Probing the Decision Boundaries of In-context Learning in Large Language Models

Siyan Zhao, Tung Nguyen, Aditya Grover

Department of Computer Science University of California Los Angeles {siyanz,tungnd,adityag}@cs.ucla.edu

Abstract

In-context learning is a key paradigm in large language models (LLMs) that enables them to generalize to new tasks and domains by simply prompting these models with a few exemplars without explicit parameter updates. Many attempts have been made to understand in-context learning in LLMs as a function of model scale, pretraining data, and other factors. In this work, we propose a new mechanism to probe and understand in-context learning from the lens of decision boundaries for in-context binary classification. Decision boundaries are straightforward to visualize and provide important information about the qualitative behavior of the inductive biases of standard classifiers. To our surprise, we find that the decision boundaries learned by current LLMs in simple binary classification tasks are often irregular and non-smooth, regardless of linear separability in the underlying task. This paper investigates the factors influencing these decision boundaries and explores methods to enhance their generalizability. We assess various approaches, including training-free and fine-tuning methods for LLMs, the impact of model architecture, and the effectiveness of active prompting techniques for smoothing decision boundaries in a data-efficient manner. Our findings provide a deeper understanding of in-context learning dynamics and offer practical improvements for enhancing robustness and generalizability of in-context learning. ¹

1 Introduction

Recent language models, such as GPT-3+ [Brown et al., 2020, Achiam et al., 2023], have demonstrated the ability to scale performance with increased training dataset size and model capacity through the simple objective of next token prediction [Kaplan et al., 2020]. A key emergent behavior of these transformer-based models is in-context learning, which allows the model to learn tasks by conditioning on a sequence of demonstrations without explicit training [Wei et al., 2022]. This unique capability allows LLMs to adapt seamlessly to new tasks, often achieving superior performance in few-shot settings [Brown et al., 2020]. Despite significant successes, the underlying mechanisms of how in-context learning works remain partially understood.

Recent attempts to understand in-context learning have focused on various aspects. From a theoretical standpoint, studies by Von Oswald et al. [2023] and Dai et al. [2023] have linked the mechanisms of in-context learning to gradient descent, suggesting that transformers can emulate optimization processes. On the practical side, research has investigated the impact of different factors on in-context learning. Works by Min et al. [2022b] and Shi et al. [2023] reveal that accurate demonstrations are not essential for effective in-context learning. On the other hand, factors such as the prompt structure and model size [Wei et al., 2023, Webson and Pavlick, 2022], or the order of in-context examples [Chen et al., 2024] greatly affect outcomes. More recently, with the development of LLMs supporting longer

38th Conference on Neural Information Processing Systems (NeurIPS 2024).

Our code is released at https://github.com/siyan-zhao/ICL_decision_boundary.

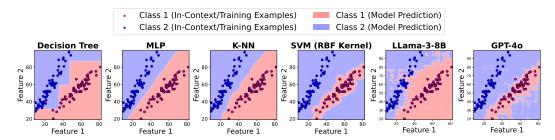


Figure 1: Decision boundaries of LLMs and traditional machine learning models on a linearly separable binary classification task. The background colors represent the model's predictions, while the points represent the in-context or training examples. LLMs exhibit non-smooth decision boundaries compared to the classical models. See Appendix E for model hyperparameters.

context lengths up to 10M [Reid et al., 2024], studies have shown that in-context learning performance improves with significant number of demonstrations [Agarwal et al., 2024, Bertsch et al., 2024], where the performance can be comparable to fine-tuning on the same amount of demonstrations. Additionally, works by Garg et al. [2022], Nguyen and Grover [2022] have demonstrated that small transformers trained from scratch can learn unseen function classes in-context from examples.

In contrast to existing approaches, our study introduces a fresh perspective by viewing in-context learning in large language models (LLMs) as a unique machine learning algorithm. This conceptual framework enables us to leverage a classical tool from machine learning – analyzing decision boundaries in binary classification tasks. By visualizing these decision boundaries, both in linear and non-linear contexts, we gain invaluable insights into the performance and behavior of in-context learning. This method allows us to probe the inductive biases and generalization capabilities of LLMs and offers a unique assessment of the robustness of their in-context learning performance. Consequently, this approach provides a comprehensive means to qualitatively analyze the underlying mechanisms that govern in-context learning and suggest ways to improve its performance in LLMs.

