

---

# DECISION POTENTIAL SURFACE: A THEORETICAL AND PRACTICAL APPROXIMATION OF LLM’S DECISION BOUNDARY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Decision boundary, the subspace of inputs where a machine learning model assigns equal classification probabilities to two classes, is pivotal in revealing core model properties and interpreting behaviors. While analyzing the decision boundary of large language models (LLMs) has raised increasing attention recently, constructing it for mainstream LLMs remains computationally infeasible due to the enormous vocabulary-sequence sizes and the auto-regressive nature of LLMs. To address this issue, in this paper we propose *Decision Potential Surface (DPS)*, a new notion for analyzing LLM decision boundary. DPS is defined on the confidences in distinguishing different sampling sequences for each input, which naturally captures the *potential* of decision boundary. We prove that the zero-height isohypse in DPS is equivalent to the decision boundary of an LLM, with enclosed regions representing decision regions. By leveraging DPS, for the first time in the literature, we propose an approximate decision boundary construction algorithm, namely  $K$ -DPS, which only requires  $K$ -finite times of sequence sampling to approximate an LLM’s decision boundary with negligible error. We theoretically derive the upper bounds for the absolute error, expected error, and the error concentration between  $K$ -DPS and the ideal DPS, demonstrating that such errors can be trade-off with sampling times. Our results are empirically validated by extensive experiments across various LLMs and corpora.

## 1 INTRODUCTION

With the rapid advancement and remarkable success of large language models (LLMs), understanding their underlying mechanisms and behaviors has become increasingly critical (Wang et al., 2023; Conmy et al., 2023; Elhage et al., 2021; Ameisen et al., 2025; Sharkey et al., 2025; Allen-Zhu & Li, 2023; Liang et al., 2025; 2024). A key approach to demystifying the “black box” of state-of-the-art AI models involves analyzing the *decision boundary* (Rosenblatt, 1958), a fundamental concept for elucidating the characteristics of machine learning (ML) models. For LLMs, decision boundaries provide valuable insights into critical phenomena, including reasoning (Yang et al., 2025b), in-context learning (Zhao et al., 2024), hallucination (Mayne et al., 2025), memorization (Li et al., 2025), and so on.

As a foundational concept in machine learning, the decision boundary represents a subspace of inputs where a model assigns equal probability to two distinct classification outcomes (Rosenblatt, 1958). Extensive theoretical and empirical studies (Lee & Oommen, 1997; Turner & Ghosh, 1996; Goodfellow et al., 2015; Madry et al., 2018; Gu et al., 2017) have demonstrated that the properties of decision boundaries reveal critical attributes of machine learning models, including performance, robustness, and generalization. Consequently, constructing and leveraging decision boundaries for LLMs become a powerful and promising approach to enhance *almost all downstream analyses* of their behavior and capabilities.

Unfortunately, analyzing decision boundaries of LLMs incurs significantly greater complexity than on deep neural networks (DNNs) (Karimi et al., 2020; Karimi & Tang, 2020; Li et al., 2019; Lee & Landgrebe, 1997; Mickisch et al., 2020; Yousefzadeh & O’Leary, 2019). Unlike classification tasks with a limited number of classes (Lee & Oommen, 1997; Turner & Ghosh, 1996; Goodfellow

---

054 et al., 2015; Madry et al., 2018; Gu et al., 2017), LLMs predict a single token from an expansive  
055 vocabulary, often exceeding 100,000 tokens. Moreover, their autoregressive nature (Bengio et al.,  
056 2003; Radford et al., 2018) requires iterative token predictions to generate complete sequences,  
057 which further compounds the complexity of modeling decision boundaries. For instance, a Qwen-3  
058 model (with 8 billion parameters) (Yang et al., 2025a) supports sequences up to 32,768 tokens with  
059 a vocabulary of 151,936, resulting in approximately  $10^{169,790}$  decision regions! Such an enormous  
060 scale renders trivial attempts on decision-boundary-based analysis or visualization computation-  
061 ally infeasible. Prior studies (Zhao et al., 2024; Yang et al., 2025b; Mayne et al., 2025; Li et al.,  
062 2025), despite their valuable contributions to their specific motivating tasks, unfortunately sidestep  
063 this critical challenge. They either simplify the problem to toy scenarios, such as binary classifica-  
064 tion (Zhao et al., 2024; Mayne et al., 2025), or use the decision boundary concept metaphorically  
065 without constructing it (Yang et al., 2025b; Li et al., 2025). Consequently, the haunting questions  
066 remain unanswered — **What constitutes an LLM’s decision boundary, and is there a universal  
067 and yet efficient algorithm to construct it?**

068 To address these questions, we propose a principled strategy for modeling the decision boundaries  
069 of LLMs, which yields theoretical guarantees, computational tractability, and interpretability simul-  
070 taneously. Inspired by the existing decision boundaries for multi-class classification, we treat gen-  
071 erative language models as a composite multi-class classification task. As trivial solutions cannot  
072 model the complex decision boundaries for such tasks, we introduce a novel concept, namely *De-*  
073 *cision Potential Surface (DPS)*, to facilitate decision boundary analysis. It is a landscape in which  
074 every point encodes the *competition potential* among candidate outputs, quantified by a *decision*  
075 *potential function (DPF)*. We theoretically demonstrate that the zero-height *isohypse* of the DPS  
076 corresponds to the decision boundary, with the enclosed regions representing decision regions.

077 By examining the definition of DPS, we surprisingly discover that enumerating the entire output  
078 space is unnecessary for computing the DPF. Instead, a sufficient sampling already captures the  
079 “competition potential”. We therefore approximate the LLM’s decision boundary with only  $K$ -finite  
080 ( $K \ll$  realistic classification count) sequence sampling, yielding  $K$ -DPS and keeping the theoretical  
081 error within a provably small bound. We establish the error bound, expected error bound, and error  
082 concentration between the ideal DPS and  $K$ -DPS, demonstrating that  $K$ -DPS offers a favorable  
083 trade-off between approximation accuracy and computational cost. Finally, we conduct extensive  
084 experiments on open-source LLMs to evaluate the empirical performance of our method.

085 To the best of our knowledge, this is the first study on constructing decision boundaries for LLMs.  
086 Moreover, our proposed *decision potential surface (DPS)* framework is the first to provide a practical  
087 approximation of decision boundaries with theoretical guarantees. Our contributions are as follows:

- 088 • We formalize the definition of an LLM’s decision boundary as a composite multi-class  
089 classification task and analyze its fundamental properties.
- 090 • We introduce the concepts of the *Decision Potential Function (DPF)* and *Decision Potential*  
091 *Surface (DPS)*. We prove that the *isohypses* of the decision potential surface represent the  
092 marginal decision boundaries of LLMs, with the zero-height isohypse equivalent to the  
093 decision boundary.
- 094 • We propose  $K$ -DPS, an efficient and bounded approximation of the ideal DPS that requires  
095 only finite sampling for each input. We theoretically and empirically establish the error  
096 bounds of this approximation relative to the ideal DPS and quantify the trade-off between  
097 approximation error and sampling size.

## 099 2 RELATED WORKS

101 **Decision Boundary Analysis on Machine Learning Models.** The earliest exploration of decision  
102 boundaries in neural networks dates back to the era of linear classifiers and shallow architectures.  
103 Rosenblatt (1958) introduced the first linear decision boundary for binary classification, where a  
104 hyperplane separates input samples into two classes. For shallow feedforward neural networks  
105 (FFNNs) with non-linear activations (e.g., sigmoid, ReLU), subsequent work quantified how hidden  
106 layers enable non-linear decision boundaries: Lee & Oommen (1997) proposed a feature extrac-  
107 tion method that maps input data to a space aligned with FFNN decision boundaries, showing that  
boundary curvature correlates with model capacity and classification accuracy. For ensemble neural

---

networks, Turner & Ghosh (1996) made a pivotal contribution: they developed a theoretical framework that connects the stability of decision boundaries (relative to the optimal boundary defined by Bayes theory) to the overall error performance of the ensemble. Their work showed that linearly combining multiple well-designed, unbiased neural classifiers can reduce fluctuations in decision boundaries, which in turn lowers the extra error that goes beyond the theoretical minimum (i.e., Bayes error). A critical insight established by this finding is that the geometric characteristics of decision boundaries are closely tied to how well a model performs, which sheds light on employing decision boundary to explain the properties of neural networks.

**Decision Boundary Analysis on Neural Networks.** In recent years, researchers extended boundary analysis to convolutional neural networks (CNNs) and transformers for computer vision (CV) tasks. Goodfellow et al. (2015) revealed a key vulnerability of deep CNNs: their decision boundaries are locally linear in high-dimensional input spaces, making them susceptible to adversarial examples. Madry et al. (2018) further formalized this by proving that robust training (e.g., adversarial training) “smooths” decision boundaries, reducing local linearity and adversarial susceptibility. Similarly, Gu et al. (2017) focused on backdoor attacks in CNNs, linking them to hidden “trapdoors” in decision boundaries. Such attacks involve planting a small, specific pattern that shifts the boundary and forces misclassification for triggered inputs. Lee & Landgrebe (1997) laid the groundwork by introducing decision boundary feature extraction, highlighting the boundary’s role in characterizing network behavior before deep learning. Later, Yousefzadeh & O’Leary (2019) examined trained networks’ decision boundaries, analyzing how architectural elements (e.g., depth and activation functions) and training data influence boundary shape, complexity, and stability, providing insights into network task performance. Mickisch et al. (2020) conducted an empirical study on deep network boundaries across CV tasks, including image classification and object detection. Via quantitative and qualitative analysis, they explored boundary behavior near correct and misclassified samples and adversarial examples, bridging theory-practice gaps. In the same year, Karimi & Tang (2020) reviewed boundary research challenges such as high input dimensionality, complex architectures, and limited visualization tools, and opportunities, including advanced math, innovative visualization, and robustness enhancements. Karimi et al. (2020) complementarily proposed metrics like smoothness, curvature, and class separation to quantify boundaries, enabling cross-model comparisons and standardized analysis for deep learning interpretability.

**Decision Boundary Analysis on Large Language Models.** Extending decision boundary analysis from traditional neural networks (e.g., CNNs) to LLMs, researchers have begun exploring how this concept illuminates LLMs’ decision-making mechanisms and limitations. Zhao et al. (2024) probed the decision boundaries of in-context learning in LLMs, shedding light on how contextual information shapes boundary formation and decision outputs. Li et al. (2025) surveyed LLMs’ knowledge boundaries, laying a foundation for linking knowledge scope to decision boundaries. Mayne et al. (2025) revealed LLMs’ ignorance of their own decision boundaries and the unreliability of self-generated counterfactual explanations. Yang et al. (2025b) proposed BARREL, a boundary-aware reasoning framework to enhance LLM factuality via boundary awareness.

However, existing work fails to address LLM decision boundary analysis core challenges: Zhao et al. (2024) and Mayne et al. (2025) simplify to binary classification toy scenarios, deviating from LLMs’ real multi-token autoregressive prediction; Yang et al. (2025b) and Li et al. (2025) use the decision boundary concept metaphorically without concrete construction methods, unable to tackle exponential decision region complexity. This leaves a lack of universal, computationally feasible, and theoretically grounded boundary-capturing approaches, hindering LLM interpretation and optimization. Therefore, this paper aims to propose new decision boundary theory for addressing high-dimensional complexity and construction barriers, enabling accurate, efficient, and interpretable boundary modeling aligned with LLMs’ characteristics.

### 3 PRELIMINARIES: DECISION BOUNDARY OF LANGUAGE MODELS

#### 3.1 DECISION BOUNDARY AND ITS GEOMETRY ON CLASSIFICATION MODELS

We begin our theoretical analysis with traditional classification models and aim to extend the insights to generative language models.

Consider a neural network  $f : \mathbb{R}^d \rightarrow \mathbb{R}^M$  that maps an input sample  $\mathbf{x} \in \mathbb{R}^d$  to a predicted probability distribution over  $M$  classes, where the set of classes can be denoted as  $\mathcal{M} = \{1, 2, \dots, M\}$ .  $M > 2$ . Our goal is to characterize the properties of  $f$  under a specific input data distribution  $\mathcal{D} \subseteq \mathbb{R}^d$ . Without loss of generality, we decompose  $f$  into three components: (i) A representation module  $h = f_r(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  that maps the input  $\mathbf{x}$  to a latent representation  $h$ . (ii) A linear classification head  $z = f_{\text{cls}}(h) = W_{\text{cls}}h + b_{\text{cls}} : \mathbb{R}^{d'} \rightarrow \mathbb{R}^M$ , where  $W_{\text{cls}} \in \mathbb{R}^{M \times d'}$  and  $b_{\text{cls}} \in \mathbb{R}^M$  are learnable parameters, projecting the representation  $h$  into classification logits  $z$ . (iii) A nonlinear normalization function  $P = \sigma(z) : \mathbb{R}^M \rightarrow \mathbb{R}^M$ , which transforms the logits into a probability distribution  $P = [p_1, p_2, \dots, p_M]$ , where  $0 \leq p_i \leq 1$  for  $i = 1, \dots, M$  and  $\sum_{i=1}^M p_i = 1$ . The final predicted class for  $\mathbf{x}$  is determined by  $\arg \max_i p_i$ . Then, the decision boundary of the neural network is defined as follows.

**Definition 3.1** (Decision Boundary of  $f$ ). The decision boundary of a neural network  $f$  under an input distribution  $\mathcal{D}$  is the set of inputs  $\mathbf{x} \in \mathcal{D}$  for which at least two classes in  $\mathcal{M} = \{1, 2, \dots, M\}$  have equal and maximal prediction probabilities. Formally, we denote this set as  $\mathcal{B}_M^{(f, \mathcal{D})}$ , defined by:

$$\mathcal{B}_M^{(f, \mathcal{D})} = \left\{ \mathbf{x} \in \mathcal{D} \mid \exists m, n \in \mathcal{M}, m \neq n, \text{ such that } p_m = p_n \text{ and } p_m \geq \max_{o \in \mathcal{M} \setminus \{m, n\}} p_o \right\}, \quad (1)$$

where  $p_i = P[i] = \sigma(f_{\text{cls}}(f_r(\mathbf{x}))) [i]$  is the predicted probability for class  $i$ .

Based on Definition 3.1, we characterize the decision boundary for multi-class classification scenarios as follows.

**Theorem 3.2** (Properties of Multi-Class Classification Boundary). *For multi-class classification ( $M > 2$ ), the decision boundary of  $f$  can be expressed as:*

$$\mathcal{B}_M^{(f, \mathcal{D})} = \bigcup_{1 \leq m < n \leq M} \mathcal{B}_{mn}, \quad (2)$$

$$\mathcal{B}_{mn} = \{ \mathbf{x} \mid (w_m - w_n)h + (b_m - b_n) = 0, z_m = z_n \geq z_o \forall o \neq m, n, h = f_r(\mathbf{x}), \mathbf{x} \in \mathcal{D} \}.$$

where  $z = W_{\text{cls}}h + b_{\text{cls}}$  is the logits,  $w_m$  and  $w_n$  are the  $m$ -th and  $n$ -th rows of  $W_{\text{cls}}$ , and  $b_m, b_n$  are the corresponding entries of  $b_{\text{cls}}$ .

Geometrically,  $\mathcal{B}_M$  induces a Voronoi partition of the representation space, where each class corresponds to a Voronoi cell.

The proof of Theorem 3.2 is provided in Appendix B.1.

