional structure  $k_{\text{head}} \circ k_{\text{bb}} \circ k_{\text{emb}}$  thus become a compressed representation capturing both the predictive dependencies and the inherent conditional randomness of the language source. This suggests that learning the correct level of stochasticity is an integral part of the compression process driven by the NLL objective, contributing to the model’s ability to generate realistic and diverse text sequences.

## 6 Information Geometry of Representation and Prediction Spaces

The Markov Category framework, particularly  $(\text{Stoch}, D)$  enriched with a divergence like  $D_{\text{KL}}$ , provides a natural bridge to Information Geometry (Amari and Nagaoka, 2000; Perrone, 2023b).

This allows for a geometric analysis of the spaces involved in AR language modeling, particularly the representation space  $\mathcal{H}$  and the space of next-token distributions  $\mathcal{P}(\mathbb{V})$ .

The space  $\mathcal{P}(\mathbb{V})$  of probability distributions over the finite vocabulary  $\mathbb{V}$  forms a  $(|\mathbb{V}| - 1)$ -dimensional simplex  $\Delta^{|\mathbb{V}|-1}$ . This space possesses a well-defined Riemannian geometry induced by the Fisher-Rao information metric  $g^{\text{FR}}$ , whose components in a local coordinate system  $\xi = (\xi_1, \dots, \xi_{|\mathbb{V}|-1})$  for a distribution  $p_\xi \in \mathcal{P}(\mathbb{V})$  are given by:

$$g_{ij}^{\text{FR}}(\xi) = \sum_{w \in \mathbb{V}} p_\xi(w) \frac{\partial \log p_\xi(w)}{\partial \xi_i} \frac{\partial \log p_\xi(w)}{\partial \xi_j} = \mathbb{E}_{W \sim p_\xi} \left[ \frac{\partial \log p_\xi(W)}{\partial \xi_i} \frac{\partial \log p_\xi(W)}{\partial \xi_j} \right]. \quad (6.1)$$

This metric quantifies the local distinguishability between nearby probability distributions, measuring the distance in terms of expected squared log-likelihood ratio gradients. The geometry of  $\mathcal{P}(\mathbb{V})$  also includes dual affine connections ( $\pm\alpha$ -connections) related to the KL divergence, providing a richer dually flat structure (Amari and Nagaoka, 2000).

The LM Head kernel  $k_{\text{head}} : (\mathcal{H}, \mathcal{B}(\mathcal{H})) \rightarrow (\mathbb{V}, \mathcal{P}(\mathbb{V}))$  corresponds to a deterministic mapping from a hidden state  $h \in \mathcal{H} \cong \mathbb{R}^{d_{\text{model}}}$  to a probability distribution  $p_h := k_{\text{head}}(h, \cdot) \in \mathcal{P}(\mathbb{V})$ . Let  $g_{\text{head}} : \mathcal{H} \rightarrow \mathcal{P}(\mathbb{V})$  denote this mapping,  $p_h = g_{\text{head}}(h)$ . Typically, this involves a linear layer followed by softmax:  $g_{\text{head}}(h) = \text{softmax}(Wh)$  where  $W \in \mathbb{R}^{|\mathbb{V}| \times d_{\text{model}}}$ . This mapping  $g_{\text{head}}$  allows us to pull back the geometric structure from  $\mathcal{P}(\mathbb{V})$  onto the representation space  $\mathcal{H}$ .

Specifically, the Fisher-Rao metric  $g^{\text{FR}}$  on  $\mathcal{P}(\mathbb{V})$  induces a (generally degenerate) Riemannian metric tensor  $g^* = g_{\text{head}}^* g^{\text{FR}}$  on  $\mathcal{H}$ . At a point  $h \in \mathcal{H}$ , the components of this pullback metric are given by:

$$g_{ab}^*(h) = \sum_{i,j} g_{ij}^{\text{FR}}(g_{\text{head}}(h)) \frac{\partial (g_{\text{head}}(h))_i}{\partial h_a} \frac{\partial (g_{\text{head}}(h))_j}{\partial h_b}, \quad a, b \in \{1, \dots, d_{\text{model}}\}, \quad (6.2)$$

where  $h_a, h_b$  are coordinates of  $h \in \mathcal{H}$ , and  $(g_{\text{head}}(h))_i, (g_{\text{head}}(h))_j$  represent local coordinates of the output distribution  $p_h \in \mathcal{P}(\mathbb{V})$  (e.g., probabilities of specific tokens, possibly excluding one due to the sum-to-one constraint). The term  $\frac{\partial (g_{\text{head}}(h))_i}{\partial h_a}$  is the Jacobian of the LM head map  $g_{\text{head}}$  evaluated at  $h$ .

Let  $J(h)$  denote this Jacobian matrix ( $|\mathbb{V}| - 1 \times d_{\text{model}}$  or  $|\mathbb{V}| \times d_{\text{model}}$  depending on coordinates). Then  $g^*(h) = J(h)^\top g^{\text{FR}}(g_{\text{head}}(h)) J(h)$ . The significance of this pullback metric  $g^*$  lies in its connection to the local distinguishability of output distributions under perturbations of the input hidden state, as measured by divergences like KL divergence.

**Theorem 6.1** (Pullback Metric and Local Divergence). Let  $g_{\text{head}} : \mathcal{H} \rightarrow \mathcal{P}(\mathbb{V})$  be the smooth map corresponding to the LM head kernel. Let  $h \in \mathcal{H}$  and  $v \in T_h \mathcal{H} \cong \mathcal{H}$ . Consider the distributions  $p_h = g_{\text{head}}(h)$  and  $p_{h+\epsilon v} = g_{\text{head}}(h + \epsilon v)$  for small  $\epsilon$ . The KL divergence between these output distributions, for small  $\epsilon$ , is locally approximated by the quadratic form defined by the pullback metric  $g^*(h)$ :

$$D_{\text{KL}}(p_{h+\epsilon v} \parallel p_h) = \frac{1}{2} \epsilon^2 g^*(h)(v, v) + O(\epsilon^3) \quad (6.3)$$

where  $g^*(h)(v, v) = \sum_{a,b=1}^{d_{\text{model}}} g_{ab}^*(h) v_a v_b$ . A similar relationship holds for symmetric KL divergence and, more generally, for any  $f$ -divergence  $D_f$  where  $f$  is sufficiently smooth around 1 with  $f''(1) > 0$ .

*Proof.* Let  $\xi$  be a local coordinate system for  $\mathcal{P}(\mathbb{V})$  around  $p_h$ . The KL divergence between two nearby distributions  $p_\xi$  and  $p_{\xi'}$  can be expanded around  $p_\xi$  as (Amari and Nagaoka, 2000):

$$D_{\text{KL}}(p_{\xi'} \| p_\xi) = \frac{1}{2} \sum_{i,j} g_{ij}^{\text{FR}}(\xi) (\xi'_i - \xi_i) (\xi'_j - \xi_j) + O(\|\xi' - \xi\|^3).$$

Let  $\xi(h)$  denote the coordinates of  $p_h = g_{\text{head}}(h)$ . For  $p_{h+\epsilon v}$ , the coordinates are  $\xi(h + \epsilon v)$ . By Taylor expansion in  $\epsilon$ :

$$\xi_i(h + \epsilon v) = \xi_i(h) + \epsilon \sum_{a=1}^{d_{\text{model}}} \frac{\partial \xi_i}{\partial h_a}(h) v_a + O(\epsilon^2).$$

Thus,  $\xi_i(h + \epsilon v) - \xi_i(h) = \epsilon J_{ia}(h) v_a + O(\epsilon^2)$ , where  $J(h)$  is the Jacobian matrix of the map  $h \mapsto \xi(h)$  (i.e., the Jacobian of  $g_{\text{head}}$  in local coordinates  $\xi$ ). Substituting this into the KL expansion:

$$\begin{aligned} D_{\text{KL}}(p_{h+\epsilon v} \| p_h) &= \frac{1}{2} \sum_{i,j} g_{ij}^{\text{FR}}(\xi(h)) \left( \epsilon \sum_a J_{ia}(h) v_a \right) \left( \epsilon \sum_b J_{jb}(h) v_b \right) + O(\epsilon^3) \\ &= \frac{1}{2} \epsilon^2 \sum_{a,b} \left( \sum_{i,j} J_{ia}(h) g_{ij}^{\text{FR}}(\xi(h)) J_{jb}(h) \right) v_a v_b + O(\epsilon^3) \\ &= \frac{1}{2} \epsilon^2 \sum_{a,b} (J(h)^\top g^{\text{FR}}(\xi(h)) J(h))_{ab} v_a v_b + O(\epsilon^3). \end{aligned}$$

The term  $J(h)^\top g^{\text{FR}}(\xi(h)) J(h)$  is precisely the matrix representation of the pullback metric  $g^*(h)$  in the standard coordinates of  $\mathcal{H} \cong \mathbb{R}^{d_{\text{model}}}$ , derived from equation (6.2). Thus,  $D_{\text{KL}}(p_{h+\epsilon v} \| p_h) = \frac{1}{2} \epsilon^2 g^*(h)(v, v) + O(\epsilon^3)$ . The result for other well-behaved  $f$ -divergences follows from their similar second-order expansion involving  $g^{\text{FR}}$ .  $\square$

This theorem formally establishes that the pullback metric  $g^*$  measures how sensitive the output distribution  $p_h$  is to infinitesimal changes in the hidden state  $h$ , where sensitivity is gauged by the local divergence (specifically, KL divergence, relating to the Fisher-Rao metric) in the output space  $\mathcal{P}(\mathbb{V})$ .

**Remark 6.2** (Pullback Metric as Expected Score Outer Product). The Fisher-Rao metric  $g^{\text{FR}}$  is the expected outer product of the score function  $\nabla_\xi \log p_\xi(W)$ . This property pulls back to  $\mathcal{H}$ . Let  $p_h(w) = k_{\text{head}}(h, \{w\})$ . The score vector for token  $w$  with respect to the representation is  $\nabla_h \log p_h(w) \in \mathcal{H}$ . The pullback metric tensor is precisely the expected outer product of this score:

$$g^*(h) = \mathbb{E}_{W \sim p_h} [(\nabla_h \log p_h(W)) (\nabla_h \log p_h(W))^\top]. \quad (6.4)$$

This directly connects the information geometry of  $\mathcal{H}$  to the sensitivity of log-probabilities to changes in the representation  $h$ . This score vector  $\nabla_h \log p_h(W)$  is analogous to that used in score-based generative models, but here taken with respect to the conditioning variable  $h$ .

The rank of the pullback metric depends on the dimensions of the spaces involved.



**Proposition 6.3** (Rank of the Pullback Metric). The rank of the pullback Fisher-Rao metric  $g^*(h)$  at a point  $h \in \mathcal{H}$  is bounded by the minimum of the representation dimension and the dimension of the probability simplex:

$$\text{rank}(g^*(h)) \leq \min(d_{\text{model}}, |\mathbb{V}| - 1). \quad (6.5)$$

*Proof.* The pullback metric  $g^*(h)$  is defined as  $g^*(h) = J(h)^\top g^{\text{FR}}(g_{\text{head}}(h))J(h)$ , where  $J(h)$  is the Jacobian of the map  $g_{\text{head}} : \mathcal{H} \rightarrow \mathcal{P}(\mathbb{V})$  (represented in appropriate local coordinates). The dimension of  $\mathcal{H}$  is  $d_{\text{model}}$ , and the dimension of  $\mathcal{P}(\mathbb{V})$  is  $d_{\text{prob}} = |\mathbb{V}| - 1$ . The Jacobian  $J(h)$  is a  $d_{\text{prob}} \times d_{\text{model}}$  matrix. The Fisher-Rao metric  $g^{\text{FR}}$  at  $g_{\text{head}}(h)$  is a  $d_{\text{prob}} \times d_{\text{prob}}$  positive definite matrix (and thus has rank  $d_{\text{prob}}$ ). Using the property that  $\text{rank}(A^\top B A) = \text{rank}(A)$  if  $B$  is positive definite, we have  $\text{rank}(g^*(h)) = \text{rank}(J(h))$ . The rank of a matrix is bounded by its dimensions, so

$$\text{rank}(g^*(h)) \leq \min(d_{\text{model}}, d_{\text{prob}}) = \min(d_{\text{model}}, |\mathbb{V}| - 1).$$

□

## 6.1 Interpretation and Implications

This geometric perspective provides several insights:

- **Sensitivity Analysis:** The quadratic form  $g^*(h)(v, v)$  quantifies the local distinguishability (via KL divergence, [equation \(6.3\)](#)) between the output distributions  $p_h$  and  $p_{h+ev}$ . It measures how sensitive the model’s prediction is to perturbations of the hidden state  $h$  in a direction  $v$ . Directions  $v$  with large  $g^*(h)(v, v)$  correspond to changes in  $h$  that significantly alter the output distribution.
- **Functional Anisotropy:** In modern LMs, the representation dimension is much smaller than the vocabulary size ( $d_{\text{model}} \ll |\mathbb{V}|$ ). By Prop. 6.3,  $\text{rank}(g^*(h)) \leq d_{\text{model}}$ . This means the pullback metric  $g^*(h)$  is necessarily a degenerate tensor on  $\mathcal{H}$ , but its rank is limited by the representation’s capacity. More importantly, the metric is highly *anisotropic*. This geometrically formalizes “functional anisotropy”: the learned mapping is highly sensitive in some directions of  $\mathcal{H}$  (those with large eigenvalues for  $g^*(h)$ ) but relatively insensitive in others. This is not a flaw but a feature, indicating learned specialization.
- **Spectrum of Sensitivity:** The eigenvalues and eigenvectors of the matrix for  $g^*(h)$  reveal the principal directions of predictive sensitivity in  $\mathcal{H}$ . Directions with large eigenvalues are those where small changes in  $h$  induce large changes (geometrically measured by  $g^{\text{FR}}$ ) in the predicted distribution  $p_h$ . These are the directions the model has learned are most important for prediction.

If two conditional distributions  $p_{H_t|s_1}$  and  $p_{H_t|s_2}$  are supported on regions of  $\mathcal{H}$  that map to distinct regions in  $\mathcal{P}(\mathbb{V})$  via  $g_{\text{head}}$ , the distance between these regions in  $\mathcal{P}(\mathbb{V})$  (measured, e.g., by integrated Fisher-Rao distance or KL divergence) contributes to  $\text{RepDiv}_D(s_1 \| s_2)$ . The geometry induced by  $g^*$  characterizes the local separation capability. Training aims to shape the encoder ( $k_{\text{bb}} \circ k_{\text{emb}}$ ) and the LM head  $k_{\text{head}}$  such that contexts with different predictive futures ( $P_{\text{data}}(\cdot|s_1)$  vs  $P_{\text{data}}(\cdot|s_2)$ ) are mapped to representations  $h_t$  whose images under  $g_{\text{head}}$  are appropriately separated in  $\mathcal{P}(\mathbb{V})$ , implicitly structuring the manifold  $(\mathcal{H}, g^*)$ .

## 7 NLL as Implicit Spectral Contrastive Learning

A central thesis of this work is that the simple objective of minimizing the negative log-likelihood (NLL) of the next token (equation (5.2)) implicitly functions as a powerful form of contrastive learning. While lacking the explicit positive/negative pairs of standard contrastive methods, we prove that NLL optimization inherently structures the learned representation space  $\mathcal{H}$  according to predictive similarity. It achieves this by implicitly solving a spectral objective that aligns the geometry of representations with the underlying predictive structure of the data  $P_{\text{data}}(\cdot|x)$ , a principle we formalize by connecting NLL to the eigenspectrum of a predictive similarity operator (HaoChen et al., 2021; Tan et al., 2024).

Let  $f_{\text{enc}} : \mathbb{V}^* \rightarrow \mathcal{H}$  denote the deterministic encoder mapping a context sequence  $x = \mathbf{w}_{<t}$  to its hidden representation  $h_x = f_{\text{enc}}(x)$ , implemented by the composition  $k_{\text{bb}} \circ k_{\text{emb}}$ . Let  $g_{\text{head}} : \mathcal{H} \rightarrow \mathcal{P}(\mathbb{V})$  be the deterministic mapping from the hidden state to the next-token distribution, corresponding to the LM head kernel  $k_{\text{head}}$ , such that  $p_{\theta}(\cdot|x) = g_{\text{head}}(h_x)$ . The training objective is to minimize the expected KL divergence over the context distribution  $\mu_{\text{ctx}} = p_{W_{<t}}$ :

$$\mathcal{L}(\theta) = \mathbb{E}_{x \sim \mu_{\text{ctx}}} [D_{\text{KL}}(P_{\text{data}}(\cdot|x) \| g_{\text{head}}(f_{\text{enc}}(x)))] \quad (7.1)$$

where  $P_{\text{data}}(\cdot|x)$  represents the true conditional distribution of the next token given context  $x$ , assumed to be derived from the data-generating process.

Successful optimization of  $\mathcal{L}(\theta)$  drives the model’s output distribution  $p_{\theta}(\cdot|x) = g_{\text{head}}(h_x)$  towards the target distribution  $P_{\text{data}}(\cdot|x)$  in the sense of minimizing average KL divergence. As we argue below, this fundamental requirement indirectly imposes geometric constraints on the distribution of representations  $h_x = f_{\text{enc}}(x)$  in  $\mathcal{H}$ .

### 7.1 Constraint on Output Distribution Approximation

Minimizing the NLL loss (equation (7.1)) directly forces the model’s predicted distribution  $p_{\theta}(\cdot|x)$  to closely approximate the target distribution  $P_{\text{data}}(\cdot|x)$ . This closeness can be measured not only by KL divergence but also by other standard metrics on probability distributions, due to well-known inequalities relating them.

**Theorem 7.1** (Output Distribution Approximation Constraint). Assume the model parameters  $\theta$  yield a small average KL divergence  $\mathcal{L}(\theta) = \mathbb{E}_{x \sim \mu_{\text{ctx}}} [D_{\text{KL}}(P_{\text{data}}(\cdot|x) \| p_{\theta}(\cdot|x))]$ , where  $p_{\theta}(\cdot|x) = g_{\text{head}}(f_{\text{enc}}(x))$ . Let  $d_{\text{out}}$  be a metric on  $\mathcal{P}(\mathbb{V})$  satisfying a Pinsker-type inequality, such as  $d_{\text{out}}(p, q)^k \leq C \cdot D_{\text{KL}}(p \| q)$  for some constants  $k, C > 0$ . Examples include Hellinger distance ( $d_H, k = 2, C = 1/2$ ) and Total Variation distance ( $d_{\text{TV}}, k = 2, C = 1$ ). Then, the expected  $d_{\text{out}}$ -distance between the model’s prediction and the true distribution is also small:

$$\mathbb{E}_{x \sim \mu_{\text{ctx}}} [d_{\text{out}}(P_{\text{data}}(\cdot|x), p_{\theta}(\cdot|x))^k] \leq C \cdot \mathcal{L}(\theta). \quad (7.2)$$

Consequently, if the model fits the data well ( $\mathcal{L}(\theta)$  is small), the distance between the model’s output distributions for any two contexts,  $d_{\text{out}}(p_{\theta}(\cdot|x), p_{\theta}(\cdot|x'))$ , must approximate the distance between the true target distributions,  $d_{\text{out}}(P_{\text{data}}(\cdot|x), P_{\text{data}}(\cdot|x'))$ . Specifically, contexts with predictively dissimilar target distributions (large LHS below) must yield model output distributions that are

also separated:

$$d_{\text{out}}(p_{\theta}(\cdot|x), p_{\theta}(\cdot|x')) \approx d_{\text{out}}(P_{\text{data}}(\cdot|x), P_{\text{data}}(\cdot|x')). \quad (7.3)$$

*Proof Sketch.* The inequality [equation \(7.2\)](#) arises directly from the assumed Pinsker-type inequality  $d_{\text{out}}(p, q)^k \leq C \cdot D_{\text{KL}}(p||q)$ . Letting  $p = P_{\text{data}}(\cdot|x)$  and  $q = p_{\theta}(\cdot|x)$ , we have  $d_{\text{out}}(P_{\text{data}}(\cdot|x), p_{\theta}(\cdot|x))^k \leq C \cdot D_{\text{KL}}(P_{\text{data}}(\cdot|x)||p_{\theta}(\cdot|x))$ . Taking the expectation over  $x \sim \mu_{\text{ctx}}$  on both sides yields the result. For the second part, the triangle inequality states  $d_{\text{out}}(p_x^{\text{data}}, p_{x'}^{\text{data}}) \leq d_{\text{out}}(p_x^{\text{data}}, p_x^{\theta}) + d_{\text{out}}(p_x^{\theta}, p_{x'}^{\theta}) + d_{\text{out}}(p_{x'}^{\theta}, p_{x'}^{\text{data}})$ . Let  $\epsilon_x = d_{\text{out}}(p_x^{\text{data}}, p_x^{\theta})$ . If  $\mathcal{L}(\theta)$  is small,  $\mathbb{E}_x[\epsilon_x^k]$  is small. By Markov's inequality,  $\mathbb{P}(\epsilon_x \geq \delta) \leq \mathbb{E}[\epsilon_x^k]/\delta^k$ , implying  $\epsilon_x$  is small with high probability for typical  $x$ . Rearranging the triangle inequality bounds  $d_{\text{out}}(p_x^{\theta}, p_{x'}^{\theta})$  by  $d_{\text{out}}(p_x^{\text{data}}, p_{x'}^{\text{data}}) \pm (\epsilon_x + \epsilon_{x'})$ , leading to [equation \(7.3\)](#). (Full proof in [Appendix A.3](#))  $\square$

This theorem formalizes the intuition that minimizing the NLL objective forces the model's predictions to mirror the structure of the true predictive distributions, specifically in terms of their pairwise distances.