### 3.2 DECISION BOUNDARY FOR LARGE LANGUAGE MODELS

An LLM  $f : \mathcal{V}^{N_q} \rightarrow \mathcal{V}^{N_r}$  generates a sequence of tokens  $\mathbf{y} = [y_1, \dots, y_{N_r}]$ , where each token  $y_t \in \mathcal{V} = \{1, 2, \dots, V\}$  is drawn from a vocabulary of size  $V$ , conditioned on an input prompt  $\mathbf{x} = [x_1, \dots, x_{N_q}] \in \mathcal{V}^{N_q}$ .  $N_q$  and  $N_r$  are sequence length of the input and the generated texts. At each generation step  $t$ , the LLM predicts the next token  $y_t$  based on the prompt and previously generated tokens, i.e.,  $y_t \sim P_f(y_t | \mathbf{x}, y_1, \dots, y_{t-1})$ . This single step generation can be viewed as a multi-class classification over  $\mathcal{V}$ , and thus, the single-token decision boundary follows Theorem 3.2. When defining the decision boundary for the entire sequence  $\mathbf{y} \in \mathcal{V}^{N_r}$ , we need to firstly model the joint probability of the sequence under the autoregressive process. We derive the decision boundary of LLMs from that of multi-classification, as shown below.

**Theorem 3.3** (Decision Boundary of Language Models). *The decision boundary of an LLM  $f$  under an input text distribution  $\mathcal{D}' \subseteq \bigcup_{n_q=1}^{N_q} \mathcal{V}^{n_q}$  is the set of prompts  $\mathbf{x} \in \mathcal{D}'$  that lead to **equal generation probabilities** for at least two distinct sequences  $\mathbf{y}_v, \mathbf{y}_w \in \mathcal{V}^{N_r}$ , with their probabilities being maximal. Formally, the decision boundary  $\mathcal{B}_{llm}^{(f, \mathcal{D}')}$  is:*

$$\mathcal{B}_{llm}^{(f, \mathcal{D}')} = \bigcup_{\mathbf{y}_v \neq \mathbf{y}_w \in \mathcal{V}^{N_r}} \mathcal{B}_{llm, v, w}, \quad (3)$$

$$\mathcal{B}_{llm, v, w} = \left\{ \mathbf{x} \in \mathcal{D}' \mid P_f(\mathbf{y}_v | \mathbf{x}) = P_f(\mathbf{y}_w | \mathbf{x}) \geq \max_{\mathbf{y}_u \in \mathcal{V}^{N_r} \setminus \{\mathbf{y}_v, \mathbf{y}_w\}} P_f(\mathbf{y}_u | \mathbf{x}) \right\},$$

where  $P_f(\mathbf{y} | \mathbf{x}) = \prod_{t=1}^{N_r} P_f(y_t | \mathbf{x}, y_1, \dots, y_{t-1})$  is the joint probability of generating sequence  $\mathbf{y}$  given prompt  $\mathbf{x}$ .

The proof of Theorem 3.3 is provided in Appendix B.2.

While Theorem 3.3 provides a concise and intuitive definition of the decision boundary for LLMs, analyzing or computing this boundary could be computationally impossible in practice. As analyzed in Section 1, the primary challenge stems from the large vocabulary size and the autoregressive nature of sequence generation, i.e., for a generation of length  $N_r$ , the total number of possible sequences is  $V^{N_r}$ , leading to **an exponential growth in the classification space**. Specifically, the decision boundary defined in Equation 3 involves comparing pairs of sequences  $\mathbf{y}_v, \mathbf{y}_w \in \mathcal{V}^{N_r}$ , resulting in up to  $\binom{V^{N_r}}{2} \approx \frac{(V^{N_r})^2}{2}$  pairwise comparisons. This is neither computationally feasible nor interpretable on subsequent visualization.

Given such intractability of directly analyzing the decision boundary defined in Theorem 3.3, a new strategy for constructing the decision boundary of large language models is essential. Specifically, we hope that this new construction satisfies the following criteria: First, it must be *theoretically rigorous*, meaning the construction should be equivalent to or provide a bounded approximation of the decision boundary defined in Theorem 3.3, ensuring consistency with the formal definition of the boundary separating prompts that yield different output sequences. Second, the method should be *practical*, meaning it must be computationally efficient and feasible for implementation, enabling the modeling of decision boundaries for industrial-scale LLMs with large vocabularies and long generation lengths. Third, the method should be *interpretable*, meaning the constructed decision boundary should explicitly capture key properties of LLMs (e.g., curvature), and provide interpretable insights into phenomena observed in LLM behavior, such as output variability or robustness.

In the next section, we will introduce an approximation procedure for the decision boundary defined in Theorem 3.3, addressing these criteria to enable practical and meaningful analysis of LLMs.

## 4 DECISION POTENTIAL SURFACE: SIMPLIFYING DECISION BOUNDARY ANALYSIS ON LARGE LANGUAGE MODELS

In this section, we introduce the *Decision Potential Surface (DPS)*, a novel concept for analyzing the decision boundaries of LLMs by representing the *decision potential* of generated sequences as a surface over the input manifold. In Section 4.1, we formally define DPS and establish its relationship with the standard decision boundary formulation in LLMs. In Section 4.2, we propose  $K$ -grained DPS ( $K$ -DPS), a practical approximation of DPS, and theoretically derive its error bounds with respect to the ideal DPS.

### 4.1 DECISION POTENTIAL SURFACE OF LANGUAGE MODELS

**Definition 4.1** (Decision Potential Surface of Language Models). Given an input text distribution  $\mathbf{x} \in \mathcal{D}'$  with  $\mathcal{D}' \subseteq \bigcup_{n_q=1}^{N_q} \mathcal{V}^{n_q}$  and a language model  $f : \mathcal{V}^{N_q} \rightarrow \mathcal{V}^{N_r}$  that generates an output sequence  $\mathbf{y} = f(\mathbf{x}) = \arg \max_{\mathbf{y}_s \in \mathcal{V}^{N_r}} P_f(\mathbf{y}_s | \mathbf{x})$ , we define the *decision potential function (DPF)*  $\Phi_f^\infty(\mathbf{x}) : \mathcal{D}' \rightarrow \mathbb{R}_+$  as the squared difference in log-likelihoods between the top two generated sequences under the input prompt  $\mathbf{x}$ , i.e.,

$$\begin{aligned} \Phi_f^\infty(\mathbf{x}) &= \left( \min_{\mathbf{y}_w \in \mathcal{V}^{N_r}, \mathbf{y}_w \neq \mathbf{y}_v} \left[ \max_{\mathbf{y}_v \in \mathcal{V}^{N_r}} \log P_f(\mathbf{y}_v | \mathbf{x}) - \log P_f(\mathbf{y}_w | \mathbf{x}) \right] \right)^2 \\ &= (\log P_f(\mathbf{y}_{1*} | \mathbf{x}) - \log P_f(\mathbf{y}_{2*} | \mathbf{x}))^2, \end{aligned} \quad (4)$$

where  $\mathbf{y}_{1*}, \mathbf{y}_{2*} \in \mathcal{V}^{N_r}$  denote the sequences with the highest and second-highest log-likelihoods, respectively. The *decision potential surface (DPS)* is then defined as  $\mathcal{S}^{(f, \mathcal{D}')} := \{\Phi_f^\infty(\mathbf{x}) \mid \mathbf{x} \in \mathcal{D}'\}$ .

Intuitively,  $\mathcal{S}^{(f, \mathcal{D}')}$  can be viewed as a surface representing the competitive likelihoods across all inputs, where each decision potential value  $\Phi_f^\infty(\mathbf{x})$  quantifies the *confidence* in distinguishing the most likely sequence.

Then, we define *isohypses* (i.e., contour lines) on surface  $\mathcal{S}^{(f, \mathcal{D}')}$  as follows:

**Definition 4.2** ( $\varepsilon$ -Isohypse). The  $\varepsilon$ -isohypse on the decision potential surface  $\mathcal{S}^{(f, \mathcal{D}')}$  is the set of inputs with the same decision potential value  $\varepsilon$ , i.e.,

$$\mathcal{D}'_{(\varepsilon, f)} = \{\mathbf{x} \mid \mathbf{x} \in \mathcal{D}'; \Phi_f^\infty(\mathbf{x}) = \varepsilon\}. \quad (5)$$

As a degenerate case, the zero level set of  $\mathcal{S}^{(f, \mathcal{D}')}$  exhibits the following property:

**Theorem 4.3** (0-Isohypse as the Decision Boundary). *The decision boundary of a language model  $f(\mathbf{x})$  under  $\mathcal{D}'$ , as defined in Theorem 3.3, is equivalent to the 0-isohypse, i.e.,*

$$\mathcal{B}_{\text{ltm}}^{(f, \mathcal{D}')} = \mathcal{D}'_{(0, f)} = \{\mathbf{x} \in \mathcal{D}' \mid \Phi_f^\infty(\mathbf{x}) = 0\}, \quad (6)$$

where regions separated by the 0-isohypse correspond exactly to the Voronoi cells.

We also provide the following corollary to characterize the surface structure:

**Corollary 4.4** ( $\varepsilon$ -Isohypse Gives  $\sqrt{\varepsilon}$ -nat Confidence Hierarchy). *For any  $\varepsilon > 0$ , the input space  $\mathcal{D}'$  is partitioned into three disjoint strata:*

- $\varepsilon$ -barriers:  $\mathcal{D}'_{(>\varepsilon, f)} = \{\mathbf{x} \mid \Phi_f^\infty(\mathbf{x}) > \varepsilon; \mathbf{x} \in \mathcal{D}'\}$ , where  $f(\mathbf{x})$  predicts the sequence of its region with at least  $\sqrt{\varepsilon}$  nats (natural units of information) of confidence over the next most likely sequence.
- $\varepsilon$ -well:  $\mathcal{D}'_{(<\varepsilon, f)} = \{\mathbf{x} \mid \Phi_f^\infty(\mathbf{x}) < \varepsilon; \mathbf{x} \in \mathcal{D}'\}$ , where  $f(\mathbf{x})$  has low confidence, with a margin less than  $\sqrt{\varepsilon}$  nats. As  $\varepsilon \rightarrow 0$ , this stratum converges to the 0-isohypse.
- $\varepsilon$ -isohypse:  $\mathcal{D}'_{(\varepsilon, f)} = \{\mathbf{x} \in \mathcal{D}' \mid \Phi_f^\infty(\mathbf{x}) = \varepsilon\}$ , representing the contour where the confidence margin is exactly  $\sqrt{\varepsilon}$  nats.

Proofs are provided in Appendix B.3 and B.4, respectively.

Unfortunately, computing the decision boundary or visualizing the potential surface based on Definition 4.1 and Theorem 4.3 remains computationally infeasible, as evaluating  $\Phi_f^\infty(\mathbf{x})$  in Equation 4 requires considering all possible sequences in  $\mathcal{V}^{N_r}$ , resulting in a computational complexity the same as before.

Fortunately, as Equation 4 depends only on the log-likelihoods of the top two sequences, we can propose an efficient approximation with a modest error, detailed in the next subsection.

## 4.2 $K$ -GRAINED DECISION POTENTIAL SURFACE AND ITS PROPERTIES

We introduce  $K$ -grained decision potential surface for approximating  $\mathcal{S}^{(f, \mathcal{D}')}$ :

**Definition 4.5** ( $K$ -Grained Decision Potential Surface). Given  $\mathbf{x} \in \mathcal{D}'$  and a language model  $f(\mathbf{x})$ , we define the  $K$ -grained potential function  $\Phi_f^K(\mathbf{x}) : \mathcal{D}' \rightarrow \mathbb{R}_+$  as

$$\begin{aligned} \Phi_f^K(\mathbf{x}) &= \left( \min_{\mathbf{y}_w \in \mathcal{Y}_K, \mathbf{y}_w \neq \mathbf{y}_v} \left[ \max_{\mathbf{y}_v \in \mathcal{Y}_K} \log P_f(\mathbf{y}_v | \mathbf{x}) - P_f(\mathbf{y}_w | \mathbf{x}) \right] \right)^2 \\ &= (\log P_f(\mathbf{y}_{1*}^K | \mathbf{x}) - \log P_f(\mathbf{y}_{2*}^K | \mathbf{x}))^2, \end{aligned} \quad (7)$$

where  $1 \ll K \ll V^{N_r}$  denotes the size of output space for each input,  $\mathcal{Y}_K = \{\mathbf{y}_v \sim P_f(\cdot | \mathbf{x}) \mid v = 1, \dots, K\}$  denotes  $K$  i.i.d. (independent and identically distributed) sampled texts,  $\mathbf{y}_{1*}^K$  and  $\mathbf{y}_{2*}^K$  denotes the top-2 generated texts owning the maximal generation logarithmic likelihood within  $\mathcal{Y}_K$ .

In this way, the computational complexity of constructing the decision boundary is reduced from  $\mathcal{O}(V^{2N_r} \cdot \mathcal{D}')$  to  $\mathcal{O}(K^2 \cdot \mathcal{D}')$ , resulting in a substantial reduction. This naturally leads to the next question: what is the error between  $\Phi_f^K(\mathbf{x})$  and  $\Phi_f^\infty(\mathbf{x})$ ? We address this by theoretically analyzing their relationship in the following theorems.

**Theorem 4.6** (Error Bound for Estimating  $\Phi_f^\infty(\mathbf{x})$  with  $\Phi_f^K(\mathbf{x})$ ). *For a fixed input  $\mathbf{x} \in \mathcal{D}'$  and a set  $\mathcal{Y}_K$  of  $K$  i.i.d. samples drawn from the language model's output distribution  $P_f(\cdot | \mathbf{x})$ , suppose the population top-2 gap satisfies  $\Delta_\infty(\mathbf{x}) = \log P_f(\mathbf{y}_{1*} | \mathbf{x}) - \log P_f(\mathbf{y}_{2*} | \mathbf{x}) \leq R_K(\mathbf{x})$ , where*

324  $R_K(\mathbf{x}) = \log P_f(\mathbf{y}_{1*}^K|\mathbf{x}) - \min_{\mathbf{y} \in \mathcal{Y}_K} \log P_f(\mathbf{y}|\mathbf{x})$  represents the log-likelihood diameter of  $\mathcal{Y}_K$ .  
 325 Then, for any  $\delta \in (0, 1)$ , the error between the sample-based decision potential  $\Phi_f^K(\mathbf{x})$  and the true  
 326 decision potential  $\Phi_f^\infty(\mathbf{x})$  satisfies:  
 327

$$328 \quad |\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \leq 2R_K^2(\mathbf{x}) \sqrt{\frac{\log(4/\delta)}{2K}}, \quad (8)$$

329 with probability at least  $1 - \delta - 2\varepsilon_{tail}$ , where  $\varepsilon_{tail} = (1 - P_f(\mathbf{y}_{1*}^K|\mathbf{x}))^K$ .

332 **Theorem 4.7** (Expected Error Bound). *Under the same conditions of Theorem 4.6, the expected*  
 333 *error between the sample-based decision potential  $\Phi_f^K(\mathbf{x})$  and the true decision potential  $\Phi_f^\infty(\mathbf{x})$  is*  
 334 *bounded as:*

$$335 \quad \mathbb{E} [|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})|] \leq 2R_K^2(\mathbf{x}) \sqrt{\frac{2\pi}{K}} + 4R_K^2(\mathbf{x})\varepsilon_{tail}, \quad (9)$$

337 where  $\varepsilon_{tail} = (1 - P_f(\mathbf{y}_{1*}^K|\mathbf{x}))^K$ .