## 7.2 Consequences for Representation Geometry

[Theorem 7.1](#) establishes that predictively dissimilar contexts  $x, x'$  must lead to distinct model output distributions  $p_{\theta}(\cdot|x), p_{\theta}(\cdot|x')$ . Since  $p_{\theta}(\cdot|x) = g_{\text{head}}(h_x)$  and  $p_{\theta}(\cdot|x') = g_{\text{head}}(h_{x'})$ , this requirement imposes constraints on the corresponding representations  $h_x = f_{\text{enc}}(x)$  and  $h_{x'} = f_{\text{enc}}(x')$ . Specifically,  $h_x$  and  $h_{x'}$  must differ in ways that are discernible by the head mapping  $g_{\text{head}}$ . The information geometry of the head mapping, captured by the pullback metric  $g^*(h)$  ([Section 6](#)), determines which differences in representation space are discernible.

**Corollary 7.2** (Implicit Representation Separation). Assume the model fits the data well ( $\mathcal{L}(\theta)$  is small). If two contexts  $x, x'$  have predictively dissimilar target distributions, meaning  $d_{\text{out}}(P_{\text{data}}(\cdot|x), P_{\text{data}}(\cdot|x'))$  is large, then their learned representations  $h_x = f_{\text{enc}}(x)$  and  $h_{x'} = f_{\text{enc}}(x')$  must differ in a way that is detectable by the head mapping. Specifically, the difference vector  $h_x - h_{x'}$  must have significant components along directions  $v$  of high predictive sensitivity (i.e., where the pullback metric quadratic form  $g^*(h)(v, v)$  is large).

$$d_{\text{out}}(P_{\text{data}}(\cdot|x), P_{\text{data}}(\cdot|x')) \text{ large} \implies \int_0^1 g^*(h_t)(h_x - h_{x'}, h_x - h_{x'}) dt \text{ is large}, \quad (7.4)$$

where the integral is over the path  $h_t = (1-t)h_{x'} + th_x$ . Conversely, if contexts  $x, x'$  are predictively similar ( $d_{\text{out}}(P_{\text{data}}(\cdot|x), P_{\text{data}}(\cdot|x'))$  is small), the NLL objective encourages their representations to be close along these same sensitive directions.

*Proof Sketch.* From [Theorem 7.1](#), if  $d_{\text{out}}(P_{\text{data}}(\cdot|x), P_{\text{data}}(\cdot|x'))$  is large, then  $d_{\text{out}}(g_{\text{head}}(h_x), g_{\text{head}}(h_{x'}))$  must also be large. The squared distance between two points in a Riemannian manifold is related to the integrated metric along a geodesic. For small distances in output space, we have  $d_{\text{out}}(p, q)^2 \approx D_{\text{KL}}(p||q)$ , which from [equation \(6.3\)](#) is related to the pullback metric  $g^*$ . A large distance between  $g_{\text{head}}(h_x)$  and  $g_{\text{head}}(h_{x'})$  implies a large integrated path length according to the pullback geometry, forcing  $h_x$  and  $h_{x'}$  to differ along directions where  $g^*$  is large. (Full proof in [Appendix A.4](#))  $\square$

This corollary establishes that NLL minimization implicitly acts like a contrastive learning objective: it pushes representations  $h_x, h_{x'}$  apart if their corresponding contexts are predictively dissimilar. This differential pressure based on predictive similarity forms the basis for our connection to spectral methods.

### 7.3 Predictive Similarity Kernels

The preceding analysis suggests that NLL shapes the representation geometry based on the *dissimilarity* between the true next-token distributions  $P_{\text{data}}(\cdot|x)$ . To connect this to spectral methods, which operate on similarity structures, we formalize the complementary notion of *predictive similarity*.

**Definition 7.3** (Predictive Similarity Kernel). Let  $p_x := P_{\text{data}}(\cdot|x)$  denote the true conditional distribution for context  $x$ . A predictive similarity kernel is a function  $K : \mathbb{V}^* \times \mathbb{V}^* \rightarrow \mathbb{R}_{\geq 0}$  quantifying the similarity between  $p_x$  and  $p_{x'}$ . Potential examples include:

- **Bhattacharyya Coefficient Kernel:**  $K_{\text{BC}}(x, x') := \text{BC}(p_x, p_{x'}) = \sum_{w \in \mathbb{V}} \sqrt{p_x(w)p_{x'}(w)}$ . This measures the cosine similarity between the square-root vectors  $(\sqrt{p_x(w)})_w$ . It relates directly to the Hellinger distance  $d_H^2(p_x, p_{x'}) = 2(1 - K_{\text{BC}}(x, x'))$  and defines a positive semidefinite kernel. High  $K_{\text{BC}}$  corresponds to low  $d_H$ .
- **Hellinger-based Kernel (Gaussian Kernel on  $\sqrt{p}$ ):**  $K_H(x, x') := \exp(-\beta d_H^2(p_x, p_{x'}))$  for some scale  $\beta > 0$ . This explicitly converts the Hellinger distance into a similarity measure via a Gaussian function, yielding a positive semidefinite kernel.
- **Expected Likelihood Kernel (Linear Kernel):**  $K_{\text{Lin}}(x, x') := \langle p_x, p_{x'} \rangle = \sum_{w \in \mathbb{V}} p_x(w)p_{x'}(w)$ . This is the standard linear kernel (inner product) between probability vectors  $p_x, p_{x'}$  and is positive semidefinite. High values indicate significant overlap between the distributions. It can be interpreted as the expected likelihood  $p_{x'}(W)$  under  $W \sim p_x$ .
- **KL-based Kernel:**  $K_{\text{KL}}(x, x') := \exp(-\beta S_{\text{KL}}(p_x, p_{x'}))$  where  $S_{\text{KL}}$  is a symmetrized KL divergence (e.g., Jensen-Shannon divergence) and  $\beta > 0$ . Constructing a positive semidefinite kernel directly from KL divergence requires symmetrization.

In general, high values of  $K(x, x')$  indicate high predictive similarity (i.e., small  $d_{\text{out}}(p_x, p_{x'})$ ). It is this similarity structure, derived from the data’s conditional probabilities, that we hypothesize the NLL objective implicitly captures within the geometry of  $\mathcal{H}$ .

### 7.4 Connection to Graph Laplacian and Dirichlet Energy Minimization

Consider an undirected graph where contexts  $x \in \mathbb{V}^*$  are nodes distributed according to  $\mu_{\text{ctx}}$ , and edge weights are given by a symmetric predictive similarity kernel  $K(x, x')$ . The quadratic form of the associated graph Laplacian corresponds to the Dirichlet energy, which measures how “smooth” a function  $\phi$  (e.g., a 1D projection of the representations) is over the graph.

$$\mathcal{E}_K(\phi) := \frac{1}{2} \iint K(x, x') (\phi(x) - \phi(x'))^2 \mu_{\text{ctx}}(dx) \mu_{\text{ctx}}(dx') = \langle \phi, \Delta_K \phi \rangle_{L^2(\mu_{\text{ctx}})}. \quad (7.5)$$

Spectral clustering aims to find embeddings (represented by functions  $\phi$ ) that minimize this energy subject to constraints, effectively mapping similar contexts close together. The NLL objective, through Corollary 7.2, exerts a related pressure.

Recent work (Park et al., 2024) connects in-context learning to Dirichlet energy minimization on a task-similarity graph. Our work shows this is a foundational principle of NLL pre-training itself, where similarity is defined by the intrinsic next-token distributions  $P_{\text{data}}(\cdot|x)$ .

**Proposition 7.4** (NLL and Implicit Dirichlet Energy Minimization). Let  $\phi_v(x) = \langle h_x, v \rangle$  be the projection of the representation  $h_x = f_{\text{enc}}(x)$  onto a direction  $v$  of high predictive sensitivity (where  $g^*(h)(v, v)$  is large). Minimizing the NLL loss  $\mathcal{L}(\theta)$  implicitly encourages configurations where the Dirichlet energy  $\mathcal{E}_K(\phi_v)$  is small.

*Proof Sketch.* The Dirichlet energy integrand  $K(x, x')(\phi_v(x) - \phi_v(x'))^2$  is large only if both  $K(x, x')$  and  $(\phi_v(x) - \phi_v(x'))^2$  are large. Corollary 7.2 shows NLL encourages  $(\phi_v(x) - \phi_v(x'))^2$  to be small when  $K(x, x')$  is high. Conversely, if  $K(x, x')$  is low,  $(\phi_v(x) - \phi_v(x'))^2$  is encouraged to be large, but the integrand is suppressed by the small  $K(x, x')$  term. Thus, NLL minimization discourages configurations that would lead to high Dirichlet energy. (Full proof in Appendix A.5.)  $\square$

This proposition formalizes the link: NLL pushes representations towards configurations favored by spectral clustering on the predictive similarity graph.

## 7.5 NLL as Spectral Objective

We now strengthen this connection, showing that NLL optimization is not merely analogous to spectral methods, but that it implicitly solves a spectral objective. This viewpoint is closer to the analysis in (Tan et al., 2024).

**Definition 7.5** (Predictive Similarity Operator). Let  $f_{\text{enc}} : \mathbb{V}^* \rightarrow \mathcal{H}$  be a fixed encoder, inducing a distribution  $\mu = p_{H_t}$  on  $\mathcal{H}$  via  $\mu_{\text{ctx}}$ . Let  $K(x, x')$  be the predictive similarity kernel. The predictive similarity operator  $M_K : L^2(\mathcal{H}, \mu) \rightarrow L^2(\mathcal{H}, \mu)$  is defined via its action on functions  $\psi : \mathcal{H} \rightarrow \mathbb{R}$ :

$$(M_K \psi)(h_x) \triangleq \mathbb{E}_{x' \sim \mu_{\text{ctx}}} [K(x, x') \psi(h_{x'})] = \int_{\mathbb{V}^*} K(x, x') \psi(f_{\text{enc}}(x')) \mu_{\text{ctx}}(dx'), \quad (7.6)$$

where  $h_x = f_{\text{enc}}(x)$ . This operator averages the function  $\psi$  over representations  $h_{x'}$  weighted by the predictive similarity  $K(x, x')$  between their originating contexts  $x, x'$  and the reference context  $x$ .

If  $K$  is symmetric,  $M_K$  is a compact self-adjoint operator whose eigenfunctions capture dominant patterns of predictive similarity. Our key result is that NLL training implicitly aligns the representation geometry with the eigenspace of  $M_K$ .

**Proposition 7.6** (NLL Objective as Implicit Spectral Alignment). Let the representation covariance matrix be  $\Sigma_H = \text{Cov}_{x \sim \mu_{\text{ctx}}} [h_x]$ . Let the head-weighted covariance be  $\Sigma_{H, g^*} = \mathbb{E}_{x \sim \mu_{\text{ctx}}} [h_x h_x^\top \otimes g^*(h_x)]$ , which captures variance along predictively sensitive directions. Minimizing the NLL loss  $\mathcal{L}(\theta)$  implicitly encourages a solution to an objective akin to maximizing the alignment between the representations and the similarity operator, subject to constraints on predictive distortion. It favors representations where directions of high variance align with the leading eigenfunctions of  $M_K$ , particularly after filtering for predictive relevance by  $g^*$ .

*Proof Sketch.* More formally, under a standard log-linear model for the head  $g_{\text{head}}$ , the NLL loss can be decomposed into an alignment term and a uniformity term. The alignment term explicitly encourages the representation  $h_x$  to align with the expected prototype vector  $\bar{g}_x = \mathbb{E}_{w \sim p_x} [g(w)]$ .

This implies that the geometry of the representations should match the geometry of the expected prototypes:  $\langle h_x, h_{x'} \rangle \propto \langle \bar{g}_x, \bar{g}_{x'} \rangle$ . If we define the predictive similarity kernel as  $K(x, x') = \langle \bar{g}_x, \bar{g}_{x'} \rangle$ , then NLL optimization encourages the Gram matrix of representations to align with the kernel matrix  $K$ . This is a spectral objective, as it drives the covariance of the representations to align with the eigenspectrum of the similarity operator  $M_K$ . This effectively performs a spectral decomposition of the predictive similarity graph, where similarity is defined by the expected next token. (A full derivation is provided in Appendix A.6).  $\square$

In summary, NLL optimization does not merely resemble spectral methods; it implicitly is a spectral method. The core objective of matching conditional probabilities imposes geometric constraints on  $\mathcal{H}$  that are equivalent to finding a low-dimensional embedding of the predictive similarity graph. This provides a rigorous, first-principles explanation for why NLL training produces representations that are so semantically and structurally organized.

## 8 Related Work

The theoretical understanding of representation learning in deep neural networks has been advanced along several parallel, yet largely disconnected, fronts. A significant challenge lies in the absence of a unified mathematical language capable of connecting a model’s compositional architecture and its training dynamics to the emergent geometric structure of its learned representations. This work is situated at the confluence of these disparate research programs, aiming to synthesize the algebraic, compositional perspective of categorical probability with the metric, differential-geometric view of information geometry. We structure our review of the related literature into two parts. First, we introduce the foundational languages of probability that our framework unifies: the synthetic view of probability rooted in category theory and the metric view rooted in information geometry. Second, we survey the tools and concepts used to analyze the geometry of learned representations, focusing on the training objectives that guide learning, the spectral methods used to measure the resulting geometry, and the optimization mechanisms that shape it.

### 8.1 The Languages of Probability

**The Synthetic View: Probability as a Category.** A burgeoning field of research seeks to reformulate probability theory on a more abstract, algebraic foundation using the language of category theory (Baez et al., 2016; Fong and Spivak, 2018). This synthetic approach, in contrast to the classical analytic approach built upon measure theory, aims to derive probabilistic concepts from a small set of powerful axioms (Fritz, 2020a; Perrone, 2023b). The central object of study in this domain is the *Markov category* (Cho and Jacobs, 2019; Fritz, 2020a). Formally, a Markov category is a symmetric monoidal category where each object is equipped with a commutative comonoid structure, consisting of morphisms that represent the abstract operations of copying and discarding information (Cho and Jacobs, 2019; Fritz, 2020a; Perrone, 2023b). The morphisms in such a category are interpreted as stochastic maps, or Markov kernels, which are probabilistic mappings between objects (Baez et al., 2016; Pardo-Guerra et al., 2025).

The pioneering work of Fritz (2020a) has established Markov categories as a robust framework for synthetic probability and statistics. A key advantage of this formalism is its generality; it provides a uniform treatment of vastly different probabilistic settings. For instance, the category **FinStoch**,



**Table 1** A Comparative Taxonomy of Theoretical Frameworks for Representation Learning. This table situates our proposed framework relative to existing theoretical paradigms, highlighting the unique synthesis of algebraic and geometric perspectives it offers to analyze the representation space directly.

Framework	Core Mathematical Object	Primary Domain of Analysis	Key Insight/Contribution	Limitation Addressed by Our Work	Key Citations
<b>Information Geometry (IG)</b>	Riemannian Manifold (Fisher-Rao Metric)	Model Parameter Space	Provides a metric for the “distance” between models; explains optimization via natural gradients.	Lacks a compositional, algebraic structure; difficult to apply to the dynamics of representations themselves.	<a href="#">Amari (2010)</a>
<b>Categorical Probability</b>	Markov Category (Symmetric Monoidal Category + Comonoids)	Abstract Probabilistic Systems	Provides a high-level, compositional language for probability, unifying disparate theories via algebraic axioms.	Inherently algebraic and non-metric; does not natively describe the shape or geometry of representation spaces.	<a href="#">Fritz (2020a)</a> , <a href="#">Perrone (2023b)</a> , <a href="#">Cho and Jacobs (2019)</a>
<b>Information Bottleneck (IB)</b>	Mutual Information	Information-Theoretic Channels	Defines representation learning as an optimal trade-off between compression and predictive relevance.	A high-level principle, not a constructive framework; does not specify the geometric mechanism for achieving the bottleneck.	<a href="#">Tishby et al. (2000a)</a>
<b>Spectral Graph Theory</b>	Graph Laplacian / Dirichlet Energy	Graph-structured Data / Activations	Quantifies the “smoothness” of functions on a graph; explains over-smoothing in GNNs and representation reorganization in LLMs.	A measurement tool, not a full-fledged theoretical framework; describes <i>what</i> geometry emerges but not <i>why</i> .	<a href="#">Yuan (2023)</a> ; <a href="#">Park et al. (2024)</a>
<b>Implicit Bias Theory</b>	Optimization Dynamics (Gradient Flow)	Model Weights & Activations	Explains generalization in overparameterized models; frames NLL as an implicit contrastive, geometry-sculpting process.	Provides a mechanism but lacks a unified language to describe the resulting geometric and probabilistic structures.	<a href="#">Gunasekar et al. (2018)</a>
<b>Our Framework (Proposed)</b>	<b>Geometric Markov Category</b>	<b>Representation Space</b>	<b>Unifies the algebraic structure of probability with the metric geometry of representations.</b>	-	<b>This Paper</b>

with finite sets as objects and stochastic matrices as morphisms, and the category `BorelStoch`, with standard Borel spaces as objects and their corresponding Markov kernels as morphisms, are both canonical examples of Markov categories ([Fritz, 2020a](#)). This high level of abstraction allows for the proof of fundamental statistical theorems—such as the Fisher-Neyman factorization theorem and Kolmogorov’s zero-one law—in a purely diagrammatic and synthetic manner, avoiding the low-level complexities of measure theory ([Fritz, 2020a](#); [Fritz and Rischel, 2020](#)). As [Fritz \(2020a\)](#) argues, relying on measure theory is akin to programming in machine code, whereas the categorical approach provides a higher-level language that facilitates reasoning about complex, compositional systems.

This line of inquiry is not merely a formal exercise; it is directly motivated by the challenges of understanding modern machine learning systems ([Yuan, 2023](#)). The compositional structure of deep neural networks, finds a natural description in the language of category theory ([Fong and Spivak, 2018](#); [Yuan, 2023](#); [Pardo-Guerra et al., 2025](#)).

**The Metric View: Probability as a Manifold.** In parallel to the algebraic developments in category theory, the field of *information geometry* (IG) has provided a powerful differential geometric lens for studying machine learning (Amari and Nagaoka, 2000). Foundational work by Amari and Nagaoka (2000) demonstrated that a parametric family of probability distributions can be viewed as a smooth manifold endowed with a canonical Riemannian metric—the Fisher-Rao metric—and a pair of dually-coupled affine connections. This geometric structure is not arbitrary; it can be intrinsically derived from statistical divergence functions, such as the Kullback-Leibler (KL) divergence, which serves as a measure of dissimilarity between distributions.

When applying these geometric tools to deep learning, it is crucial to draw a distinction between IG and the related field of *geometric deep learning* (GDL). GDL is primarily concerned with generalizing neural network architectures to operate on data that resides in non-Euclidean domains, such as graphs or manifolds; its focus is the geometry of the *input data space*. In contrast, IG has traditionally been used to analyze the geometry of the *parameter space* of a model. By viewing the set of all possible model parameters as a manifold, IG provides sophisticated tools for understanding the dynamics of training, offering a more nuanced perspective on optimization than that afforded by standard  $\ell_1$  or  $\ell_2$  regularization (Amari and Nagaoka, 2000).

**Towards Categorical Information Geometry.** While the algebraic and geometric approaches have largely evolved independently, a new frontier is emerging at their intersection. Recent work has begun to explicitly forge a synthesis, aiming to create a *categorical information geometry* (Perrone, 2023b). This research program, led by researchers such as Perrone (2023b), seeks to enrich the abstract, compositional structures of Markov categories with the metric and quantitative notions central to information theory and geometry, such as entropy and divergence (Perrone, 2023a,b).

This emerging synthesis recognizes that a complete theoretical picture requires both the compositional language of categories and the metric language of geometry. However, to date, the applications of this nascent field have focused primarily on reformulating abstract probability theory. The critical connection to the analysis of *learned representations* in practical, large-scale deep learning models remains largely unexplored. This work aims to bridge that gap, demonstrating that a categorical information geometry provides the ideal framework for analyzing the structure of the representation spaces sculpted by the learning process.