339 **Corollary 4.8** (Concentration Bound). *Under the same conditions as Theorem 4.6, for any  $\lambda > 0$ ,*  
 340 *the tail probability of the error satisfies:*

$$341 \quad \Pr (|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \geq \lambda) \leq 4 \exp\left(-\frac{K\lambda^2}{2R_K^4(\mathbf{x})}\right) + 2\varepsilon_{tail}, \quad (10)$$

342 where  $\varepsilon_{tail} = (1 - P_f(\mathbf{y}_{1*}^K|\mathbf{x}))^K$ .

343 Proofs are provided in Appendix B.5, B.6 and B.7, respectively.

344 From the above theorems we can see that the estimation error contracts with  $K$  at the familiar  $1/\sqrt{K}$   
 345 rate, mirroring the decay of an empirical mean, while the residual tail probability  $\varepsilon_{tail}$  is rendered  
 346 exponentially negligible by the fact that the model’s top-1 sentence carries almost all mass, so the  
 347 chance that it is missed in  $K$  *i.i.d.* draws is effectively zero even for modest  $K$ . Moreover, the  
 348 common factor  $R_K(\mathbf{x})$  is a worst-case log-likelihood diameter whose width is dictated by the most  
 349 unlikely sentence that happens to be sampled, and a sharper analysis could clearly replace this global  
 350 spread with a more aggressive local gap (e.g.,  $R'_K(\mathbf{x}) = \log P_f(\mathbf{y}_{1*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{3*}^K|\mathbf{x})$  with  $\mathbf{y}_{3*}^K$   
 351 the third most likely generated sentence in  $\mathcal{Y}_K$ ) without inflating the bound, thereby tightening the  
 352 constants while preserving the same decay.

## 353 5 EMPIRICAL ANALYSIS

### 354 5.1 SETTINGS

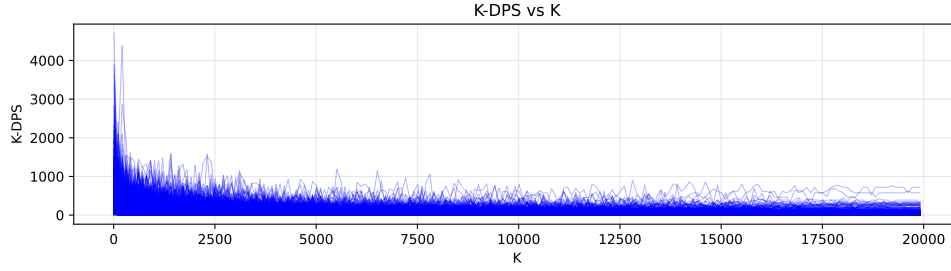
355 **Datasets and Models.** We utilize both pre-training corpora and supervised fine-tuning (SFT)  
 356 datasets to simulate the input data distribution for constructing decision boundaries and the decision  
 357 potential surface. For the pre-training corpus, we select Wikipedia Mini (Ridder & Schilling, 2025),  
 358 an unsupervised text corpus containing a condensed version of Wikipedia articles. For supervised  
 359 fine-tuning, we employ Tulu-3-SFT-MIX (Lambert et al., 2025), OpenO1-SFT (Xia et al., 2025a),  
 360 HH-RLHF (Ganguli et al., 2022), and Alpaca (Taori et al., 2023), all of which are widely used in  
 361 academic and industrial settings. We use Llama3.2-1B (Grattafiori et al., 2024) as the backbone in  
 362 experiments.

363 **Implementation Details.** For sampling, we use *nucleus sampling* in our  $K$ -DPS implementation,  
 364 with the clipping probability  $p$  set to 0.9. In subsequent experiments, each data point is repeated five  
 365 times. The experiments are conducted on  $4 \times 80\text{GB}$  Nvidia Tesla H100 GPUs.

### 366 5.2 INFLUENCE OF SAMPLING GRAIN $K$

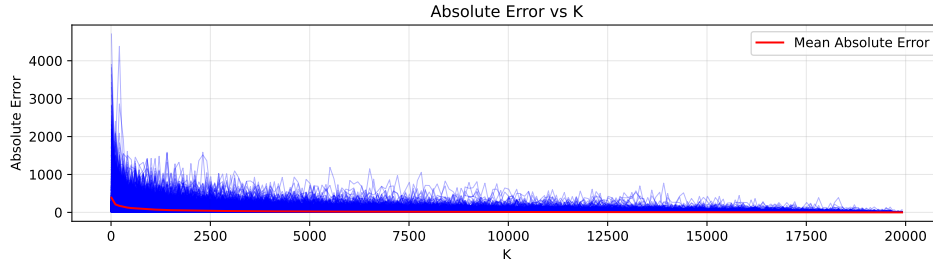
367 We first evaluate the impact of the key hyperparameter, the sampling grain  $K$ , on the  $K$ -DPS value  
 368 and the absolute errors between  $K$ -DPS and the ideal DPS. Specifically, we set  $K = 20,000$ ,  
 369 a sufficiently large value, to simulate and approximate the ideal DPS. We then compute the  $K$ -  
 370 DPS values by varying  $K$  from 10 to 20,000 to illustrate how the decision potential value  $\Phi_f^K(\cdot)$   
 371

378  
379  
380  
381  
382  
383  
384  
385  
386



387 Figure 1: Effect of sampling size  $K$  on the values of decision potential function, with each blue point representing the  $K$ -DPS value for a single input sample. Each blue line represents a trend of  $K$ -DPS for one input sample.

391  
392  
393  
394  
395  
396  
397  
398  
399  
400



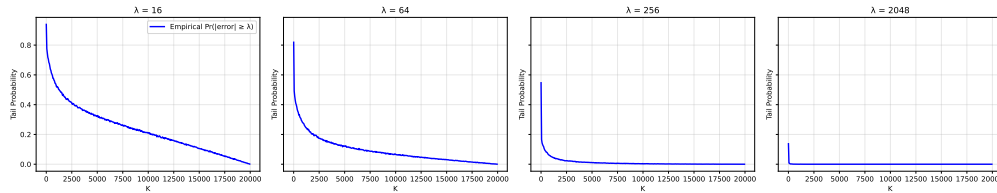
401 Figure 2: Effect of sampling size  $K$  on the absolute error between the reference  $K$ -DPS (computed with  $K = 20,000$ ) and  $K$ -DPS values for varying  $K$ . Each blue line represents a trend of absolute error across input samples.

405  
406  
407  
408  
409  
410  
411  
412

converges to the ideal potential  $\Phi_f^\infty(\cdot)$ . Similarly, we calculate the absolute errors of  $\Phi_f^K(\cdot)$  across different settings of  $K$ . As shown in Figure 1, the potential values rapidly converge to their true values (represented by horizontal lines in the tails), indicating that a relatively small  $K$  can yield a highly accurate decision potential surface. Moreover, by examining the errors defined in Equation 8, as depicted in Figure 2, we observe that both the absolute error for individual samples and the empirically average error decrease to zero, confirming the effectiveness of  $K$ -DPS. Figures 1 and 2 also serve as valuable references for selecting appropriate  $K$  values.

### 413 5.3 EMPIRICAL CONCENTRATION BIAS

414  
415  
416  
417  
418  
419  
420  
421  
422



423 Figure 3: Empirical concentration experiments with different  $\lambda$  values.

424  
425  
426  
427  
428  
429  
430  
431

We also present an empirical study of concentration experiments, focusing on the trend of sample probabilities for inputs with a decision potential error exceeding a given fixed value  $\lambda$  across various sampling sizes  $K$ . As shown in Figure 3, we evaluate the tail probability for  $K$  values ranging from 10 to 20,000, with  $\lambda$  set to 16, 64, 256, and 2048. These  $\lambda$  values represent the geometric errors between the approximate and ideal DPS values. It is noteworthy to emphasize that even a  $\lambda$  value of 256 is not excessively large or insignificant, as our decision potential function  $\Phi_f^K(\cdot)$  is defined as the **square** of logarithmic errors, as specified in Equation 7.

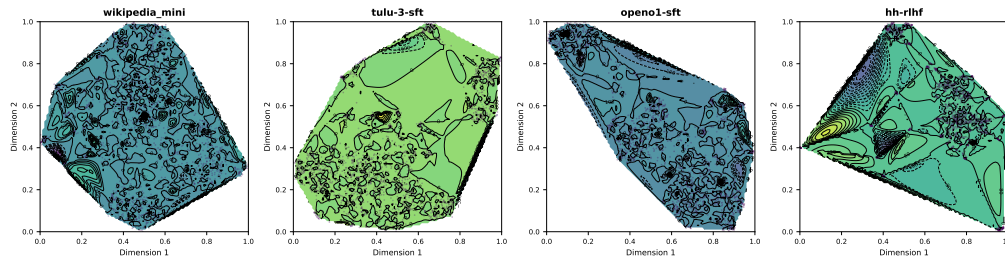


Figure 4: Contour visualization of the  $K$ -DPS ( $K = 2,500$ ) for Llama-3.2-1B on four datasets. Region colors represent the decision potential values. Black lines denote isohypses, with the 0-isohypse indicating decision boundary. Cubic interpolation is applied to construct the mesh grid, with visualizations using linear and nearest interpolation shown in Figures 10 and 11.

From Figure 3, we observe that the tail probabilities exhibit an exponential decrease, indicating that the likelihood of exceeding a given error bound diminishes significantly with a linear increase in the sampling size  $K$ . Specifically, Figure 3 demonstrates that a sampling size of 10,000 ensures an absolute error below 64 with 90% confidence and an error below 256 with 99% probability. These results align closely with our absolute error analysis presented in Figure 2.

#### 5.4 VISUALIZING DECISION POTENTIAL SURFACE

While this paper primarily focuses on the error analysis of LLMs’ decision boundary construction, our proposed  $K$ -DPS can also be used to intuitively visualize both the decision boundary and the decision potential surface of an LLM under a given input distribution, as detailed in this section.

**Settings.** For visualization, we construct a low-dimensional representation of the original input distribution  $\mathcal{D}'$ , typically in two dimensions to facilitate human understanding. First, we extract the last hidden state of an input  $\mathbf{x}$  from the LLM as the original embedding of the input point. Next, we apply UMAP with 100 neighbors and a minimum distance of 0.2 for dimensionality reduction. Finally, we normalize the reduced embeddings to the range  $[0, 1]$  to construct the decision potential surface visualization. For interpolation, we evaluate nearest, linear, and cubic interpolation methods to approximate the  $K$ -DPS values on a mesh grid.

**Visualization Results.** As shown in Figure 4, we visualize the decision potential surface of the pretrained Llama-3.2-1B model on four corpora: Wikipedia Mini, Tulu-3-SFT, OpenO1-SFT, and HH-RLHF, using cubic interpolation. The heights of the isohypses are marked in the figures. The sampling size of  $K = 2,500$ . From Figure 4, it is evident that most contours are at zero height, indicating that  $K$ -DPS effectively captures the decision boundary of large language models. Consequently, properties of the decision boundary (i.e., the 0-isohypse), such as curvature, location, and density, can be readily analyzed to facilitate interpretation and analysis of LLMs for future studies.

However, due to limitations in the interpolation strategy, the decision potential surface may be invalid in regions with sparse input points. For instance, the top-left regions of the second and fourth subfigures show  $\Phi^K$  values significantly below zero, reflecting interpolation errors due to the absence of input samples in these areas. Additional visualizations with other interpolation methods are provided in the Appendix C for reference.

## 6 CONCLUSION

In this study, we explore the construction of decision boundaries for large language models. We identify the primary challenges as stemming from the expansive vocabulary size and the exponential growth in generated token sequences. To overcome these obstacles, we propose the concept of decision potential surface to quantify the confidence of language models in their decisions. We theoretically demonstrate that the zero-height isohypse on this surface corresponds to the decision boundary, and its approximated implementation substantially reduces the computational complexity of constructing the decision boundary. Through rigorous theoretical and empirical analyses, we evaluate the errors and validate the effectiveness of the proposed method.

---

## REPRODUCIBILITY STATEMENT

As a theoretical study, we have clearly articulated all assumptions, definitions, theorems, and proofs in the main text and the Appendix. For the empirical results, we have provided the source code<sup>1</sup> to facilitate straightforward reproduction of the experimental findings. We welcome any additional suggestions to enhance the reproducibility of this work.

## REFERENCES

- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.2, knowledge manipulation. *arXiv preprint arXiv:2309.14402*, 2023.
- Emmanuel Ameisen, Jack Lindsey, Adam Pearce, Wes Gurnee, Nicholas L Turner, Brian Chen, Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, et al. Circuit tracing: Revealing computational graphs in language models. *Transformer Circuits Thread*, 6, 2025.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic language model. *J. Mach. Learn. Res.*, 3:1137–1155, 2003. URL <https://jmlr.org/papers/v3/bengio03a.html>.
- Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/34e1dbe95d34d7ebaf99b9bcaeb5b2be-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/34e1dbe95d34d7ebaf99b9bcaeb5b2be-Abstract-Conference.html).
- 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. *Transformer Circuits Thread*, 1(1):12, 2021.
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630, 2024.
- Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, Andy Jones, Sam Bowman, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Nelson Elhage, Sheer El-Showk, Stanislav Fort, Zac Hatfield-Dodds, Tom Henighan, Danny Hernandez, Tristan Hume, Josh Jacobson, Scott Johnston, Shauna Kravec, Catherine Olsson, Sam Ringer, Eli Tran-Johnson, Dario Amodei, Tom Brown, Nicholas Joseph, Sam McCandlish, Chris Olah, Jared Kaplan, and Jack Clark. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned, 2022. URL <https://arxiv.org/abs/2209.07858>.
- Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koepl, Preslav Nakov, and Iryna Gurevych. A survey of confidence estimation and calibration in large language models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 6577–6595, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.366. URL <https://aclanthology.org/2024.naacl-long.366/>.
- Jiahui Geng, Qing Li, Herbert Woitschlaeger, Zongxiong Chen, Yuxia Wang, Preslav Nakov, Hans-Arno Jacobsen, and Fakhri Karray. A comprehensive survey of machine unlearning techniques for large language models. *CoRR*, abs/2503.01854, 2025. doi: 10.48550/ARXIV.2503.01854. URL <https://doi.org/10.48550/arXiv.2503.01854>.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *3rd International Conference on Learning Representations (ICLR 2015)*, San Diego, CA, USA, 2015.

---

<sup>1</sup><https://github.com/liangzid/DPS>

---

540 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
541 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan,  
542 Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Ko-  
543 renev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava  
544 Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux,  
545 Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret,  
546 Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius,  
547 Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary,  
548 Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab  
549 AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco  
550 Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind That-  
551 tai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Kore-  
552 vaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra,  
553 Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-  
554 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,  
555 Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jong-  
556 soo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala,  
557 Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid  
558 El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren  
559 Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin,  
560 Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi,  
561 Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew  
562 Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar  
563 Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev,  
564 Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan  
565 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan,  
566 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon  
567 Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit  
568 Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan  
569 Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell,  
570 Seohyun Sonia Kim, Sergey Edunov, Shao-liang Nie, Sharan Narang, Sharath Rparathy, Sheng  
571 Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer  
572 Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman,  
573 Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mi-  
574 haylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor  
575 Kerkez, Vincent Gouget, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei  
576 Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang  
577 Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Gold-  
578 schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning  
579 Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papanikos, Aaditya Singh,  
580 Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria,  
581 Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein,  
582 Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew  
583 Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie  
584 Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,  
585 Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leon-  
586 hardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu  
587 Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Mon-  
588 talvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao  
589 Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia  
590 Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide  
591 Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le,  
592 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily  
593 Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers,  
Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni,  
Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia  
Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan,  
Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harri-  
son Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj,

---

594 Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James  
595 Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-  
596 nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang,  
597 Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Jun-  
598 jie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy  
599 Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang,  
600 Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell,  
601 Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa,  
602 Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias  
603 Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L.  
604 Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike  
605 Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari,  
606 Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan  
607 Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong,  
608 Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent,  
609 Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar,  
610 Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Ro-  
611 driguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,  
612 Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin  
613 Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon,  
614 Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ra-  
615 maswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha,  
616 Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal,  
617 Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satter-  
618 field, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj  
619 Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo  
620 Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook  
621 Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Ku-  
622 mar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov,  
623 Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiao-  
624 jian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,  
625 Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao,  
626 Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhao-  
627 duo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL  
628 <https://arxiv.org/abs/2407.21783>.  
629

630 Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the  
631 machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.  
632

633 Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and  
634 Lei Ma. Look before you leap: An exploratory study of uncertainty analysis for large language  
635 models. *IEEE Transactions on Software Engineering*, 51(2):413–429, February 2025. ISSN  
636 2326-3881. doi: 10.1109/tse.2024.3519464. URL [http://dx.doi.org/10.1109/TSE.](http://dx.doi.org/10.1109/TSE.2024.3519464)  
637 [2024.3519464](http://dx.doi.org/10.1109/TSE.2024.3519464).  
638

639 Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez,  
640 Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer  
641 El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bow-  
642 man, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna  
643 Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom  
644 Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Kap-  
645 lan. Language models (mostly) know what they know, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2207.05221)  
646 [abs/2207.05221](https://arxiv.org/abs/2207.05221).  
647

648 Hamid Karimi and Jiliang Tang. Decision boundary of deep neural networks: Challenges and op-  
649 portunities. In *Proceedings of the 13th International Conference on Web Search and Data Min-*  
650 *ing*, WSDM '20, pp. 919–920, New York, NY, USA, 2020. Association for Computing Machin-  
651 ery. ISBN 9781450368223. doi: 10.1145/3336191.3372186. URL [https://doi.org/10.](https://doi.org/10.1145/3336191.3372186)  
652 [1145/3336191.3372186](https://doi.org/10.1145/3336191.3372186).

---

648 Hamid Karimi, Tyler Derr, and Jiliang Tang. Characterizing the decision boundary of deep neural  
649 networks, 2020. URL <https://arxiv.org/abs/1912.11460>.

650

651 Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for  
652 uncertainty estimation in natural language generation, 2023. URL <https://arxiv.org/abs/2302.09664>.

653

654 Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahma-  
655 man, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Ma-  
656 lik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris  
657 Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Ha-  
658 jishirzi. Tulu 3: Pushing frontiers in open language model post-training, 2025. URL <https://arxiv.org/abs/2411.15124>.

659

660 Chulhee Lee and D.A. Landgrebe. Decision boundary feature extraction for neural networks. *IEEE*  
661 *Transactions on Neural Networks*, 8(1):75–83, 1997. doi: 10.1109/72.554193.

662

663 Sungzoon Lee and B. John Oommen. Decision boundary boundary feature extraction for neural  
664 networks. *IEEE Transactions on Neural Networks*, 8(4):865–875, 1997.

665

666 Moxin Li, Yong Zhao, Wenxuan Zhang, Shuaiyi Li, Wenya Xie, See-Kiong Ng, Tat-Seng Chua,  
667 and Yang Deng. Knowledge boundary of large language models: A survey. In Wanxiang Che,  
668 Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the*  
669 *63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,  
670 pp. 5131–5157, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN  
671 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.256. URL <https://aclanthology.org/2025.acl-long.256/>.

672

673 Yu Li, Lizhong Ding, and Xin Gao. On the decision boundary of deep neural networks, 2019. URL  
674 <https://arxiv.org/abs/1808.05385>.

675

676 Zi Liang, Haibo Hu, Qingqing Ye, Yaxin Xiao, and Haoyang Li. Why are my prompts leaked?  
677 unraveling prompt extraction threats in customized large language models. *arXiv preprint*  
*arXiv:2408.02416*, 2024.

678

679 Zi Liang, Haibo Hu, Qingqing Ye, Yaxin Xiao, and Ronghua Li. Does low rank adaptation lead  
680 to lower robustness against training-time attacks?, 2025. URL <https://arxiv.org/abs/2505.12871>.

681

682 Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantifi-  
683 cation for black-box large language models, 2024. URL <https://arxiv.org/abs/2305.19187>.

684

685 Sijia Liu, Yuanshun Yao, Jinghan Jia, Stephen Casper, Nathalie Baracaldo, Peter Hase, Yuguang  
686 Yao, Chris Yuhao Liu, Xiaojun Xu, Hang Li, Kush R. Varshney, Mohit Bansal, Sanmi Koyejo,  
687 and Yang Liu. Rethinking machine unlearning for large language models. *Nat. Mac. Intell.*, 7(2):  
688 181–194, 2025a. doi: 10.1038/S42256-025-00985-0. URL <https://doi.org/10.1038/s42256-025-00985-0>.

689

690 Xiaoou Liu, Tiejun Chen, Longchao Da, Chacha Chen, Zhen Lin, and Hua Wei. Uncertainty quan-  
691 tification and confidence calibration in large language models: A survey. In *Proceedings of*  
692 *the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V.2, KDD '25*,  
693 pp. 6107–6117, New York, NY, USA, 2025b. Association for Computing Machinery. ISBN  
694 9798400714542. doi: 10.1145/3711896.3736569. URL <https://doi.org/10.1145/3711896.3736569>.

695

696 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
697 Towards deep learning models resistant to adversarial attacks. In *6th International Conference on*  
698 *Learning Representations (ICLR 2018)*, Vancouver, BC, Canada, 2018.

699

700 Harry Mayne, Ryan Othniel Kearns, Yushi Yang, Andrew M. Bean, Eoin Delaney, Chris Russell, and  
701 Adam Mahdi. LLMs don’t know their own decision boundaries: The unreliability of self-generated  
counterfactual explanations, 2025. URL <https://arxiv.org/abs/2509.09396>.

---

702 David Mickisch, Felix Assion, Florens Greßner, Wiebke Günther, and Mariele Motta. Under-  
703 standing the decision boundary of deep neural networks: An empirical study, 2020. URL  
704 <https://arxiv.org/abs/2002.01810>.  
705

706 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language under-  
707 standing by generative pre-training. 2018.

708 Fabian Ridder and Malte Schilling. The hallurag dataset: Detecting closed-domain hallucinations  
709 in rag applications using an llm’s internal states, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2412.17056)  
710 [2412.17056](https://arxiv.org/abs/2412.17056).  
711

712 Frank Rosenblatt. The perceptron: A probabilistic model for information storage and organization  
713 in the brain. *Psychological Review*, 65(6):386–408, 1958.

714 Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lindsey, Jeff Wu, Lucius Bushnaq, Nicholas  
715 Goldowsky-Dill, Stefan Heimersheim, Alejandro Ortega, Joseph Isaac Bloom, Stella Biderman,  
716 Adrià Garriga-Alonso, Arthur Conmy, Neel Nanda, Jessica Rumbelow, Martin Wattenberg, Nandi  
717 Schoots, Joseph Miller, Eric J. Michaud, Stephen Casper, Max Tegmark, William Saunders, David  
718 Bau, Eric Todd, Atticus Geiger, Mor Geva, Jesse Hoogland, Daniel Murfet, and Tom McGrath.  
719 Open problems in mechanistic interpretability. *CoRR*, abs/2501.16496, 2025. doi: 10.48550/  
720 ARXIV.2501.16496. URL <https://doi.org/10.48550/arXiv.2501.16496>.  
721

722 Ola Shorinwa, Zhiting Mei, Justin Lidard, Allen Z. Ren, and Anirudha Majumdar. A survey on  
723 uncertainty quantification of large language models: Taxonomy, open research challenges, and  
724 future directions. *ACM Comput. Surv.*, 58(3), September 2025. ISSN 0360-0300. doi: 10.1145/  
725 3744238. URL <https://doi.org/10.1145/3744238>.

726 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
727 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.  
728 [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.

729 K. Turner and J. Ghosh. Analysis of decision boundaries in linearly combined neural classifiers.  
730 *Pattern Recognition*, 29(2):295–307, 1996.  
731

732 Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. In-  
733 terpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The*  
734 *Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda,*  
735 *May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=NpsVSN6o4ul)  
736 [NpsVSN6o4ul](https://openreview.net/forum?id=NpsVSN6o4ul).

737 Shijie Xia, Yiwei Qin, Xuefeng Li, Yan Ma, Run-Ze Fan, Steffi Chern, Haoyang Zou, Fan Zhou,  
738 Xiangkun Hu, Jiahe Jin, Yanheng He, Yixin Ye, Yixiu Liu, and Pengfei Liu. Generative ai act  
739 ii: Test time scaling drives cognition engineering, 2025a. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2504.13828)  
740 [2504.13828](https://arxiv.org/abs/2504.13828).

741 Zhiqiu Xia, Jinxuan Xu, Yuqian Zhang, and Hang Liu. A survey of uncertainty estimation methods  
742 on large language models, 2025b. URL <https://arxiv.org/abs/2503.00172>.  
743

744 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
745 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,  
746 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin  
747 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,  
748 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui  
749 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang  
750 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger  
751 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan  
752 Qiu. Qwen3 technical report, 2025a. URL <https://arxiv.org/abs/2505.09388>.

753 Junxiao Yang, Jinzhe Tu, Haoran Liu, Xiaoce Wang, Chujie Zheng, Zhexin Zhang, Shiyao Cui,  
754 Caishun Chen, Tiantian He, Hongning Wang, Yew-Soon Ong, and Minlie Huang. Barrel:  
755 Boundary-aware reasoning for factual and reliable llms, 2025b. URL [https://arxiv.org/](https://arxiv.org/abs/2505.13529)  
[abs/2505.13529](https://arxiv.org/abs/2505.13529).

---

756 Jin Yao, Eli Chien, Minxin Du, Xinyao Niu, Tianhao Wang, Zezhou Cheng, and Xiang Yue. Ma-  
757 chine unlearning of pre-trained large language models. In Lun-Wei Ku, Andre Martins, and  
758 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-  
759 putational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16,  
760 2024*, pp. 8403–8419. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.  
761 ACL-LONG.457. URL <https://doi.org/10.18653/v1/2024.acl-long.457>.

762 Roozbeh Yousefzadeh and Dianne P O’Leary. Investigating decision boundaries of trained neural  
763 networks, 2019. URL <https://arxiv.org/abs/1908.02802>.

764

765 Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From catas-  
766 trophic collapse to effective unlearning. *CoRR*, abs/2404.05868, 2024. doi: 10.48550/ARXIV.  
767 2404.05868. URL <https://doi.org/10.48550/arXiv.2404.05868>.

768 Siyan Zhao, Tung Nguyen, and Aditya Grover. Probing the decision boundaries of in-  
769 context learning in large language models. In A. Globerson, L. Mackey, D. Bel-  
770 grave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural In-  
771 formation Processing Systems*, volume 37, pp. 130408–130432. Curran Associates, Inc.,  
772 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/  
773 file/eb5dd4476448c44e55a759a985b3bbec-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/eb5dd4476448c44e55a759a985b3bbec-Paper-Conference.pdf).

774

775 Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial  
776 attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

---

## 810 A LLM USAGE

811 It is used for error checking, proofreading, result visualization, and code optimization.

## 814 B PROOFS

### 815 B.1 PROOF OF THEOREM 3.2

816 *Proof. Part I: Proof of Equation 2.*

817 We aim to characterize the decision boundary  $\mathcal{B}_M^{(f, \mathcal{D})}$  for a neural network  $f : \mathbb{R}^d \rightarrow \mathbb{R}^M$  in the  
818 multi-class classification setting ( $M > 2$ ) under an input distribution  $\mathcal{D} \subseteq \mathbb{R}^d$ . The network is  
819 decomposed as  $f = \sigma \circ f_{\text{cls}} \circ f_r$ , where:

- 820 •  $f_r : \mathbb{R}^d \rightarrow \mathbb{R}^{d_h}$  maps the input  $\mathbf{x}$  to a latent representation  $h = f_r(\mathbf{x})$ ,
- 821 •  $f_{\text{cls}} : \mathbb{R}^{d_h} \rightarrow \mathbb{R}^M$  is a linear classification head,  $f_{\text{cls}}(h) = W_{\text{cls}}h + b_{\text{cls}}$ , with  $W_{\text{cls}} \in \mathbb{R}^{M \times d_h}$ ,  
822  $b_{\text{cls}} \in \mathbb{R}^M$ ,
- 823 •  $\sigma : \mathbb{R}^M \rightarrow \mathbb{R}^M$  is the softmax function,  $\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^M e^{z_j}}$ , producing probabilities  $P =$   
824  $[p_1, p_2, \dots, p_M]$  with  $\sum_{i=1}^M p_i = 1$ .

825 By Definition 3.1, the decision boundary  $\mathcal{B}_M^{(f, \mathcal{D})}$  is the set of inputs  $\mathbf{x} \in \mathcal{D}$  such that there exist at  
826 least two classes  $m, n \in \mathcal{M} = \{1, 2, \dots, M\}$ ,  $m \neq n$ , with equal and maximal probabilities:

$$827 p_m = p_n \geq \max_{o \in \mathcal{M} \setminus \{m, n\}} p_o. \quad (11)$$

828 Since  $\sigma$  is the softmax function,  $p_m = \sigma(z)_m = \frac{e^{z_m}}{\sum_{j=1}^M e^{z_j}}$ , the condition  $p_m = p_n$  implies:

$$829 \frac{e^{z_m}}{\sum_{j=1}^M e^{z_j}} = \frac{e^{z_n}}{\sum_{j=1}^M e^{z_j}} \implies e^{z_m} = e^{z_n} \implies z_m = z_n. \quad (12)$$

830 The logits are given by  $z = W_{\text{cls}}h + b_{\text{cls}}$ , so:

$$831 z_m = w_m h + b_m, \quad z_n = w_n h + b_n, \quad (13)$$

832 where  $w_m, w_n$  are the  $m$ -th and  $n$ -th rows of  $W_{\text{cls}}$ , and  $b_m, b_n$  are the corresponding entries of  $b_{\text{cls}}$ .  
833 Thus,  $z_m = z_n$  implies:

$$834 (w_m - w_n)h + (b_m - b_n) = 0. \quad (14)$$

835 Additionally, for  $p_m = p_n$  to be maximal, we require  $p_m \geq p_o$  for all  $o \neq m, n$ , which implies:

$$836 \frac{e^{z_m}}{\sum_{j=1}^M e^{z_j}} \geq \frac{e^{z_o}}{\sum_{j=1}^M e^{z_j}} \implies e^{z_m} \geq e^{z_o} \implies z_m \geq z_o, \quad \forall o \neq m, n. \quad (15)$$

837 Since  $z_m = z_n$ , this becomes:

$$838 z_m = z_n \geq z_o, \quad \forall o \neq m, n. \quad (16)$$

839 In the representation space, this translates to:

$$840 (w_m - w_o)h + (b_m - b_o) \geq 0, \quad (w_n - w_o)h + (b_n - b_o) \geq 0, \quad \forall o \neq m, n. \quad (17)$$

841 For each pair  $m, n \in \mathcal{M}$ ,  $1 \leq m < n \leq M$ , define:

$$842 \mathcal{B}_{mn} = \{h \in \mathbb{R}^{d_h} \mid (w_m - w_n)h + (b_m - b_n) = 0, z_m = z_n \geq z_o \forall o \neq m, n, h = f_r(\mathbf{x}), \mathbf{x} \in \mathcal{D}\}. \quad (18)$$

843 The decision boundary is the union of all such pairwise boundaries:

$$844 \mathcal{B}_M^{(f, \mathcal{D})} = \bigcup_{1 \leq m < n \leq M} \mathcal{B}_{mn}. \quad (19)$$