## 8.2 The Geometry of Learning and Representation

**Objectives of Learning Representations.** A guiding principle for understanding the purpose of representation learning is the *Information Bottleneck* (IB) theory, introduced by Tishby et al. (2000b). The IB principle posits that an optimal representation  $T$  of some input data  $X$  should be a bottleneck that is maximally informative about a relevant target variable  $Y$  while being maximally compressed with respect to the input  $X$  (Tishby et al., 2000b; Shwartz-Ziv and Tishby, 2017). This trade-off between predictive accuracy and compressional complexity provides a powerful, high-level objective for representation learning.

The IB framework has been particularly influential in the theoretical analysis of deep learning. It led to the hypothesis that the training of deep neural networks proceeds in two distinct phases: an initial fitting phase, where the mutual information  $I(T; Y)$  between the representation and the target increases, followed by a “compression” phase, where the mutual information  $I(X; T)$  between

the input and the representation decreases (Shwartz-Ziv and Tishby, 2017). While the universality of this two-phase dynamic is a subject of ongoing debate, the core intuition—that effective training involves not just memorization but also a form of structured compression—provides a compelling motivation for investigating the geometry of the learned representations.

**Measurements of Representation Geometry.** To move from high-level principles to concrete analysis, we require quantitative tools to probe the geometric properties of the high-dimensional activation spaces within a neural network. A particularly effective set of tools for this purpose comes from spectral graph theory. Given a graph with adjacency matrix  $A$  and degree matrix  $D$ , the *Graph Laplacian* is defined as  $L = D - A$  (Berahmand et al., 2025). For any function  $f$  defined on the nodes of the graph (e.g., a feature activation), the quadratic form  $f^\top L f$  defines the graph’s *Dirichlet energy*. This quantity measures the smoothness of the function with respect to the graph structure; a low Dirichlet energy indicates that connected nodes have similar function values (Park et al., 2024).

This concept has a rich history in machine learning, forming the basis of spectral clustering algorithms (Berahmand et al., 2025) and, more recently, being used to analyze and mitigate the over-smoothing problem in Graph Neural Networks (GNNs). Over-smoothing occurs when stacking many GNN layers causes the representations of all nodes to converge to an indistinguishable point, a phenomenon characterized by the Dirichlet energy of the representations collapsing to zero.

Most critically for the present work, this classical geometric tool has been shown to be deeply relevant to the internal dynamics of the most advanced models. A recent study demonstrates that during in-context learning, Large Language Models (LLMs) dynamically reorganize their internal concept representations in a manner that explicitly minimizes the Dirichlet energy with respect to an implicit graph structure defined by the context (Park et al., 2024). This groundbreaking result elevates Dirichlet energy from a tool for analyzing explicit graphs to a general principle governing the emergent geometry of learned representations. This trend towards spectral analysis is further evidenced by other recent work using methods like Centered Kernel Alignment (CKA) to track representation dynamics and spectral editing of activations (SEA) to control model behavior (Qiu et al., 2024).

**Mechanisms of Learning Representations.** The geometric structures observed in learned representations are not accidental; they are a direct consequence of the implicit biases of the training algorithm. For modern, highly overparameterized models, the optimization process itself, typically driven by variants of gradient descent, imparts an *implicit bias* or *implicit regularization* on the final solution (Vardi, 2023). Even when multiple parameter settings can achieve zero training error, the optimization algorithm preferentially converges to a “simple” solution that generalizes well. For linear models trained on classification tasks, this bias often corresponds to finding the maximum-margin separator, a classic geometric concept (Gunasekar et al., 2018; Soudry et al., 2018; Chizat and Bach, 2020).

A crucial modern insight is that training with the standard Negative Log-Likelihood (NLL) objective can be understood as a form of *implicit contrastive learning*. The NLL loss for a sample  $(x, y_{\text{true}})$ , given by  $-\log P(y_{\text{true}}|x) = -\log \frac{\exp(z_{\text{true}})}{\sum_j \exp(z_j)}$ , is minimized by simultaneously increasing the logit  $z_{\text{true}}$  for the correct class and decreasing the logits  $z_j$  for all incorrect classes. This dynamic is functionally

equivalent to a contrastive loss that pulls the representation of  $x$  towards a “positive” target (the true class) while pushing it away from all negative targets (the incorrect classes). This connection is profound, as it links the vast body of work on self-supervised contrastive learning—which explicitly aims to sculpt a geometrically structured representation space—to standard supervised training. It provides the causal mechanism for why NLL training produces well-structured, geometrically organized representations: the objective itself is implicitly performing a contrastive optimization that enforces this structure.

## 9 Conclusion

In this work, we introduced a mathematical framework for analyzing the auto-regressive generation step in language models, leveraging the expressive power of Markov Categories. By modeling the process as a composition of Markov kernels,  $k_{\text{gen},\theta} = k_{\text{head}} \circ k_{\text{bb}} \circ k_{\text{emb}}$ , we established a foundation for a compositional, information-theoretic analysis. This allowed us to formally quantify the information surplus in hidden states, providing a clear theoretical rationale for the success of modern speculative decoding techniques like EAGLE.

More importantly, our framework provides a unified, first-principles explanation for the remarkable effectiveness of the negative log-likelihood (NLL) objective. We showed that NLL training is not merely about predicting the next token; it is a powerful structure-learning algorithm in disguise. We proved that minimizing NLL forces the model to: (1) achieve optimal data compression by learning the intrinsic conditional stochasticity of the data, a process we measured with categorical entropy; and (2) implicitly perform spectral contrastive learning. By analyzing the information geometry of the prediction head via the pullback Fisher-Rao metric, we demonstrated that NLL sculpts the representation space, aligning its geometry with the eigenspectrum of a predictive similarity operator. This explains how NLL learns semantically organized representations without explicit contrastive pairs.

This compositional, probabilistic, and information-geometric perspective offers a principled, mathematically grounded alternative to purely empirical or heuristic analysis, unifying concepts from information theory, geometry, and spectral methods to reveal the deep principles driving the success of large language models.

## References

- Shun-ichi Amari. Information geometry in optimization, machine learning and statistical inference. *Frontiers of Electrical and Electronic Engineering in China*, 5(3):241–260, 2010.
- Shun-ichi Amari and Hiroshi Nagaoka. *Methods of information geometry*, volume 191. American Mathematical Soc., 2000.
- John C Baez, Brendan Fong, and Blake S Pollard. A compositional framework for markov processes. *Journal of Mathematical Physics*, 57(3):033301, 2016.
- Kamal Berahmand, Farid Saberi-Movahed, Razieh Sheikhpour, Yuefeng Li, and Mahdi Jalili. A comprehensive survey on spectral clustering with graph structure learning. *arXiv preprint arXiv:2501.13597*, 2025.

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Advances in neural information processing systems*, volume 33, pages 1877–1901, 2020.
- Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Conference on learning theory*, pages 1305–1338. PMLR, 2020.
- Kenta Cho and Bart Jacobs. Disintegration and bayesian inversion via string diagrams. *Mathematical Structures in Computer Science*, 29(7):938–971, 2019.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for transformer circuits, 2021.
- Brendan Fong and David I Spivak. Seven sketches in compositionality: An invitation to applied category theory. *arXiv preprint arXiv:1803.05316*, 2018.
- Tobias Fritz. A synthetic approach to markov kernels, conditional independence and theorems on sufficient statistics. *Advances in Mathematics*, 370:107239, 2020a.
- Tobias Fritz. A synthetic approach to markov kernels, conditional independence and theorems on sufficient statistics. *Advances in Mathematics*, 370:107239, 2020b.
- Tobias Fritz and Eigil Fjeldgren Rischel. Infinite products and zero-one laws in categorical probability. *Compositionality*, 2:3, 2020. doi: 10.32408/compositionality-2-3.
- Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. Better & faster large language models via multi-token prediction. *arXiv preprint arXiv:2404.19737*, 2024.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, pages 1832–1841. PMLR, 2018.
- Jeff Z. HaoChen, Hao Chen, Chen Wei, Adrien Gaidon, and Tengyu Ma. Provable Guarantees for Self-Supervised Deep Learning with Spectral Contrastive Loss. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pages 14239–14250, 2021.
- John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, 2019.
- Olav Kallenberg and Olav Kallenberg. *Foundations of modern probability*, volume 2. Springer, 1997.
- Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. Eagle: Speculative sampling requires rethinking feature uncertainty. *arXiv preprint arXiv:2401.15077*, 2024.
- Christopher D Manning, Kevin Clark, John Hewitt, Urvashi Khandelwal, and Omer Levy. Emergent

- linguistic structure in artificial neural networks trained by self-supervision. *Proceedings of the National Academy of Sciences*, 117(48):30046–30054, 2020.
- Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in neural information processing systems*, volume 29, 2016.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- Sebastian Pardo-Guerra, Johnny Jingze Li, Kalyan Basu, and Gabriel A Silva. Neural networks and markov categories. *AppliedMath*, 5(3):93, 2025.
- Core Francisco Park, Andrew Lee, Ekdeep Singh Lubana, Yongyi Yang, Maya Okawa, Kento Nishi, Martin Wattenberg, and Hidenori Tanaka. Iclr: In-context learning of representations. *arXiv preprint arXiv:2501.00070*, 2024.
- Paolo Perrone. Markov categories and entropy. *IEEE Transactions on Information Theory*, 70(3):1671–1692, 2023a.
- Paolo Perrone. Categorical information geometry. In *International Conference on Geometric Science of Information*, pages 268–277. Springer, 2023b.
- Yifu Qiu, Zheng Zhao, Yftah Ziser, Anna Korhonen, Edoardo Maria Ponti, and Shay Cohen. Spectral editing of activations for large language model alignment. *Advances in Neural Information Processing Systems*, 37:56958–56987, 2024.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*, 2017.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57, 2018.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- Zhiqian Tan, Yifan Zhang, Jingqin Yang, and Yang Yuan. Contrastive Learning Is Spectral Clustering On Similarity Graph. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024.
- Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000a.
- Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000b.
- Gal Vardi. On the implicit bias in deep-learning algorithms. *Communications of the ACM*, 66(6):86–93, 2023.



- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, volume 30, 2017.
- Yang Yuan. On the power of foundation models. In *International Conference on Machine Learning*, pages 40519–40530. PMLR, 2023.