864 **Part II: Voronoi Cells.**

865 Each  $\mathcal{B}_{mn}$  is a  $(d_h - 1)$ -dimensional hyperplane in  $\mathbb{R}^{d_h}$  defined by  $(w_m - w_n)h + (b_m - b_n) = 0$ ,  
 866 restricted to points where  $z_m = z_n \geq z_o$ . Geometrically, the classification region for class  $i$  is:

867 
$$\mathcal{R}_i = \{h \in \mathbb{R}^{d_h} \mid w_i h + b_i > w_j h + b_j, \forall j \neq i, h = f_r(\mathbf{x}), \mathbf{x} \in \mathcal{D}\}. \quad (20)$$

868 These regions are convex polytopes, as they are defined by the intersection of half-spaces  $(w_i -$   
 869  $w_j)h + (b_i - b_j) > 0$ . The boundaries between  $\mathcal{R}_m$  and  $\mathcal{R}_n$  occur where  $(w_m - w_n)h + (b_m - b_n) = 0$   
 870 and  $z_m = z_n \geq z_o$ , forming  $\mathcal{B}_{mn}$ . The collection  $\{\mathcal{R}_i\}_{i=1}^M$  partitions the representation space, and  
 871 the hyperplanes  $\mathcal{B}_{mn}$  form the boundaries of a Voronoi-like partition, where each  $\mathcal{R}_i$  is a Voronoi  
 872 cell corresponding to class  $i$ .  
 873

874 This completes the proof. □

875 **B.2 PROOF OF THEOREM 3.3**

876 *Proof.* We aim to characterize the decision boundary  $\mathcal{B}_{llm}^{(f, \mathcal{D}' )}$  of an LLM  $f : \mathcal{V}^{N_q} \rightarrow \mathcal{V}^{N_r}$  under an  
 877 input text distribution  $\mathcal{D}' \subseteq \bigcup_{n_q=1}^{N_q} \mathcal{V}^{n_q}$ . The LLM generates a sequence  $\mathbf{y} = [y_1, \dots, y_{N_r}] \in \mathcal{V}^{N_r}$ ,  
 878 where  $\mathcal{V} = \{1, 2, \dots, V\}$  is the vocabulary, conditioned on a prompt  $\mathbf{x} \in \mathcal{D}'$ . The joint probability  
 879 of generating  $\mathbf{y}$  is:

880 
$$P_f(\mathbf{y}|\mathbf{x}) = \prod_{t=1}^{N_r} P_f(y_t|\mathbf{x}, y_1, \dots, y_{t-1}),$$

881 where  $P_f(y_t|\mathbf{x}, y_1, \dots, y_{t-1})$  is the probability of predicting token  $y_t$  at step  $t$ , modeled as a multi-  
 882 class classification over  $\mathcal{V}$ .  
 883

884 Based on  $\mathbf{y}^* = \arg \max_{\mathbf{y} \in \mathcal{V}^{N_r}} P_f(\mathbf{y}|\mathbf{x})$  and Definition 3.1, the decision boundary  $\mathcal{B}_{llm}^{(f, \mathcal{D}' )}$  is the set  
 885 of prompts  $\mathbf{x} \in \mathcal{D}'$  where at least two distinct sequences  $\mathbf{y}_v, \mathbf{y}_w \in \mathcal{V}^{N_r}$  have equal and maximal  
 886 joint probabilities:  
 887

888 
$$P_f(\mathbf{y}_v|\mathbf{x}) = P_f(\mathbf{y}_w|\mathbf{x}) \geq \max_{\mathbf{y}_u \in \mathcal{V}^{N_r} \setminus \{\mathbf{y}_v, \mathbf{y}_w\}} P_f(\mathbf{y}_u|\mathbf{x}).$$

889 For each pair of distinct sequences  $\mathbf{y}_v, \mathbf{y}_w \in \mathcal{V}^{N_r}$ , define:

890 
$$\mathcal{B}_{llm, vw} = \left\{ \mathbf{x} \in \mathcal{D}' \mid P_f(\mathbf{y}_v|\mathbf{x}) = P_f(\mathbf{y}_w|\mathbf{x}) \geq \max_{\mathbf{y}_u \in \mathcal{V}^{N_r} \setminus \{\mathbf{y}_v, \mathbf{y}_w\}} P_f(\mathbf{y}_u|\mathbf{x}) \right\}.$$

891 The decision boundary is the union over all such pairs:

892 
$$\mathcal{B}_{llm}^{(f, \mathcal{D}' )} = \bigcup_{\mathbf{y}_v \neq \mathbf{y}_w \in \mathcal{V}^{N_r}} \mathcal{B}_{llm, vw}.$$

893 To show this, consider the autoregressive process. For a prompt  $\mathbf{x}$ , the probability  $P_f(\mathbf{y}|\mathbf{x})$  depends  
 894 on the token probabilities at each step. Obviously, the predicted sequence  $\mathbf{y}^*$  maximizes  $P_f(\mathbf{y}|\mathbf{x})$ .  
 895 The decision boundary occurs when two sequences  $\mathbf{y}_v$  and  $\mathbf{y}_w$  have equal probabilities, and no other  
 896 sequence has a higher probability. This implies:

897 
$$P_f(\mathbf{y}_v|\mathbf{x}) = \prod_{t=1}^{N_r} P_f(y_{v,t}|\mathbf{x}, y_{v,1}, \dots, y_{v,t-1}) = \prod_{t=1}^{N_r} P_f(y_{w,t}|\mathbf{x}, y_{w,1}, \dots, y_{w,t-1}) = P_f(\mathbf{y}_w|\mathbf{x}),$$

898 and for all  $\mathbf{y}_u \neq \mathbf{y}_v, \mathbf{y}_w$ :

899 
$$P_f(\mathbf{y}_v|\mathbf{x}) \geq P_f(\mathbf{y}_u|\mathbf{x}).$$

900 Since each token prediction is a multi-class classification (as in Theorem 3.2), the boundary for  
 901 a single token  $y_t$  is defined by equal probabilities for the top tokens. For the full sequence, the  
 902 boundary  $\mathcal{B}_{llm, vw}$  corresponds to prompts  $\mathbf{x}$  where the joint probabilities align, which may occur  
 903 when the log-probabilities differ at some steps but sum to the same value. The maximality condition  
 904 ensures that  $\mathbf{y}_v$  and  $\mathbf{y}_w$  are the top sequences.  
 905

906 This completes the proof. □

### B.3 PROOF OF THEOREM 4.3

*Proof.* We aim to prove that the decision boundary  $\mathcal{B}_{llm}^{(f, \mathcal{D}')}$  defined in Theorem 3.3 is equivalent to the 0-isohypse  $\mathcal{D}'_{(0,f)}$  on the decision potential surface  $\mathcal{S}^{(f, \mathcal{D}')}$ , and that the regions separated by this boundary correspond exactly to the Voronoi cells in the token-combined classification definition.

Recall from Theorem 3.3 that the decision boundary is

$$\mathcal{B}_{llm}^{(f, \mathcal{D}')} = \bigcup_{\mathbf{y}_v \neq \mathbf{y}_w \in \mathcal{V}^{N_r}} \mathcal{B}_{llm, vw}, \quad (21)$$

where

$$\mathcal{B}_{llm, vw} = \left\{ \mathbf{x} \in \mathcal{D}' \mid P_f(\mathbf{y}_v | \mathbf{x}) = P_f(\mathbf{y}_w | \mathbf{x}) \geq \max_{\mathbf{y}_u \in \mathcal{V}^{N_r} \setminus \{\mathbf{y}_v, \mathbf{y}_w\}} P_f(\mathbf{y}_u | \mathbf{x}) \right\}. \quad (22)$$

This boundary consists of prompts  $\mathbf{x}$  where at least two distinct sequences  $\mathbf{y}_v$  and  $\mathbf{y}_w$  have equal and maximal joint probabilities, leading to ambiguity in the predicted output sequence.

From Definition 4.1, the decision potential function is

$$\Phi_f^\infty(\mathbf{x}) = (\log P_f(\mathbf{y}_1 | \mathbf{x}) - \log P_f(\mathbf{y}_2 | \mathbf{x}))^2, \quad (23)$$

where  $\mathbf{y}_1, \mathbf{y}_2 \in \mathcal{V}^{N_r}$  are the sequences with the highest and second-highest log-likelihoods, respectively. The 0-isohypse is defined as

$$\mathcal{D}'_{(0,f)} = \{ \mathbf{x} \in \mathcal{D}' \mid \Phi_f^\infty(\mathbf{x}) = 0 \}. \quad (24)$$

By definition,  $\Phi_f^\infty(\mathbf{x}) = 0$  if and only if  $\log P_f(\mathbf{y}_1 | \mathbf{x}) = \log P_f(\mathbf{y}_2 | \mathbf{x})$ , which implies  $P_f(\mathbf{y}_1 | \mathbf{x}) = P_f(\mathbf{y}_2 | \mathbf{x})$ . Since  $\mathbf{y}_1$  and  $\mathbf{y}_2$  are the top two sequences by log-likelihood, this equality ensures that

$$P_f(\mathbf{y}_1 | \mathbf{x}) = P_f(\mathbf{y}_2 | \mathbf{x}) \geq P_f(\mathbf{y}_u | \mathbf{x}), \quad \forall \mathbf{y}_u \neq \mathbf{y}_1, \mathbf{y}_2, \quad (25)$$

satisfying the maximality condition in Theorem 3.3. Thus,  $\mathbf{x} \in \mathcal{D}'_{(0,f)}$  if and only if  $\mathbf{x} \in \mathcal{B}_{llm}^{(f, \mathcal{D}')}$ , establishing the set equivalence

$$\mathcal{B}_{llm}^{(f, \mathcal{D}')} = \mathcal{D}'_{(0,f)}. \quad (26)$$

Geometrically, the regions separated by the 0-isohypse are the connected components of  $\mathcal{D}' \setminus \mathcal{D}'_{(0,f)}$ , where each region corresponds to prompts for which a unique sequence  $\mathbf{y}_i$  has the highest log-likelihood ( $\Phi_f^\infty(\mathbf{x}) > 0$ ). These regions are exactly the Voronoi cells in the sequence-level classification framework of Theorem 3.3, as each cell consists of prompts yielding the same maximal sequence. The 0-isohypse forms the boundaries between these cells, partitioning the prompt space  $\mathcal{D}'$  into regions of unambiguous predictions.

This completes the proof.  $\square$

### B.4 PROOF OF COROLLARY 4.4

*Proof.* We aim to show that for any  $\varepsilon > 0$ , the input space  $\mathcal{D}'$  is partitioned into three disjoint strata based on the value of the decision potential function  $\Phi_f^\infty(\mathbf{x})$ :

$$\mathcal{D}' = \mathcal{D}'_{(>\varepsilon, f)} \sqcup \mathcal{D}'_{(<\varepsilon, f)} \sqcup \mathcal{D}'_{(\varepsilon, f)}, \quad (27)$$

where  $\sqcup$  denotes disjoint union.

From Definition 4.2, the  $\varepsilon$ -isohypse is

$$\mathcal{D}'_{(\varepsilon, f)} = \{ \mathbf{x} \in \mathcal{D}' \mid \Phi_f^\infty(\mathbf{x}) = \varepsilon \}, \quad (28)$$

and the other strata are defined as

$$\mathcal{D}'_{(>\varepsilon, f)} = \{ \mathbf{x} \in \mathcal{D}' \mid \Phi_f^\infty(\mathbf{x}) > \varepsilon \}, \quad \mathcal{D}'_{(<\varepsilon, f)} = \{ \mathbf{x} \in \mathcal{D}' \mid \Phi_f^\infty(\mathbf{x}) < \varepsilon \}. \quad (29)$$

Since  $\Phi_f^\infty : \mathcal{D}' \rightarrow \mathbb{R}_+$  is a continuous function (assuming log-likelihoods are continuous in the prompt space), these sets are disjoint and their union covers  $\mathcal{D}'$ .

972 •  $\varepsilon$ -confident regions: For  $\mathbf{x} \in \mathcal{D}'_{(>\varepsilon,f)}$ ,  $\Phi_f^\infty(\mathbf{x}) > \varepsilon$ , so

$$973 \quad |\log P_f(\mathbf{y}_1|\mathbf{x}) - \log P_f(\mathbf{y}_2|\mathbf{x})| > \sqrt{\varepsilon}. \quad (30)$$

974 Since  $\mathbf{y}_1$  has the highest log-likelihood,  $\log P_f(\mathbf{y}_1|\mathbf{x}) - \log P_f(\mathbf{y}_2|\mathbf{x}) > \sqrt{\varepsilon}$ , meaning the model  
975 predicts  $\mathbf{y}_1$  with at least  $\sqrt{\varepsilon}$  nats (natural units of information) of confidence over the next most  
976 likely sequence  $\mathbf{y}_2$ .  
977

978 •  $\varepsilon$ -uncertain regions: For  $\mathbf{x} \in \mathcal{D}'_{(<\varepsilon,f)}$ ,  $\Phi_f^\infty(\mathbf{x}) < \varepsilon$ , so

$$979 \quad |\log P_f(\mathbf{y}_1|\mathbf{x}) - \log P_f(\mathbf{y}_2|\mathbf{x})| < \sqrt{\varepsilon}. \quad (31)$$

980 Here, the model has low confidence, with a margin less than  $\sqrt{\varepsilon}$  nats between the top two sequences.  
981 As  $\varepsilon \rightarrow 0$ ,  $\Phi_f^\infty(\mathbf{x}) \rightarrow 0$ , so  $\mathcal{D}'_{(<\varepsilon,f)}$  converges to the 0-isohypse  $\mathcal{D}'_{(0,f)}$ , where the margin is zero.

982 •  $\varepsilon$ -isohypse: For  $\mathbf{x} \in \mathcal{D}'_{(\varepsilon,f)}$ ,  $\Phi_f^\infty(\mathbf{x}) = \varepsilon$ , so the confidence margin is exactly  $\sqrt{\varepsilon}$  nats, forming the  
983 contour that separates confident and uncertain regions.