# Appendix

<b>A</b>	<b>Full Proofs of Theorems</b>	<b>29</b>
A.1	Proof of Theorem 5.1 (NLL Minimization as Average KL Minimization) . . . . .	29
A.2	Proof of Theorem 5.2 (Convergence of Average Categorical Entropy) . . . . .	30
A.3	Proof of Theorem 7.1 (Output Distribution Approximation Constraint) . . . . .	31
A.4	Proof of Corollary 7.2 (Implicit Representation Separation) . . . . .	32
A.5	Proof of Proposition 7.4 (NLL Objective and Implicit Dirichlet Energy Minimization) . . . . .	33
A.6	Proof of Proposition 7.6 (NLL Objective as Implicit Spectral Alignment) . . . . .	33

## A Full Proofs of Theorems

### A.1 Proof of Theorem 5.1 (NLL Minimization as Average KL Minimization)

Let  $p_x(\cdot) := k_{\text{data}}(x, \cdot)$  denote the true conditional probability distribution  $P_{\text{data}}(\cdot|x)$  for context  $x = \mathbf{w}_{<t}$ . Let  $q_{x,\theta}(\cdot) = k_{\text{gen},\theta}(x, \cdot)$  denote the model's conditional probability distribution  $P_\theta(\cdot|x)$ . The context distribution is  $p_{W_{<t}}$ .

The cross-entropy loss is defined as:

$$\begin{aligned} L_{\text{CE}}(\theta) &= -\mathbb{E}_{(x,w) \sim P_{\text{data}}} [\log q_{x,\theta}(w)] \\ &= -\mathbb{E}_{x \sim p_{W_{<t}}} [\mathbb{E}_{W \sim p_x(\cdot)} [\log q_{x,\theta}(W)]] \\ &= -\mathbb{E}_{x \sim p_{W_{<t}}} \left[ \sum_{w \in \mathbb{V}} p_x(w) \log q_{x,\theta}(w) \right] \end{aligned}$$

The average KL divergence is defined as:

$$\begin{aligned} \mathcal{L}_{\text{KL}}(\theta) &= \mathbb{E}_{x \sim p_{W_{<t}}} [D_{\text{KL}}(p_x(\cdot) \parallel q_{x,\theta}(\cdot))] \\ &= \mathbb{E}_{x \sim p_{W_{<t}}} \left[ \sum_{w \in \mathbb{V}} p_x(w) \log \frac{p_x(w)}{q_{x,\theta}(w)} \right] \\ &= \mathbb{E}_{x \sim p_{W_{<t}}} \left[ \sum_{w \in \mathbb{V}} p_x(w) \log p_x(w) - \sum_{w \in \mathbb{V}} p_x(w) \log q_{x,\theta}(w) \right] \\ &= \mathbb{E}_{x \sim p_{W_{<t}}} [-H(p_x(\cdot))] - \mathbb{E}_{x \sim p_{W_{<t}}} \left[ \sum_{w \in \mathbb{V}} p_x(w) \log q_{x,\theta}(w) \right] \\ &= -H(W_t|W_{<t})_{\text{data}} + L_{\text{CE}}(\theta) \end{aligned}$$

where  $H(p_x(\cdot))$  is the Shannon entropy of the distribution  $p_x(\cdot)$ , and  $H(W_t|W_{<t})_{\text{data}} = \mathbb{E}_{x \sim p_{W_{<t}}} [H(p_x(\cdot))]$  is the average conditional Shannon entropy of the data generating process.

Rearranging gives:

$$L_{\text{CE}}(\theta) = \mathcal{L}_{\text{KL}}(\theta) + H(W_t|W_{<t})_{\text{data}}$$

Since  $H(W_t|W_{<t})_{\text{data}}$  is a property of the data distribution and does not depend on the model parameters  $\theta$ , minimizing  $L_{\text{CE}}(\theta)$  with respect to  $\theta$  is equivalent to minimizing  $\mathcal{L}_{\text{KL}}(\theta)$ .

The KL divergence  $D_{\text{KL}}(p \parallel q) \geq 0$  for any probability distributions  $p, q$ , with equality if and only if  $p = q$ . Therefore, the average KL divergence  $\mathcal{L}_{\text{KL}}(\theta) = \mathbb{E}_{x \sim p_{W_{<t}}} [D_{\text{KL}}(p_x(\cdot) \parallel q_{x,\theta}(\cdot))]$  is also non-negative, as it is an expectation of non-negative values.

The minimum value  $\mathcal{L}_{\text{KL}}(\theta) = 0$  is achieved if and only if the integrand is zero  $p_{W_{<t}}$ -almost everywhere. That is,  $D_{\text{KL}}(p_x(\cdot) \parallel q_{x,\theta^*}(\cdot)) = 0$  for  $p_{W_{<t}}$ -almost every  $x$ . This occurs if and only if  $p_x(\cdot) = q_{x,\theta^*}(\cdot)$  for  $p_{W_{<t}}$ -almost every  $x$ . In terms of kernels, this means  $k_{\text{data}}(x, \cdot) = k_{\text{gen},\theta^*}(x, \cdot)$  for  $p_{W_{<t}}$ -almost every  $x$ .

If the model class  $\{k_{\text{gen},\theta}\}$  contains  $k_{\text{data}}$ , say  $k_{\text{data}} = k_{\text{gen},\theta_{\text{true}}}$ , then choosing  $\theta^* = \theta_{\text{true}}$  achieves  $\mathcal{L}_{\text{KL}}(\theta^*) = 0$ , which is the minimum possible value.  $\square$

## A.2 Proof of Theorem 5.2 (Convergence of Average Categorical Entropy)

We want to show that  $\lim_{n \rightarrow \infty} \bar{\mathcal{H}}_D(k_{\text{head},n}; p_{H_t,\theta_n}) = \bar{\mathcal{H}}_D(k_{\text{head},\theta^*}; p_{H_t,\theta^*})$ .

Recall the definition:

$$\bar{\mathcal{H}}_D(k_{\text{head},\theta}; p_{H_t,\theta}) = \mathbb{E}_{h \sim p_{H_t,\theta}} [\Psi_D(h, k_{\text{head},\theta}(h, \cdot))],$$

where  $\Psi_D(h, p) := D_{\mathbb{V} \otimes \mathbb{V}}(\sum_{w \in \mathbb{V}} p(w) \delta_{(w,w)} \parallel p \otimes p)$ , and  $p = k_{\text{head},\theta}(h, \cdot)$ .

Let  $X_n$  be the random variable  $\Psi_D(H_n, k_{\text{head},n}(H_n, \cdot))$  where  $H_n \sim p_{H_t,\theta_n}$ . We want to show  $\lim_{n \rightarrow \infty} \mathbb{E}[X_n] = \mathbb{E}[X^*]$ , where  $X^* = \Psi_D(H^*, k_{\text{head},\theta^*}(H^*, \cdot))$  with  $H^* \sim p_{H_t,\theta^*}$ .

We are given: (i)  $k_{\text{head},n}(h, \cdot) \rightarrow k_{\text{head},\theta^*}(h, \cdot)$  in a suitable topology (e.g., total variation) for  $p_{H_t,\theta^*}$ -almost every  $h$ . Let's denote this  $p_n(h) \rightarrow p^*(h)$ .

(ii)  $p_{H_t,\theta_n} \Rightarrow p_{H_t,\theta^*}$  (weak convergence). This means  $\int g(h) p_{H_t,\theta_n}(dh) \rightarrow \int g(h) p_{H_t,\theta^*}(dh)$  for all bounded continuous functions  $g : \mathcal{H} \rightarrow \mathbb{R}$ .

(iii) The function  $\Psi_D(h, p)$  is continuous and bounded in  $p$  (with respect to the topology in (i)) for relevant  $h$ . Since  $\mathbb{V}$  is finite, standard divergences like KL and TV are continuous functions of the probability vectors  $p \in \Delta^{|\mathbb{V}|-1}$ . The map  $p \mapsto \sum p(w) \delta_{(w,w)}$  and  $p \mapsto p \otimes p$  are also continuous. Thus,  $p \mapsto \Psi_D(h, p)$  is continuous for fixed  $h$ . Boundedness also holds for typical divergences on finite spaces. Let  $M$  be an upper bound:  $|\Psi_D(h, p)| \leq M$ .

Let  $\Phi_n(h) = \Psi_D(h, k_{\text{head},n}(h, \cdot))$  and  $\Phi^*(h) = \Psi_D(h, k_{\text{head},\theta^*}(h, \cdot))$ . From (i) and the continuity part of (iii), we have  $\Phi_n(h) \rightarrow \Phi^*(h)$  for  $p_{H_t,\theta^*}$ -almost every  $h$ .

We want to show  $\lim_{n \rightarrow \infty} \int \Phi_n(h) p_{H_t,\theta_n}(dh) = \int \Phi^*(h) p_{H_t,\theta^*}(dh)$ .

Consider the difference:

$$\begin{aligned} |\mathbb{E}[X_n] - \mathbb{E}[X^*]| &= \left| \int \Phi_n(h) p_{H_t,\theta_n}(dh) - \int \Phi^*(h) p_{H_t,\theta^*}(dh) \right| \\ &\leq \left| \int \Phi_n(h) p_{H_t,\theta_n}(dh) - \int \Phi^*(h) p_{H_t,\theta_n}(dh) \right| \\ &\quad + \left| \int \Phi^*(h) p_{H_t,\theta_n}(dh) - \int \Phi^*(h) p_{H_t,\theta^*}(dh) \right| \\ &= \left| \int (\Phi_n(h) - \Phi^*(h)) p_{H_t,\theta_n}(dh) \right| + \left| \int \Phi^*(h) p_{H_t,\theta_n}(dh) - \int \Phi^*(h) p_{H_t,\theta^*}(dh) \right| \end{aligned}$$

The second term converges to 0 as  $n \rightarrow \infty$  due to the weak convergence (ii), provided  $\Phi^*(h)$  is bounded and continuous. While  $\Phi^*(h)$  might not be continuous in  $h$ , if it is bounded and continuous  $p_{H_t,\theta^*}$ -almost everywhere, weak convergence is often sufficient. Let's assume  $\Phi^*(h)$  behaves well enough (e.g., is bounded and continuous almost everywhere w.r.t. the limiting measure  $p_{H_t,\theta^*}$ ) for  $\int \Phi^*(h) p_{H_t,\theta_n}(dh) \rightarrow \int \Phi^*(h) p_{H_t,\theta^*}(dh)$ . (This is sometimes known as the Generalized Continuous Mapping Theorem or Portmanteau Theorem).

For the first term, we have  $\Phi_n(h) \rightarrow \Phi^*(h)$  for  $p_{H_t,\theta^*}$ -almost every  $h$ . We also have the bound  $|\Phi_n(h) - \Phi^*(h)| \leq |\Phi_n(h)| + |\Phi^*(h)| \leq 2M$  from the boundedness assumption (iii). We can use a variant of the Dominated Convergence Theorem adapted for converging measures. Since  $p_{H_t,\theta_n} \Rightarrow p_{H_t,\theta^*}$

$p_{H_t, \theta^*}$  and  $\Phi_n \rightarrow \Phi^*$  pointwise a.e. (w.r.t.  $p_{H_t, \theta^*}$ ), and the sequence  $\Phi_n$  is uniformly bounded, we can conclude that  $\int (\Phi_n(h) - \Phi^*(h)) p_{H_t, \theta_n}(dh) \rightarrow 0$ . Combining these, we get  $\lim_{n \rightarrow \infty} \mathbb{E}[X_n] = \mathbb{E}[X^*]$ .

For the final part: if the model class is expressive such that  $k_{\text{gen}, \theta^*} = k_{\text{data}}$  (meaning  $\mathcal{L}_{\text{KL}}(\theta^*) = 0$ ), then the model perfectly matches the data generating process almost everywhere. If we assume the data process can be similarly factorized  $k_{\text{data}} = k_{\text{head, data}} \circ k_{\text{enc, data}}$ , then matching  $k_{\text{gen}, \theta^*} = k_{\text{data}}$  implies that the components must match (up to potential identifiability issues, e.g., transformations between the encoder output and head input that cancel out). Under reasonable assumptions (e.g., the factorization is unique in the relevant sense), we would have  $k_{\text{head}, \theta^*} \approx k_{\text{head, data}}$  and the distribution induced by the encoder  $k_{\text{bb}} \circ k_{\text{emb}}$  would approximate the distribution of the “true” internal state feeding into  $k_{\text{head, data}}$ , i.e.,  $p_{H_t, \theta^*} \approx p_{H_t, \text{data}}$ . Therefore,  $\bar{\mathcal{H}}_D(k_{\text{head}, \theta^*}; p_{H_t, \theta^*}) \approx \bar{\mathcal{H}}_D(k_{\text{head, data}}; p_{H_t, \text{data}})$ .  $\square$

### A.3 Proof of Theorem 7.1 (Output Distribution Approximation Constraint)

We are given the average KL divergence loss:

$$\mathcal{L}(\theta) = \mathbb{E}_{x \sim \mu_{\text{ctx}}} [D_{\text{KL}}(P_{\text{data}}(\cdot|x) \| p_{\theta}(\cdot|x))]$$

where  $p_{\theta}(\cdot|x) = g_{\text{head}}(f_{\text{enc}}(x))$  and  $\mu_{\text{ctx}}$  is the distribution over contexts. We are also given a metric  $d_{\text{out}}$  on  $\mathcal{P}(\mathbb{V})$  satisfying a Pinsker-type inequality:

$$d_{\text{out}}(p, q)^k \leq C \cdot D_{\text{KL}}(p \| q)$$

for some constants  $k, C > 0$ . Examples include Hellinger distance  $d_H$  ( $k = 2, C = 1/2$ ) and Total Variation distance  $d_{\text{TV}}$  ( $k = 2, C = 1$ ).

Let  $p = P_{\text{data}}(\cdot|x)$  and  $q = p_{\theta}(\cdot|x)$  for a specific context  $x$ . Applying the inequality yields:

$$d_{\text{out}}(P_{\text{data}}(\cdot|x), p_{\theta}(\cdot|x))^k \leq C \cdot D_{\text{KL}}(P_{\text{data}}(\cdot|x) \| p_{\theta}(\cdot|x)).$$

Now, we take the expectation of both sides with respect to the context distribution  $x \sim \mu_{\text{ctx}}$ . Since expectation is linear and the inequality holds pointwise for each  $x$ , we get:

$$\begin{aligned} \mathbb{E}_{x \sim \mu_{\text{ctx}}} [d_{\text{out}}(P_{\text{data}}(\cdot|x), p_{\theta}(\cdot|x))^k] &\leq \mathbb{E}_{x \sim \mu_{\text{ctx}}} [C \cdot D_{\text{KL}}(P_{\text{data}}(\cdot|x) \| p_{\theta}(\cdot|x))] \\ \mathbb{E}_{x \sim \mu_{\text{ctx}}} [d_{\text{out}}(P_{\text{data}}(\cdot|x), p_{\theta}(\cdot|x))^k] &\leq C \cdot \mathbb{E}_{x \sim \mu_{\text{ctx}}} [D_{\text{KL}}(P_{\text{data}}(\cdot|x) \| p_{\theta}(\cdot|x))] \\ \mathbb{E}_{x \sim \mu_{\text{ctx}}} [d_{\text{out}}(P_{\text{data}}(\cdot|x), p_{\theta}(\cdot|x))^k] &\leq C \cdot \mathcal{L}(\theta). \end{aligned}$$

This establishes the first part of the theorem, [equation \(7.2\)](#).

For the second part, consider any two contexts  $x, x'$ . Let  $p_x^{\text{data}} = P_{\text{data}}(\cdot|x)$ ,  $p_{x'}^{\text{data}} = P_{\text{data}}(\cdot|x')$ ,  $p_x^{\theta} = p_{\theta}(\cdot|x)$ , and  $p_{x'}^{\theta} = p_{\theta}(\cdot|x')$ . The triangle inequality for the metric  $d_{\text{out}}$  states:

$$d_{\text{out}}(A, C) \leq d_{\text{out}}(A, B) + d_{\text{out}}(B, C)$$

Applying this twice:

$$\begin{aligned} d_{\text{out}}(p_x^{\text{data}}, p_{x'}^{\text{data}}) &\leq d_{\text{out}}(p_x^{\text{data}}, p_x^{\theta}) + d_{\text{out}}(p_x^{\theta}, p_{x'}^{\text{data}}) \\ &\leq d_{\text{out}}(p_x^{\text{data}}, p_x^{\theta}) + (d_{\text{out}}(p_x^{\theta}, p_{x'}^{\theta}) + d_{\text{out}}(p_{x'}^{\theta}, p_{x'}^{\text{data}})) \\ &= d_{\text{out}}(p_x^{\text{data}}, p_x^{\theta}) + d_{\text{out}}(p_x^{\theta}, p_{x'}^{\theta}) + d_{\text{out}}(p_{x'}^{\theta}, p_{x'}^{\text{data}}). \end{aligned}$$

Let  $\epsilon_x = d_{\text{out}}(p_x^{\text{data}}, p_x^\theta)$  and  $\epsilon_{x'} = d_{\text{out}}(p_{x'}^\theta, p_{x'}^{\text{data}})$ . If the model fits the data well,  $\mathcal{L}(\theta)$  is small. From [equation \(7.2\)](#),  $\mathbb{E}_{x \sim \mu_{\text{ctx}}}[\epsilon_x^k] \leq C\mathcal{L}(\theta)$ , meaning the expected error (to the power  $k$ ) is small. By Markov's inequality, for any  $\delta > 0$ ,

$$\mathbb{P}(\epsilon_x^k \geq \delta^k) \leq \frac{\mathbb{E}[\epsilon_x^k]}{\delta^k} \leq \frac{C\mathcal{L}(\theta)}{\delta^k}.$$

Thus,  $\mathbb{P}(\epsilon_x \geq \delta)$  is small if  $\mathcal{L}(\theta)$  is small, implying that for a vast majority of contexts  $x$  drawn from  $\mu_{\text{ctx}}$ , the individual error  $\epsilon_x$  is small. Therefore, for typical pairs  $(x, x')$ , both  $\epsilon_x$  and  $\epsilon_{x'}$  are small.

Rearranging the triangle inequality gives:

$$\begin{aligned} d_{\text{out}}(p_x^\theta, p_{x'}^\theta) &\geq d_{\text{out}}(p_x^{\text{data}}, p_{x'}^{\text{data}}) - (\epsilon_x + \epsilon_{x'}) \\ d_{\text{out}}(p_x^\theta, p_{x'}^\theta) &\leq d_{\text{out}}(p_x^{\text{data}}, p_{x'}^{\text{data}}) + (\epsilon_x + \epsilon_{x'}) \end{aligned}$$

When  $\epsilon_x$  and  $\epsilon_{x'}$  are small, these inequalities show that the distance between the model's output distributions,  $d_{\text{out}}(p_x^\theta, p_{x'}^\theta)$ , must be close to the distance between the true data distributions,  $d_{\text{out}}(p_x^{\text{data}}, p_{x'}^{\text{data}})$ . In particular, if  $d_{\text{out}}(p_x^{\text{data}}, p_{x'}^{\text{data}})$  is large (predictively dissimilar contexts), then  $d_{\text{out}}(p_x^\theta, p_{x'}^\theta)$  must also be large, as the difference is bounded by small error terms.  $\square$

#### A.4 Proof of Corollary 7.2 (Implicit Representation Separation)

From Theorem 7.1, if  $\mathcal{L}(\theta)$  is small, then for typical contexts  $x, x'$ :

$$d_{\text{out}}(g_{\text{head}}(h_x), g_{\text{head}}(h_{x'})) \approx d_{\text{out}}(P_{\text{data}}(\cdot|x), P_{\text{data}}(\cdot|x')).$$

If contexts  $x, x'$  are predictively dissimilar, the RHS is large, implying the LHS must also be large. Let  $p_x = g_{\text{head}}(h_x)$  and  $p_{x'} = g_{\text{head}}(h_{x'})$  be two points in the manifold  $(\mathcal{P}(\mathbb{V}), g^{\text{FR}})$ . The distance between them, under the geometry induced by the Fisher-Rao metric, is the infimum of the lengths of all paths connecting them. The length of any such path  $\gamma(t)$  from  $p_{x'}$  to  $p_x$  is given by  $L(\gamma) = \int_0^1 \sqrt{g_{\gamma(t)}^{\text{FR}}(\gamma'(t), \gamma'(t))} dt$ .

Consider the specific path induced by the straight line in representation space,  $\gamma(t) = g_{\text{head}}((1-t)h_{x'} + th_x)$ . Let  $h(t) = (1-t)h_{x'} + th_x$ . The tangent vector to this path in  $\mathcal{P}(\mathbb{V})$  is  $\gamma'(t) = J_{g_{\text{head}}}(h(t)) \cdot (h_x - h_{x'})$ , where  $J_{g_{\text{head}}}$  is the Jacobian of  $g_{\text{head}}$ . The squared length of this tangent vector is:

$$g_{\gamma(t)}^{\text{FR}}(\gamma'(t), \gamma'(t)) = (h_x - h_{x'})^\top J_{g_{\text{head}}}(h(t))^\top g_{\gamma(t)}^{\text{FR}} J_{g_{\text{head}}}(h(t)) (h_x - h_{x'}) = g_{h(t)}^*(h_x - h_{x'}, h_x - h_{x'}).$$

The total squared length of this particular path is  $(\int_0^1 \sqrt{g_{h(t)}^*(h_x - h_{x'}, h_x - h_{x'})} dt)^2$ . Since the geodesic distance  $d_{\text{FR}}(p_x, p_{x'})$  is the shortest path length, it is bounded above by this length.