984 The disjointness of the strata follows from the strict inequalities and equality defining them, and  
985 their union covers  $\mathcal{D}'$  since  $\Phi_f^\infty(\mathbf{x}) \geq 0$  for all  $\mathbf{x} \in \mathcal{D}'$ .  
986

987 This completes the proof.  $\square$

## 988 B.5 PROOF OF THEOREM 4.6

989 *Proof.* Let  $\mathbf{x} \in \mathcal{D}'$  be a fixed input, and let  $\mathcal{Y}_K = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K\}$  be a set of  $K$  *i.i.d.* samples  
990 drawn from the language model's output distribution  $P_f(\cdot|\mathbf{x})$ . The decision potential function is:

$$991 \quad \Phi_f^\infty(\mathbf{x}) = (\log P_f(\mathbf{y}_{1*}|\mathbf{x}) - \log P_f(\mathbf{y}_{2*}|\mathbf{x}))^2, \quad (32)$$

992 where  $\mathbf{y}_{1*}$  and  $\mathbf{y}_{2*}$  are the top-2 generated texts with the highest logarithmic likelihoods over the  
993 entire output space  $\mathcal{Y}^{N_r}$ , and

$$994 \quad \Phi_f^K(\mathbf{x}) = (\log P_f(\mathbf{y}_{1*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2*}^K|\mathbf{x}))^2, \quad (33)$$

995 where  $\mathbf{y}_{1*}^K$  and  $\mathbf{y}_{2*}^K$  are the top-2 generated texts within  $\mathcal{Y}_K$ . We aim to bound the error  $|\Phi_f^K(\mathbf{x}) -$   
996  $\Phi_f^\infty(\mathbf{x})|$  with probability at least  $1 - \delta - 2\varepsilon_{\text{tail}}$  for  $\delta \in (0, 1)$ .  
997

998 **Step 1: Preliminary.** Define:

$$999 \quad \begin{aligned} 1000 \quad \Delta_\infty(\mathbf{x}) &= \log P_f(\mathbf{y}_{1*}|\mathbf{x}) - \log P_f(\mathbf{y}_{2*}|\mathbf{x}), \\ 1001 \quad \Delta_K(\mathbf{x}) &= \log P_f(\mathbf{y}_{1*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2*}^K|\mathbf{x}). \end{aligned} \quad (34)$$

1002 Thus,  $\Phi_f^\infty(\mathbf{x}) = (\Delta_\infty(\mathbf{x}))^2$  and  $\Phi_f^K(\mathbf{x}) = (\Delta_K(\mathbf{x}))^2$ . The error can be expressed as:

$$1003 \quad \begin{aligned} 1004 \quad |\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| &= |(\Delta_K(\mathbf{x}))^2 - (\Delta_\infty(\mathbf{x}))^2| \\ 1005 \quad &= |\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| \cdot |\Delta_K(\mathbf{x}) + \Delta_\infty(\mathbf{x})|. \end{aligned} \quad (35)$$

1006 Since  $\mathcal{Y}_K$  is a finite,  $\mathbf{y}_{1*}^K$  and  $\mathbf{y}_{2*}^K$  are the top-2 outputs in  $\mathcal{Y}_K$ , which may not include  $\mathbf{y}_{1*}$  or  $\mathbf{y}_{2*}$ .  
1007 Define:

$$1008 \quad R_K(\mathbf{x}) = \log P_f(\mathbf{y}_{1*}^K|\mathbf{x}) - \min_{\mathbf{y} \in \mathcal{Y}_K} \log P_f(\mathbf{y}|\mathbf{x}), \quad (36)$$

1009 which represents the *diameter* of log-likelihoods in  $\mathcal{Y}_K$ .

1010 **Lemma B.1** ( $\Pr(\mathbf{y}_{1*} \notin \mathcal{Y}_K) \leq \varepsilon_{\text{tail}}$ ). *Define the tail probability  $\varepsilon_{\text{tail}}$  as:  $\varepsilon_{\text{tail}} = (1 - P_f(\mathbf{y}_{1*}^K|\mathbf{x}))^K$ .*  
1011 *Then, we have  $\Pr(\mathbf{y}_{1*} \notin \mathcal{Y}_K) \leq \varepsilon_{\text{tail}}$ ,*

1012 *A short proof of Lemma B.1:* As  $\mathbf{y}_k \in \mathcal{Y}_K$  are *i.i.d.*, we know that with  $K$ -time sampling of  $\mathbf{y} \sim$   
1013  $P_f(\cdot|\mathbf{x})$  the probability that we cannot obtain  $\mathbf{y}_{1*}$  obeys a geometric distribution, i.e.,

$$1014 \quad \Pr(\mathbf{y}_{1*} \notin \mathcal{Y}_k) = (1 - P_f(\mathbf{y}_{1*}|\mathbf{x}))^k. \quad (37)$$

1015 As  $P_f(\mathbf{y}_{1*}|\mathbf{x}) \geq P_f(\mathbf{y}_{1*}^K|\mathbf{x})$ , then we have

$$1016 \quad \Pr(\mathbf{y}_{1*} \notin \mathcal{Y}_K) = (1 - P_f(\mathbf{y}_{1*}|\mathbf{x}))^K \leq (1 - P_f(\mathbf{y}_{1*}^K|\mathbf{x}))^K = \varepsilon_{\text{tail}}, \quad (38)$$

1026 which ends the proof.

1027 Based on Lemma B.1, we know that  $\varepsilon_{\text{tail}}$  bounds the probability that the true top output  $\mathbf{y}_{1^*}$  is not  
1028 included in  $\mathcal{Y}_K$ .

1030 **Step 2: Bounding**  $|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})|$ .

1031 Since  $\Delta_K(\mathbf{x})$  is computed over a random sample, we consider using *concentration inequalities* to  
1032 bound the deviation  $|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})|$ . The log-likelihoods  $\log P_f(\mathbf{y}_k|\mathbf{x})$  for  $\mathbf{y}_k \in \mathcal{Y}_K$  are *i.i.d.*,  
1033 and they are bounded within the diameter  $R_K(\mathbf{x})$ . By *Hoeffding's inequality*, the deviation of the  
1034 sample maximum log-likelihood from its expected maximum is bounded. Specifically, for the top-1  
1035 log-likelihood  $\forall t > 0$ , we have:

$$1037 \Pr(|\log P_f(\mathbf{y}_{1^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{1^*}|\mathbf{x})| > t) \leq 2 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right). \quad (39)$$

1041 Similarly, for the second-highest log-likelihood, a similar bound applies.

1042 Combining these, we have:

$$1044 \begin{aligned} 1045 |\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| &= |(\log P_f(\mathbf{y}_{1^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2^*}^K|\mathbf{x})) - (\log P_f(\mathbf{y}_{1^*}|\mathbf{x}) - \log P_f(\mathbf{y}_{2^*}|\mathbf{x}))| \\ 1046 &= |(\log P_f(\mathbf{y}_{1^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{1^*}|\mathbf{x})) - (\log P_f(\mathbf{y}_{2^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2^*}|\mathbf{x}))|. \end{aligned} \quad (40)$$

1048 Based on the triangle inequality  $|a - b| \leq |a| + |b|$  when  $a, b \in \mathbb{R}$ , we know that

$$1050 \begin{aligned} 1051 |\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| &= |(\log P_f(\mathbf{y}_{1^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{1^*}|\mathbf{x})) - (\log P_f(\mathbf{y}_{2^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2^*}|\mathbf{x}))| \\ 1052 &\leq |\log P_f(\mathbf{y}_{1^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{1^*}|\mathbf{x})| + |\log P_f(\mathbf{y}_{2^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2^*}|\mathbf{x})|. \end{aligned} \quad (41)$$

1054 To bound  $|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})|$ , we aim to find the maximal probability for the event  $|\Delta_K(\mathbf{x}) -$   
1055  $\Delta_\infty(\mathbf{x})| < t'$  with  $t' > 0$ . Without losing generality, we set  $t' = 2t$ , where the objective can be  
1056 reformated as:

$$1058 \begin{aligned} 1059 \Pr(|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| < t') & \\ 1060 &= 1 - \Pr(|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| \geq t'), \end{aligned} \quad (42)$$

1062 where

$$1064 \begin{aligned} 1065 \Pr(|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| \geq t') & \\ 1066 &= \Pr(|\log P_f(\mathbf{y}_{1^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{1^*}|\mathbf{x})| \geq t \text{ or } |\log P_f(\mathbf{y}_{2^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2^*}|\mathbf{x})| \geq t) \\ 1067 &\leq \Pr(|\log P_f(\mathbf{y}_{1^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{1^*}|\mathbf{x})| \geq t) + \Pr(|\log P_f(\mathbf{y}_{2^*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2^*}|\mathbf{x})| \geq t) \\ 1068 &\leq 2 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right) + 2 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right) \\ 1069 &= 4 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right). \end{aligned} \quad (43)$$

1074 So we have

$$1076 \begin{aligned} 1077 \Pr(|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| < t') & \\ 1078 &= 1 - \Pr(|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| \geq t') \\ 1079 &\geq 1 - 4 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right). \end{aligned} \quad (44)$$

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

Suppose we have at least  $1 - \delta$  probability to support this event stands, we have

$$\begin{aligned}
& 1 - 4 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right) \geq 1 - \delta \\
& \Leftrightarrow 4 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right) \leq \delta \\
& \Leftrightarrow \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right) \leq \frac{\delta}{4} \\
& \Leftrightarrow -\frac{2Kt^2}{R_K^2(\mathbf{x})} \leq \log \frac{\delta}{4} \\
& \Leftrightarrow \frac{2Kt^2}{R_K^2(\mathbf{x})} \geq -\log \frac{\delta}{4} \\
& \Leftrightarrow \frac{2Kt^2}{R_K^2(\mathbf{x})} \geq \log \frac{4}{\delta} \\
& \Leftrightarrow t^2 \geq \frac{R_K^2(\mathbf{x})}{2K} \log \frac{4}{\delta} \\
& \Leftrightarrow t \geq |R_K(\mathbf{x})| \sqrt{\frac{\log(4/\delta)}{2K}} \\
& \Leftrightarrow t \geq R_K(\mathbf{x}) \sqrt{\frac{\log(4/\delta)}{2K}}.
\end{aligned} \tag{45}$$

In other words,  $\forall t > 0$  we bound  $\Pr(|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| < t)$  with probability at least  $1 - \delta$  when:

$$t = R_K(\mathbf{x}) \sqrt{\frac{\log(4/\delta)}{2K}}. \tag{46}$$

**Step 3: Bounding  $|\Delta_K(\mathbf{x}) + \Delta_\infty(\mathbf{x})|$ .**

**Assumption B.2** (Bounded Population Gap). There exists a constant  $M > 0$  such that for any  $\mathbf{x}$ , the population top-2 gap satisfies:

$$\Delta_\infty(\mathbf{x}) = \log P_f(\mathbf{y}_{1*}|\mathbf{x}) - \log P_f(\mathbf{y}_{2*}|\mathbf{x}) \leq M \tag{47}$$

Then we assume that

$$M \leq R_K(\mathbf{x}) \tag{48}$$

when  $K \gg 1$ .

This assumption is reasonable as most practical language models do not have extremely large differences between top-2 probabilities, and the probability differences between top-2 would be much smaller than the range of between the top-1 and the sample with the minimal probability in the sampling set. Now we can obtain that:

$$|\Delta_K(\mathbf{x}) + \Delta_\infty(\mathbf{x})| \leq |\Delta_K(\mathbf{x})| + |\Delta_\infty(\mathbf{x})| \leq 2 \cdot R_K(\mathbf{x}). \tag{49}$$

**Step 4: Final bound.**

Define the events

$$A = \{\mathbf{y}_{1*} \in \mathcal{Y}_K \text{ and } \mathbf{y}_{2*} \in \mathcal{Y}_K\}, \quad B = \{\mathbf{y}_{1*} \notin \mathcal{Y}_K \text{ or } \mathbf{y}_{2*} \notin \mathcal{Y}_K\}. \tag{50}$$

Lemma B.1 and a union bound give

$$\Pr(B) \leq 2\varepsilon_{\text{tail}}. \tag{51}$$

- On event  $A$  we have  $\mathbf{y}_{1*}^K = \mathbf{y}_{1*}$  and  $\mathbf{y}_{2*}^K = \mathbf{y}_{2*}$ , hence

$$\Phi_f^K(\mathbf{x}) = \Phi_f^\infty(\mathbf{x}) \implies |\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| = 0. \tag{52}$$

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

- On event  $B$  we use the worst-case gap

$$\begin{aligned}
& |\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \\
& \leq |\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| \cdot |\Delta_K(\mathbf{x}) + \Delta_\infty(\mathbf{x})| \\
& \leq R_K(\mathbf{x}) \sqrt{\frac{\log(4/\delta)}{2K}} \cdot 2R_K(\mathbf{x}) \\
& = 2R_K^2(\mathbf{x}) \sqrt{\frac{\log(4/\delta)}{2K}}.
\end{aligned} \tag{53}$$

This completes the proof. □

## B.6 PROOF OF THEOREM 4.7

*Proof.* We aim to bound the expected error  $\mathbb{E}[|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})|]$  for a fixed input  $\mathbf{x} \in \mathcal{D}'$  and a set  $\mathcal{Y}_K = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K\}$  of  $K$  *i.i.d.* samples drawn from the language model's output distribution  $P_f(\cdot|\mathbf{x})$ . Recall that:

$$\Phi_f^\infty(\mathbf{x}) = (\log P_f(\mathbf{y}_{1*}|\mathbf{x}) - \log P_f(\mathbf{y}_{2*}|\mathbf{x}))^2, \quad \Phi_f^K(\mathbf{x}) = (\log P_f(\mathbf{y}_{1*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2*}^K|\mathbf{x}))^2, \tag{54}$$

where  $\mathbf{y}_{1*}, \mathbf{y}_{2*}$  are the top-2 outputs over the entire output space  $\mathcal{Y}^{N_r}$ , and  $\mathbf{y}_{1*}^K, \mathbf{y}_{2*}^K$  are the top-2 outputs in  $\mathcal{Y}_K$ . Define:

$$\Delta_\infty(\mathbf{x}) = \log P_f(\mathbf{y}_{1*}|\mathbf{x}) - \log P_f(\mathbf{y}_{2*}|\mathbf{x}), \quad \Delta_K(\mathbf{x}) = \log P_f(\mathbf{y}_{1*}^K|\mathbf{x}) - \log P_f(\mathbf{y}_{2*}^K|\mathbf{x}), \tag{55}$$

so that  $\Phi_f^\infty(\mathbf{x}) = (\Delta_\infty(\mathbf{x}))^2$ ,  $\Phi_f^K(\mathbf{x}) = (\Delta_K(\mathbf{x}))^2$ , and the error is:

$$|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| = |(\Delta_K(\mathbf{x}))^2 - (\Delta_\infty(\mathbf{x}))^2| = |\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| \cdot |\Delta_K(\mathbf{x}) + \Delta_\infty(\mathbf{x})|. \tag{56}$$

By Assumption B.2,  $|\Delta_\infty(\mathbf{x})| \leq R_K(\mathbf{x})$ , where  $R_K(\mathbf{x}) = \log P_f(\mathbf{y}_{1*}^K|\mathbf{x}) - \min_{\mathbf{y} \in \mathcal{Y}_K} \log P_f(\mathbf{y}|\mathbf{x})$  is the log-likelihood diameter of  $\mathcal{Y}_K$ . Also,  $|\Delta_K(\mathbf{x})| \leq R_K(\mathbf{x})$ , so:

$$|\Delta_K(\mathbf{x}) + \Delta_\infty(\mathbf{x})| \leq |\Delta_K(\mathbf{x})| + |\Delta_\infty(\mathbf{x})| \leq 2R_K(\mathbf{x}). \tag{57}$$

Thus, the error is bounded by:

$$|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \leq |\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| \cdot 2R_K(\mathbf{x}). \tag{58}$$

We compute the expectation:

$$\mathbb{E}[|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})|] \leq 2R_K(\mathbf{x}) \cdot \mathbb{E}[|\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})|]. \tag{59}$$

Define  $Z = |\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})|$ . From the proof of Theorem 4.6 (Equation 43), Hoeffding's inequality gives:

$$\Pr(Z \geq t) \leq 4 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right). \tag{60}$$

The expectation of  $Z$  is:

$$\mathbb{E}[Z] = \int_0^\infty \Pr(Z \geq t) dt \leq \int_0^\infty 4 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right) dt. \tag{61}$$

Substitute  $u = \frac{2Kt^2}{R_K^2(\mathbf{x})}$ , so  $t = R_K(\mathbf{x})\sqrt{\frac{u}{2K}}$ ,  $dt = \frac{R_K(\mathbf{x})}{2\sqrt{2K}\sqrt{u}} du$ . Then:

$$\mathbb{E}[Z] \leq \int_0^\infty 4e^{-u} \cdot \frac{R_K(\mathbf{x})}{2\sqrt{2K}\sqrt{u}} du = \frac{2R_K(\mathbf{x})}{\sqrt{2K}} \int_0^\infty \frac{e^{-u}}{\sqrt{u}} du. \tag{62}$$

Since  $\int_0^\infty \frac{e^{-u}}{\sqrt{u}} du = \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$ , we have:

$$\mathbb{E}[Z] \leq \frac{2R_K(\mathbf{x})}{\sqrt{2K}} \cdot \sqrt{\pi} = R_K(\mathbf{x})\sqrt{\frac{2\pi}{K}}. \tag{63}$$

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

Thus:

$$\mathbb{E}[|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})|] \leq 2R_K(\mathbf{x}) \cdot R_K(\mathbf{x}) \sqrt{\frac{2\pi}{K}} = 2R_K^2(\mathbf{x}) \sqrt{\frac{2\pi}{K}}. \quad (64)$$

To account for event  $B = \{\mathbf{y}_{1*} \notin \mathcal{Y}_K \text{ or } \mathbf{y}_{2*} \notin \mathcal{Y}_K\}$  with  $\Pr(B) \leq 2\varepsilon_{\text{tail}}$  (from Lemma B.1 and union bound), we note that on event  $A = \{\mathbf{y}_{1*} \in \mathcal{Y}_K \text{ and } \mathbf{y}_{2*} \in \mathcal{Y}_K\}$ , the error is zero. Thus, we add a conservative term for event  $B$ , where the error is at most  $2R_K^2(\mathbf{x})$  (since  $|\Delta_K|, |\Delta_\infty| \leq R_K(\mathbf{x})$ ), so  $|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \leq 2R_K^2(\mathbf{x})$ :

$$\mathbb{E}[|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \cdot \mathbb{1}_B] \leq 2R_K^2(\mathbf{x}) \cdot \Pr(B) \leq 2R_K^2(\mathbf{x}) \cdot 2\varepsilon_{\text{tail}} = 4R_K^2(\mathbf{x})\varepsilon_{\text{tail}}, \quad (65)$$

where  $\mathbb{1}_B$  is the indicator function which is 1 only when event  $B$  occurs.

Combining both terms, the expected error is:

$$\mathbb{E}[|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})|] \leq 2R_K^2(\mathbf{x}) \sqrt{\frac{2\pi}{K}} + 4R_K^2(\mathbf{x})\varepsilon_{\text{tail}}, \quad (66)$$

where  $\varepsilon_{\text{tail}} = (1 - P_f(\mathbf{y}_{1*}^K | \mathbf{x}))^K$ . This completes the proof.  $\square$

## B.7 PROOF OF COROLLARY 4.8

*Proof.* We aim to bound the tail probability  $\Pr(|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \geq \lambda)$  for  $\lambda > 0$ . Using the same notation as in Theorem 4.7, we have:

$$|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| = |\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})| \cdot |\Delta_K(\mathbf{x}) + \Delta_\infty(\mathbf{x})|. \quad (67)$$

Since  $|\Delta_K(\mathbf{x}) + \Delta_\infty(\mathbf{x})| \leq 2R_K(\mathbf{x})$ , let  $Z = |\Delta_K(\mathbf{x}) - \Delta_\infty(\mathbf{x})|$ , so:

$$|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \leq Z \cdot 2R_K(\mathbf{x}). \quad (68)$$

Thus:

$$\Pr(|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \geq \lambda) = \Pr(Z \cdot 2R_K(\mathbf{x}) \geq \lambda) = \Pr\left(Z \geq \frac{\lambda}{2R_K(\mathbf{x})}\right). \quad (69)$$

From the proof of Theorem 4.6 (Equation 43), Hoeffding's inequality gives:

$$\Pr(Z \geq t) \leq 4 \exp\left(-\frac{2Kt^2}{R_K^2(\mathbf{x})}\right). \quad (70)$$

Set  $t = \frac{\lambda}{2R_K(\mathbf{x})}$ :

$$\Pr\left(Z \geq \frac{\lambda}{2R_K(\mathbf{x})}\right) \leq 4 \exp\left(-\frac{2K \cdot \left(\frac{\lambda}{2R_K(\mathbf{x})}\right)^2}{R_K^2(\mathbf{x})}\right) = 4 \exp\left(-\frac{K\lambda^2}{2R_K^4(\mathbf{x})}\right). \quad (71)$$

Define events  $A = \{\mathbf{y}_{1*} \in \mathcal{Y}_K \text{ and } \mathbf{y}_{2*} \in \mathcal{Y}_K\}$  and  $B = \{\mathbf{y}_{1*} \notin \mathcal{Y}_K \text{ or } \mathbf{y}_{2*} \notin \mathcal{Y}_K\}$ . On event  $A$ , the error is zero, so it does not contribute to the tail probability. On event  $B$ , with  $\Pr(B) \leq 2\varepsilon_{\text{tail}}$  (from Lemma B.1 and union bound), the tail probability is bounded by:

$$\Pr(|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \geq \lambda) \leq \Pr\left(\left\{Z \geq \frac{\lambda}{2R_K(\mathbf{x})}\right\} \cap B\right) + \Pr(A). \quad (72)$$

Since  $\Pr(A) \geq 1 - 2\varepsilon_{\text{tail}}$  and the error is zero on  $A$ , we focus on event  $B$ :

$$\Pr\left(\left\{Z \geq \frac{\lambda}{2R_K(\mathbf{x})}\right\} \cap B\right) \leq \Pr\left(Z \geq \frac{\lambda}{2R_K(\mathbf{x})}\right) + \Pr(B) \leq 4 \exp\left(-\frac{K\lambda^2}{2R_K^4(\mathbf{x})}\right) + 2\varepsilon_{\text{tail}}. \quad (73)$$

Thus, the tail probability is:

$$\Pr(|\Phi_f^K(\mathbf{x}) - \Phi_f^\infty(\mathbf{x})| \geq \lambda) \leq 4 \exp\left(-\frac{K\lambda^2}{2R_K^4(\mathbf{x})}\right) + 2\varepsilon_{\text{tail}}, \quad (74)$$

where  $\varepsilon_{\text{tail}} = (1 - P_f(\mathbf{y}_{1*}^K | \mathbf{x}))^K$ . This completes the proof.  $\square$

---

## 1242 C SUPPLEMENTAL EXPERIMENTS

### 1243 1244 C.1 IMPLIPLICATIONS OF $K$ -DPS

1245  
1246 In this subsection, we provide several preliminary, proof-of-concept applications of the  $K$ -DPS  
1247 algorithm.

#### 1248 1249 C.1.1 ALIGNMENT

1250  
1251 **Settings.** We select three groups of LLMs in our experiments with each group containing two  
1252 models before and after the human alignment. We respectively visualize the linear interpolation  
1253 visualization, nearest interpolation visualization, heatmap, and three dimensional visualization for  
1254 each model, as shown in each row. The color indicates the value of corresponding positions, where  
1255 a darker color indicates more smaller  $K$ -DPS score. We use  $K = 2, 500$  in our experiments, and  
1256 the input queries come from AdvBenchk (Zou et al., 2023).

1257  
1258 **Results.** As shown in Figure 6, the decision boundary of aligned models becomes dramatically  
1259 smoother and flatter compared to their pre-alignment counterparts when evaluated on adversarial  
1260 prompts from AdvBench. This striking smoothing effect indicates that alignment substantially re-  
1261 duces regions of high confidence in harmful outputs (i.e., sharp peaks with high  $K$ -DPS scores).  
1262 Instead, it creates broad, low- $K$ -DPS basins that strongly favor refusal. This geometric transfor-  
1263 mation directly explains both: *i) Why jailbreaks succeed on unaligned models:* they target narrow,  
1264 high-confidence “vulnerability spikes” that remain in the pre-alignment landscape; *ii) Why align-  
1265 ment mitigates most jailbreaks:* it eliminates these spikes entirely, making harmful responses prob-  
1266 abilistically unlikely across vast regions of prompt space.

1267  
1268 Moreover, by applying the same  $K$ -DPS visualization to intermediate checkpoints throughout the  
1269 alignment process, we can even track precisely how these dangerous peaks progressively flatten and  
1270 how decision boundaries move, offering the first fine-grained geometric view of how safety training  
1271 reshapes the model’s decision manifold step by step.

1272  
1273 This level of mechanistic explanation and dynamic visualization was previously infeasible with prior  
1274 probing or interpretability techniques, but becomes straightforward and highly revealing through our  
1275 DPS-based decision boundary construction.

#### 1276 1277 C.1.2 UNLEARNING

1278  
1279 **Settings.**  $K$ -DPS also enables fine-grained analysis of machine unlearning algorithms, where prior  
1280 interpretability tools offer almost no intuitive explanations.

1281  
1282 We apply our  $K$ -DPS to two representative unlearning methods: Gradient Ascent (GA) (Yao et al.,  
1283 2024) and Negative Preference Optimization (NPO) (Zhang et al., 2024). We use the standard Harry  
1284 Potter book as the forget (unlearning) corpus and a Wikipedia subset<sup>2</sup> as the retain set. During  
1285 unlearning, we continue training on the retain set using standard gradient descent (GDR) or the KL  
1286 divergence (KLR) from the original model. We set  $K = 2,000$ .

1287  
1288 **Results.** As shown in Figure 7, our  $K$ -DPS visualizations provide a far clearer picture of the side  
1289 effects of unlearning than previously possible (Liu et al., 2025a; Geng et al., 2025). Prior work could  
1290 only report that unlearning without proper retention training degrades overall performance, yet was  
1291 unable to show what form this degradation takes in the model’s internal decision process. With  $K$ -  
1292 DPS, we reveal that naïve unlearning methods (e.g., pure Gradient Ascent) can trigger catastrophic  
1293 collapse of the entire decision manifold: large portions of the prompt space that were previously  
1294 smooth become extremely jagged and fragmented, with erratic high- and low- $K$ -DPS spikes appear-  
1295 ing in regions unrelated to the forget corpus.

1296  
1297 In contrast, when retention training is included (e.g., GA-GDR and GA-KLR), the damage is sub-  
1298 stantially mitigated. Among these, GA-KLR (KL-regularized retention) preserves a decision sur-  
1299 face that is visibly the closest to the original model, albeit still perceptibly distorted. This aligns  
1300 with quantitative results in the unlearning literature showing that KL regularization best preserves  
1301 general capabilities.

---

1302  
1303 <sup>2</sup><https://huggingface.co/datasets/rag-datasets/rag-mini-wikipedia>

1296 The above observations demonstrate that K-DPS not only confirms known phenomena at a qual-  
1297 itative level but, for the first time, makes the geometric nature of “unlearning damage” directly  
1298 observable and comparable across methods.

1299 **Others.** Beyond the two core insights above, we believe our decision-boundary framework can nat-  
1300 urally extend to a wide range of important LLM phenomena that have so far resisted precise mech-  
1301 anistic analysis, including: *i*) the precise locations and shapes of jailbreak vulnerabilities in prompt  
1302 space, *ii*) the emergence and structure of memorization regions, *iii*) the geometry of hallucination-  
1303 prone areas versus high-fidelity regions, *iv*) systematic changes in the boundary during continual  
1304 learning or catastrophic forgetting, and so on.

## 1306 C.2 DISCUSSIONS

### 1308 C.2.1 WILL THE TWO HIGHEST-PROBABILITY SEQUENCES OFTEN BE NEARLY IDENTICAL?

1309 In this subsection, we dive into the analysis of the sampled responses under  $K$ -DPS, where a poten-  
1310 tial issue is that the top-2 completions might be extremely similar under some situations, sometimes  
1311 even differing by only a single token or a minor lexical variation, which may question the effective-  
1312 ness of the proposed method.

1313 In this subsection, we examine the sampled completions under  $K$ -DPS and address a potential con-  
1314 cern: in certain cases, the top-2 generations can be nearly identical, differing by only a single token,  
1315 minor lexical variations, or superficial formatting. Such high surface-level similarity might raise  
1316 questions about whether K-DPS is sufficiently sensitive to meaningful semantic discrepancies.

1317 Suggested by the visualizations (Figures 10, 6, 7, 11, and 4), the fraction of such near-identical cases  
1318 is actually not significant in practice, especially in the regions that matter most for decision-boundary  
1319 analysis. As clearly demonstrated in these figures, meaningful and sharp decision boundaries consis-  
1320 tently appear. Degenerate boundaries (which would occur if top-1 and top-2 were almost always  
1321 identical) are rarely observed across the vast majority of inputs we study, indicating that this is not  
1322 a common phenomena in the sampling.

1323 To further investigate this problem, we conduct some quantitative analysis, which confirms substan-  
1324 tial lexical and semantic divergence in K-DPS construction. Specifically, for all candidate sequences  
1325 used to construct Figure 10, we measured the normalized Levenshtein edit distance between top-1  
1326 and top-2 completions.

K-DPS Score Range	Avg. Normalized Edit Distance
< 0.1	0.15
< 0.5	0.20
< 1.0	0.23
< 5.0	0.30
< 10.0	0.32

1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335 Table 1: Average normalized edit distance between top-1 and top-2 completions as a function of  
1336  $K$ -DPS confidence score range.

1337 As shown in Table 1, even in the highest-confidence regime ( $K$ -DPS values  $\leq 0.1$ ), the top-2 se-  
1338 quences differ by 15% of tokens on average; near decision boundaries (higher  $K$ -DPS), divergence  
1339 reaches 30–32%, which is far from trivial variants and typically meaningful.

1340 Regarding the situation that the edit distance between top-2 sequences are small, we can address it  
1341 with some simple filtering strategies, such as slightly increasing the sampling temperature during  
1342 candidate generation or directly discarding pairs with near-zero edit distance.

### 1345 C.2.2 $K$ -DPS VERSUS MODEL UNCERTAINTY

1346 We notice that the construction of explicit decision boundaries in the representation space might  
1347 exhibit connections with several core research areas in LLMs, particularly confidence estimation  
1348 and uncertainty quantification (UQ) (Geng et al., 2024; Huang et al., 2025; Liu et al., 2025b; Xia  
1349 et al., 2025b; Shorinwa et al., 2025; Lin et al., 2024). These uncertainty quantification approaches

1350 typically include verbalized confidence expressed in natural language (Kadavath et al., 2022), token-  
1351 level entropy of the output distribution (Kuhn et al., 2023), and semantic entropy computed over se-  
1352 mantically equivalent clusters of multiple generations (Kuhn et al., 2023; Farquhar et al., 2024), with  
1353 the latter achieving state-of-the-art performance in hallucination detection and selective generation  
1354 tasks.

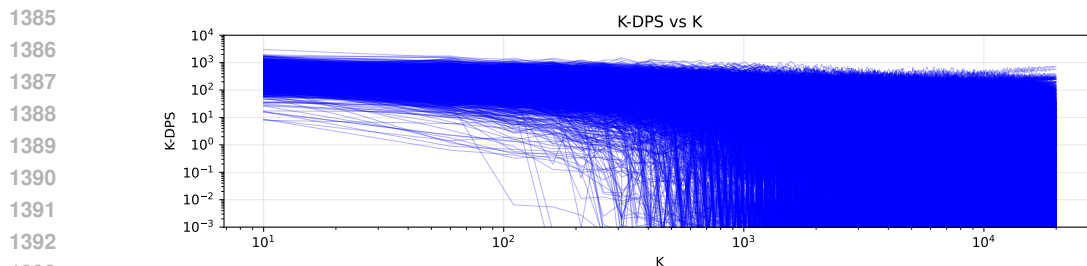
1355 While these methods also measure the certainty and the confidence of the model decision, our K-  
1356 DPS decision-boundary construction differs from them in several fundamental aspects:

1357 First, classical uncertainty quantification techniques Kadavath et al. (2022); Kuhn et al. (2023);  
1358 Farquhar et al. (2024) are essentially heuristic or sampling-based scores lacking formal theoretical  
1359 guarantees, whereas K-DPS provides provably conservative classification boundaries with explicit  
1360 error bounds. It achieves a precise and meaningful approximation of the decision boundary. Second,  
1361 existing UQ methods operate at the instance level and treat each generation independently, while  
1362 K-DPS explicitly builds and reasons over distribution-level decision boundaries, enabling global  
1363 geometric understanding of the model’s reliable support. For the usage, conventional approaches  
1364 remain largely oblivious to the location of samples relative to the empirical data manifold, whereas  
1365 K-DPS deliberately identifies and penalizes anomalous boundary samples that fall near or outside  
1366 the observed support of each semantic class. These distinctions shift the paradigm from post-hoc  
1367 uncertainty scoring to principled, boundary-aware certification of LLM generations.