A large distance  $d_{\text{out}}(p_x, p_{x'})$  (which is closely related to  $d_{\text{FR}}(p_x, p_{x'})$ ) implies that the length of any path connecting them must be large. This in turn requires that the integrand  $\sqrt{g_{h(t)}^*(h_x - h_{x'}, h_x - h_{x'})}$  must be large over a substantial portion of the path. This can only happen if the difference vector  $h_x - h_{x'}$  has significant components along directions  $v$  where the pullback metric  $g^*$  is large (i.e., directions of high predictive sensitivity).



Conversely, if contexts are predictively similar, the RHS is small, implying the LHS must be small. This means the representations  $h_x$  and  $h_{x'}$  must be mapped to nearby points in  $\mathcal{P}(\mathbb{V})$ . This is achieved if  $h_x, h_{x'}$  are close along the directions of high predictive sensitivity, as these are the directions that would otherwise cause a large separation in the output space. Hence, NLL encourages representations of predictively similar contexts to be close along sensitive directions.  $\square$

## A.5 Proof of Proposition 7.4 (NLL Objective and Implicit Dirichlet Energy Minimization)

We analyze the Dirichlet energy:

$$\mathcal{E}_K(\phi_v) = \frac{1}{2} \iint K(x, x') (\langle h_x - h_{x'}, v \rangle)^2 \mu_{ctx}(dx) \mu_{ctx}(dx')$$

where  $v$  is a direction of high predictive sensitivity. The integrand,  $I(x, x') = K(x, x') (\langle h_x - h_{x'}, v \rangle)^2$ , contributes significantly to the integral only when both of its terms are non-negligible.

Consider the two cases based on the predictive similarity kernel  $K(x, x')$ :

1. **High Predictive Similarity ( $K(x, x')$  is large):** In this case, by definition of the kernel,  $d_{\text{out}}(P_{\text{data}}(\cdot|x), P_{\text{data}}(\cdot|x'))$  is small. By the converse part of Corollary 7.2, minimizing NLL encourages the representations  $h_x$  and  $h_{x'}$  to be close along sensitive directions  $v$ . This means the term  $(\langle h_x - h_{x'}, v \rangle)^2$  is encouraged to be small. Therefore, the integrand  $I(x, x')$  is a product of a large term ( $K(x, x')$ ) and a small term, which is small.
2. **Low Predictive Similarity ( $K(x, x')$  is small):** In this case,  $d_{\text{out}}(P_{\text{data}}(\cdot|x), P_{\text{data}}(\cdot|x'))$  is large. By Corollary 7.2, NLL forces the representations  $h_x$  and  $h_{x'}$  to be far apart along sensitive directions  $v$ , making  $(\langle h_x - h_{x'}, v \rangle)^2$  large. However, the integrand  $I(x, x')$  is a product of a small term ( $K(x, x')$ ) and a large term. The small value of  $K(x, x')$  suppresses the contribution of this pair to the total energy.

In both cases, the NLL objective pushes the model towards a configuration where the integrand of the Dirichlet energy is small. By discouraging configurations where both  $K(x, x')$  and the squared difference of the projected representations are simultaneously large, NLL minimization implicitly favors representations with low Dirichlet energy when projected onto predictively sensitive directions.  $\square$

## A.6 Proof of Proposition 7.6 (NLL Objective as Implicit Spectral Alignment)

The argument that Negative Log-Likelihood (NLL) optimization implicitly solves a spectral objective follows the principles established in the literature on provable spectral contrastive learning (HaoChen et al., 2021; Tan et al., 2024). We provide a detailed derivation by decomposing the NLL loss under a log-linear model assumption, revealing its connection to a spectral alignment objective.

**Setup.** To make the analysis transparent, we make a few simplifying assumptions that are common in theoretical analyses of representation learning:

- **Log-Linear Head:** We model the LM head kernel  $k_{\text{head}}$  as a standard log-linear or softmax model. The probability of the next token  $w \in \mathbb{V}$  given a representation  $h_x \in \mathcal{H}$  is:

$$p_\theta(w|x) = g_{\text{head}}(h_x)(w) = \frac{\exp(g(w)^\top h_x)}{Z(h_x)}, \quad \text{where } Z(h_x) = \sum_{w' \in \mathbb{V}} \exp(g(w')^\top h_x).$$

Here,  $g(w) \in \mathcal{H}$  is a learnable “prototype” or “output” vector for each token  $w$  in the vocabulary.

- **Fixed Norm Representations:** We assume the representations  $h_x$  and prototypes  $g(w)$  are constrained to have a fixed norm, e.g.,  $\|h_x\| = \|g(w)\| = 1$ . This is a common simplification that removes scaling effects and allows the analysis to focus on the geometric alignment (i.e., dot products and angles) between vectors.

Under these assumptions, the NLL loss for a single data point  $(x, w)$  where  $w \sim P_{\text{data}}(\cdot|x) = p_x$  is:

$$\mathcal{L}_x(w) = -\log p_\theta(w|x) = -g(w)^\top h_x + \log Z(h_x). \quad (\text{A.1})$$

The total expected loss over the data distribution is:

$$\mathcal{L}(\theta) = \mathbb{E}_{x \sim \mu_{\text{ctx}}} \left[ \mathbb{E}_{w \sim p_x} [-g(w)^\top h_x + \log Z(h_x)] \right]. \quad (\text{A.2})$$

**Decomposing the NLL Objective.** We can decompose the total loss into two meaningful components:

$$\mathcal{L}(\theta) = \underbrace{\mathbb{E}_{x \sim \mu_{\text{ctx}}} \left[ -(\mathbb{E}_{w \sim p_x} [g(w)])^\top h_x \right]}_{\mathcal{L}_{\text{align}}} + \underbrace{\mathbb{E}_{x \sim \mu_{\text{ctx}}} \left[ \log \sum_{w' \in \mathbb{V}} \exp(g(w')^\top h_x) \right]}_{\mathcal{L}_{\text{unif}}}. \quad (\text{A.3})$$

- **Alignment Loss ( $\mathcal{L}_{\text{align}}$ ):** This term measures the alignment between a representation  $h_x$  and the expected prototype vector for its corresponding true next-token distribution. Let’s define the **context-specific expected prototype** as  $\bar{g}_x \triangleq \mathbb{E}_{w \sim p_x} [g(w)]$ . Then the alignment loss is simply  $\mathcal{L}_{\text{align}} = -\mathbb{E}_x [\bar{g}_x^\top h_x]$ . Minimizing this term, by the Cauchy-Schwarz inequality, forces the representation  $h_x$  to align with its corresponding expected prototype  $\bar{g}_x$ . At the optimum for this term, we would have  $h_x \propto \bar{g}_x$ .
- **Uniformity Loss ( $\mathcal{L}_{\text{unif}}$ ):** This term, involving the log-partition function, acts as a regularizer. It prevents the trivial solution where all representations  $h_x$  and prototypes  $g(w)$  collapse to a single point. Minimizing the overall loss requires keeping the arguments of the exponentials,  $g(w')^\top h_x$ , from being universally large. This pushes representations of different contexts apart, encouraging them to spread out over the manifold (e.g., the unit sphere). This term plays a role analogous to the negative samples in contrastive learning objectives like InfoNCE.

**Connection to Predictive Similarity and Spectral Alignment.** The structure revealed by the loss decomposition allows us to connect NLL to a spectral objective.

- **Inducing a Similarity Metric:** The alignment objective  $h_x \propto \bar{g}_x$  is key. If the model successfully optimizes this part of the loss, the geometry of the learned representations  $\{h_x\}$  will mirror the geometry of the expected prototypes  $\{\bar{g}_x\}$ . This provides a natural way to define our predictive similarity kernel from Definition 7.3. Let's define it as the inner product of the expected prototypes:

$$K(x, x') \triangleq \langle \bar{g}_x, \bar{g}_{x'} \rangle.$$

This is a valid positive semidefinite kernel (it is the Gram matrix of the vectors  $\{\bar{g}_x\}$ ). If the model learns effectively, we have the crucial correspondence:

$$\langle h_x, h_{x'} \rangle \propto \langle \bar{g}_x, \bar{g}_{x'} \rangle = K(x, x').$$

This shows that the NLL objective implicitly forces the geometric structure (inner products) of the learned representations to match the predictive similarity kernel. This is precisely the goal of many contrastive and metric learning methods.

- **The Spectral Objective:** The goal of making the representation Gram matrix  $HH^\top$  (where  $H$  is the matrix of representations  $[h_{x_1}, h_{x_2}, \dots]$ ) proportional to the predictive similarity kernel matrix  $K$  is a spectral objective. Consider the problem of finding representations  $h$  that maximize alignment with the similarity structure. This can be formulated as maximizing:

$$\max_{f_\theta} \mathbb{E}_{x, x' \sim \mu_{ctx}} [\langle h_x, h_{x'} \rangle K(x, x')].$$

This objective is maximized when the principal components of the representation distribution  $\{h_x\}$  align with the leading eigenfunctions of the predictive similarity operator  $M_K$  from Definition 7.5. The optimization dynamics of stochastic gradient descent on the NLL loss can be shown to approximate a power iteration method for finding the leading eigenvectors of the operator  $M_K$ , especially when the uniformity term ensures that the representation covariance does not collapse.

**Conclusion.** In summary, the NLL objective, when viewed through the lens of a log-linear model, decomposes into two forces:

1. An **alignment force** ( $\mathcal{L}_{\text{align}}$ ) that pulls a representation  $h_x$  towards the average prototype  $\bar{g}_x$  of its true subsequent tokens.
2. A **uniformity force** ( $\mathcal{L}_{\text{unif}}$ ) that pushes representations apart to prevent collapse.

The interplay of these two forces encourages a solution where the pairwise similarity of representations  $\langle h_x, h_{x'} \rangle$  becomes proportional to the predictive similarity  $K(x, x') = \langle \bar{g}_x, \bar{g}_{x'} \rangle$ . This is equivalent to structuring the representation space such that its covariance is aligned with the eigenspectrum of the predictive similarity operator  $M_K$ . Therefore, NLL optimization does not just incidentally create structured representations; it actively performs a spectral decomposition of the predictive similarity graph inherent in the training data. The pullback metric  $g^*$  from the main text provides a more general, coordinate-free description of the local predictive sensitivity that is here captured by the linear prototypes  $g(w)$ .  $\square$