1368 Nevertheless, we acknowledge that K-DPS and traditional uncertainty quantification methods indeed  
1369 share some core insights. Both paradigms ultimately aim to identify when an LLM’s output is unre-  
1370 liable, whether due to hallucination, out-of-distribution inputs, adversarial attacks, or memorization-  
1371 based spurious responses. Technically, they all ground their analysis in the same internal represen-  
1372 tations of the model: prior UQ approaches directly use raw logits, token probabilities, or hidden  
1373 states to compute verbalized confidence or entropy measures, whereas K-DPS leverages the DPF  
1374 as the theoretical indicator to perform boundary construction. Consequently, the decision boundary  
1375 learned by  $K$ -DPS can be interpreted as a geometrically principled extension of uncertainty sig-  
1376 nals: samples assigned high semantic entropy or low verbalized confidence often naturally fall into  
1377 low-density or boundary regions detected by  $K$ -DPS, providing a unified explanatory framework  
1378 for why existing UQ methods succeed or fail on specific examples. In practice, the two families of  
1379 approaches are highly complementary: uncertainty scores can serve as lightweight pre-filters, while  
1380 K-DPS offers stricter, certifiable analysis for LLM inference.

### 1381 C.3 EMPIRICAL ERROR ANALYSIS UNDER LOG SCALES

1382 We present log-log reconstructions of Figure 1 and Figure 2 in Figure 5 and Figure 8, respectively.



1394 Figure 5: Effect of sampling size  $K$  on the values of decision potential function, with each blue  
1395 point representing the  $K$ -DPS value for a single input sample. Each blue line represents a trend of  
1396  $K$ -DPS for one input sample.  
1397

### 1398 C.4 SUPPLEMENTAL VISUALIZATION

1399  
1400  
1401  
1402  
1403

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

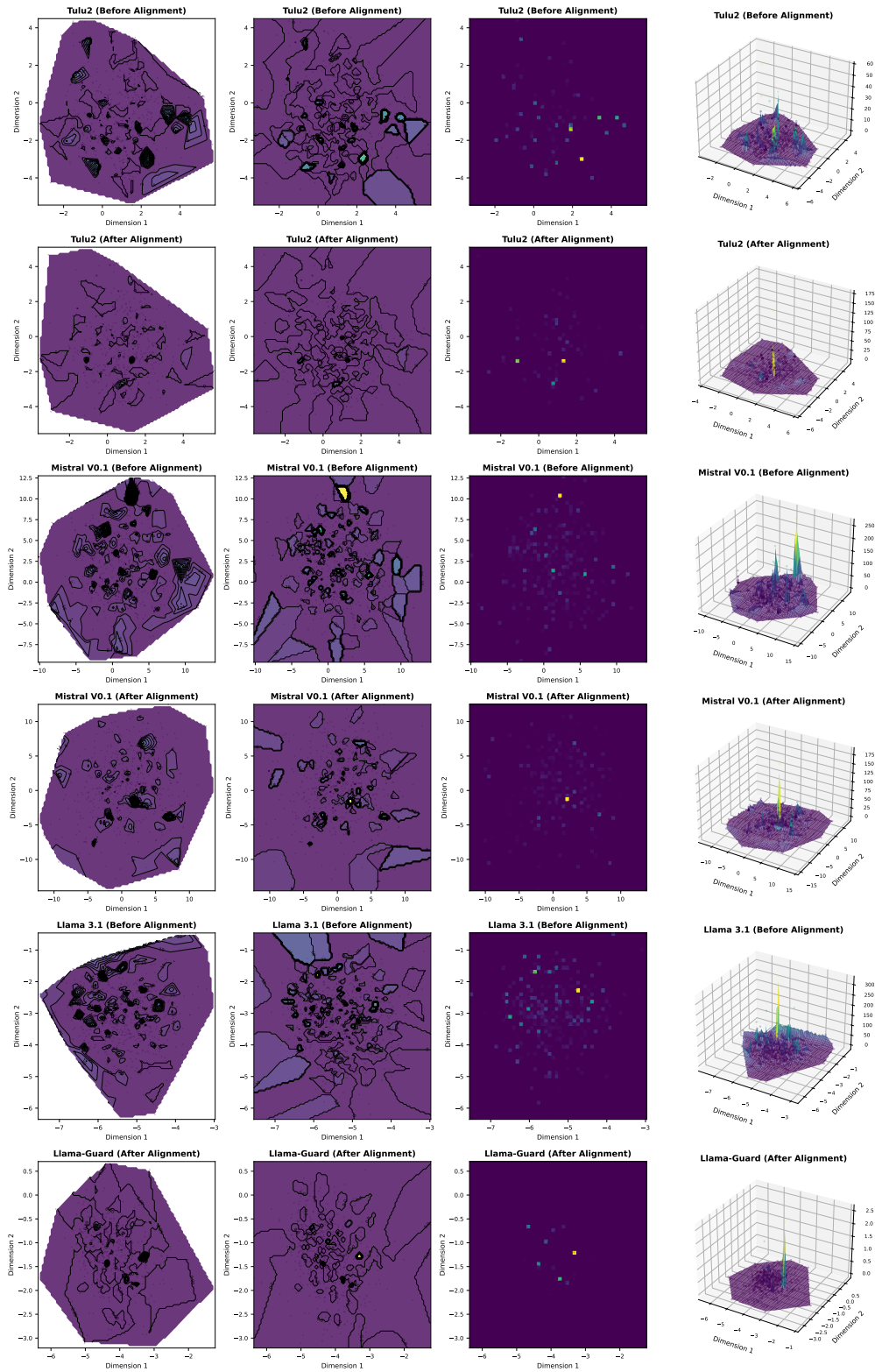


Figure 6: Visualization of decision boundaries constructed by K-DPS in the semantic embedding space, before and after preference alignment. Each row (top to bottom) presents results for Tulu-2-7B (SFT only), Tulu-2-7B-DPO, Mistral-7B-v0.1, Zephyr-7B, Meta-Llama-3.1-8B, and Llama-Guard-3-8B. Columns (left to right) display: 2D interpolation using linear and nearest-neighbor methods, heatmap, and 3D rendering of the decision boundary.

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

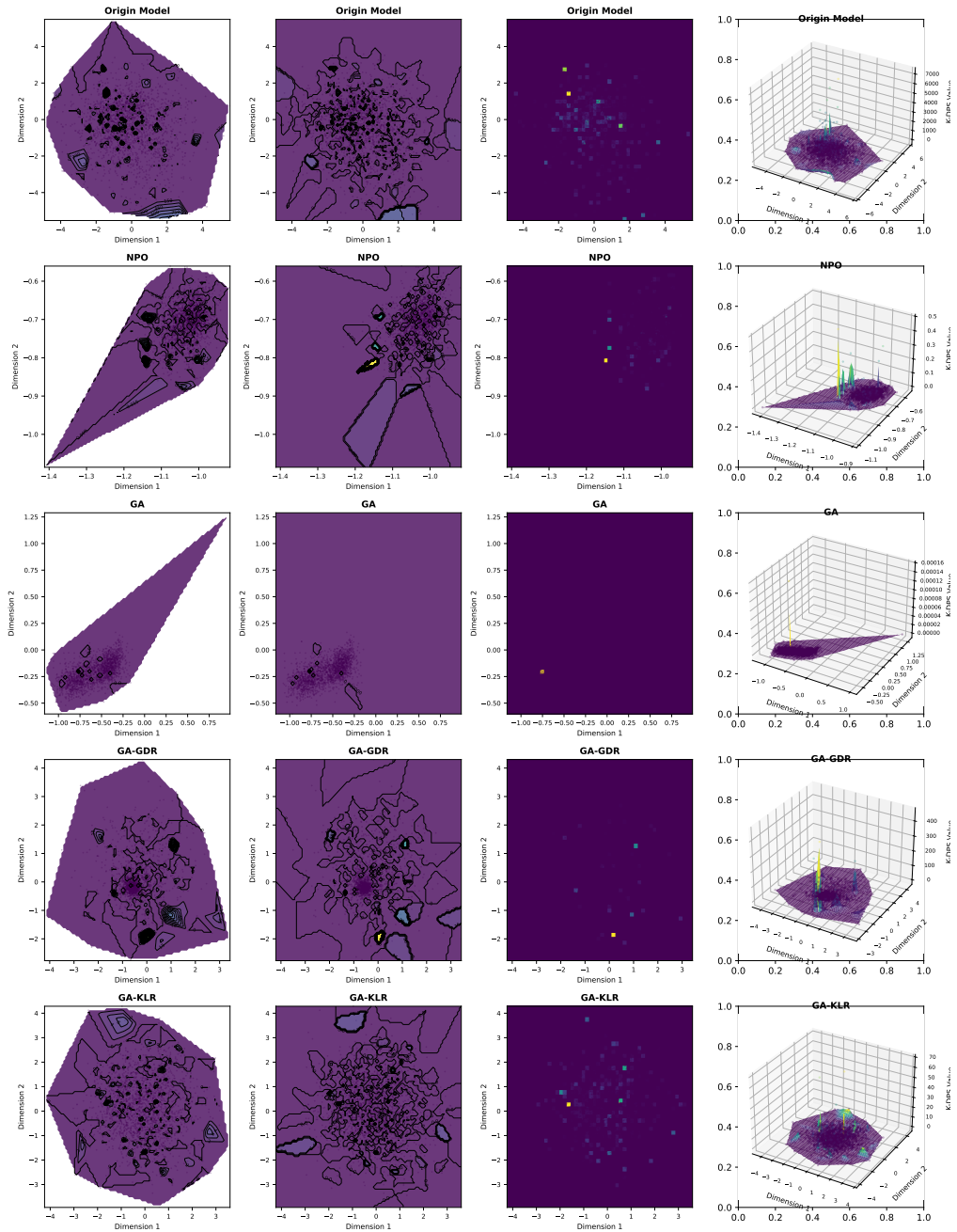


Figure 7: Visualization of decision boundaries constructed by K-DPS on machine unlearning models. Experimental settings and layout follow Figure 6; rows from top to bottom correspond to the original model and its unlearned variants using different unlearning algorithms.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

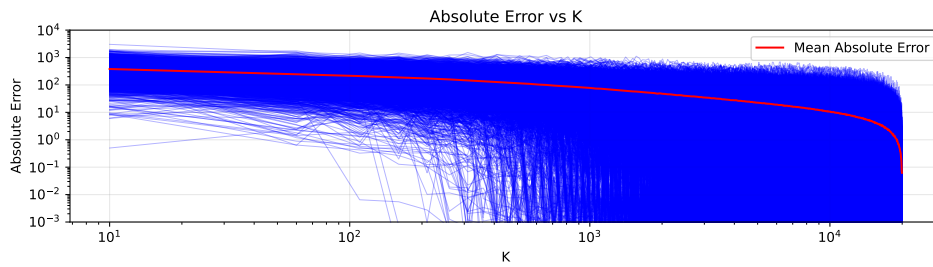


Figure 8: Effect of sampling size  $K$  on the absolute error between the reference  $K$ -DPS (computed with  $K = 20,000$ ) and  $K$ -DPS values for varying  $K$ . Each blue line represents a trend of absolute error across input samples.

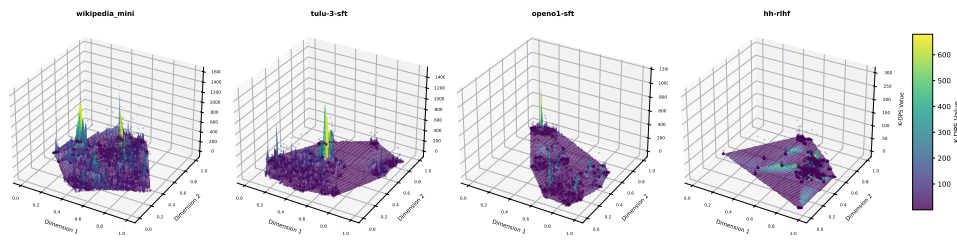


Figure 9: Three-dimensional visualization of the  $K$ -DPS ( $K = 2,500$ ) for Llama-3.2-1B on four datasets.

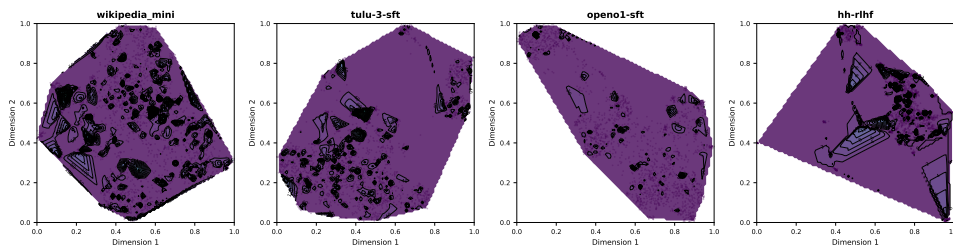


Figure 10: Contour visualization of 2,500-grained decision potential surface for Llama3.2-1B on four datasets with linear interpolation.

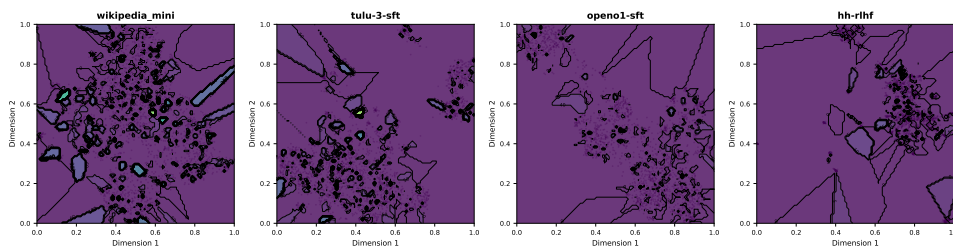


Figure 11: Contour visualization of 2,500-grained decision potential surface of Llama3.2-1B on four datasets with nearest interpolation.

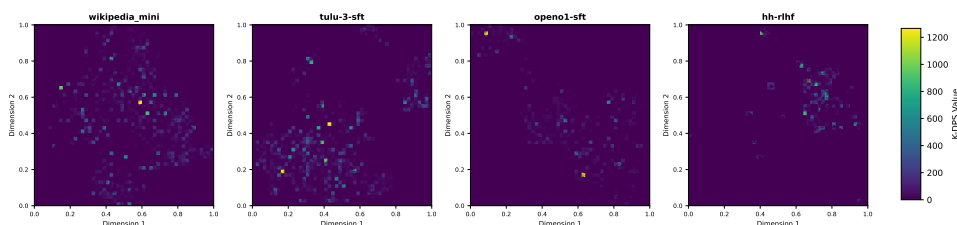


Figure 12: Heatmap visualization of the decision potential surface on Four datasets.