To our surprise, we found that the recent LLMs struggle to provide smooth decision boundaries in all the classification tasks we considered, regardless of the model size, the number and ordering of in-context examples, and semantics of the label format. This issue persists even for simple binary linear classification tasks, where classical methods such as SVM can easily achieve smooth boundaries with fewer examples as shown in Figure 1. This observation raises questions about the factors that influence the decision boundaries of LLMs. To explore this, we experimented with a series of open-source LLMs including Llama2-7b, Llama2-13b, Llama3-8b [Touvron et al., 2023], Mistral-7b [Jiang et al., 2023], pruned Llama2-1.3b [Xia et al., 2023], as well as state-of-the-art closed-source LLMs GPT-4o and GPT-3-Turbo [Brown et al., 2020]. We then explore methods to smooth the decision boundary, including fine-tuning and adaptive prompting strategies. Our work provides valuable practical insights for understanding and improving in-context learning in LLMs through a new perspective. Our contributions can be summarized as follows:

- We introduce a novel mechanism to probe and understand in-context learning in LLMs by visualizing and analyzing the decision boundaries on classification tasks.
- We demonstrate that state-of-the-art LLMs exhibit non-smooth, irregular decision boundaries even on simple linearly separable tasks, unlike classical ML models.
- We study the influence of various factors impacting decision boundary smoothness, including model size, pretraining data and objectives, number of in-context examples, quantization levels, label semantics, and order of examples.
- We identify methods to improve the smoothness of LLM decision boundaries, such as fine-tuning earlier layers, fine-tuning on synthetic tasks and uncertainty-aware active learning.

2 Background

2.1 Training Large Language Models

Large Language Models (LLMs) are trained on vast corpora of text using unsupervised learning. During training, these models learn to predict the next token in a sequence. Given a sequence

of tokens $(x_1, x_2, \dots, x_{t-1})$, the model predicts the next token x_t by maximizing the likelihood $P(x_t|x_1, x_2, \dots, x_{t-1})$. The training objective typically involves minimizing the cross-entropy loss:

$$L = -\sum_{i=1}^{N} \sum_{t=1}^{T_i} \log P(x_t | x_1, x_2, \dots, x_{t-1})$$
 (1)

where T_i is the number of tokens in the *i*-th sequence and N is the total number of sequences in the corpus. During training, teacher forcing is often employed, where the model receives the ground truth token x_t as input at each time step instead of its own prediction, enabling parallel training.

2.2 In-Context Learning in LLMs

After training, LLMs can generalize to new tasks through a mechanism known as in-context learning. Let $\mathcal{S} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ represent the set of n input-output pairs provided as examples in the prompt, where \mathbf{x}_i is an input and y_i is the corresponding output. Given a new input \mathbf{x}_{new} , the LLM is turned into a task-specific model that predicts the output \hat{y}_{new} by conditioning on the given examples: $P(\hat{y}_{\text{new}}|\mathbf{x}_{\text{new}}, \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\})$. In-context learning allows the LLM to perform tasks by leveraging the context provided by these examples, thereby inferring the task and generating appropriate responses for new inputs. This approach utilizes the model's ability to recognize patterns and apply learned knowledge without additional training or fine-tuning.

3 Methodology

We aim to better understand in-context learning in Large Language Models by investigating their decision boundaries on a series of binary classification tasks. To increase the generality of our framework, we evaluate several existing LLMs on different task distributions under different settings. We present the general framework here, and refer to Section 4 for specific experiment settings.

3.1 In-Context Classification

Consider a K-class classification task with a data distribution $p_{\text{data}}(\mathbf{x},y)$, where \mathbf{x} is the input feature and $y \in \{1,\ldots,K\}$ is the class label. To construct an in-context prompt, we sample n examples $(\mathbf{x}_i,y_i) \sim p_{\text{data}}$ for $i=1,\ldots,n$. Given a new test point \mathbf{x}_{test} , in-context learning constructs a prompt $P=(\mathbf{x}_1,y_1,\ldots,\mathbf{x}_n,y_n,\mathbf{x}_{\text{test}})$ by concatenating the n sampled examples and the test point. The prompt P is then fed to the LLM π , which predicts a class \hat{y} for \mathbf{x}_{test} .

We prompt the LLM with P and obtain its prediction for \mathbf{x}_{test} by choosing the most likely class in the next token distribution. Formally, let V denote the size of the LLM's vocabulary, and $\mathbf{l} \in \mathbb{R}^V$ be the vector of logit values for each of the tokens. To obtain a class prediction, we convert each class label i into a unique token id, say c(i) and choose the class with the maximum logit value as the predicted label for $\mathbf{x}_{\text{query}}$, i.e., $\hat{y} = \arg\max_{i \in \{1, \dots, K\}} l_{c(i)}$.

3.2 Decision Boundary Visualization

To visualize the decision boundary of a model π , we generate a grid of points covering the feature space defined by the in-context examples set \mathcal{S} . Let $\mathcal{S} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_k, y_k)\}$ represent the set of in-context examples, and $\mathbf{x}_{\min}, \mathbf{x}_{\max} \in \mathbb{R}^d$ denote the minimum and maximum values of the features in \mathcal{S} along each dimension. We create a uniform grid with G points along each dimension, resulting in a total of G^d grid points. The grid points are denoted as $\mathbf{X}_{\text{grid}} = \{\mathbf{x}_{\text{query}} \mid \mathbf{x}_{\text{query}} \in [\mathbf{x}_{\min}, \mathbf{x}_{\max}]^d, \mathbf{x}_{\text{query}} = \mathbf{x}_{\min} + i\Delta\mathbf{x}, i \in \{0, 1, \dots, G-1\}\}$ where $\Delta\mathbf{x} = \frac{1}{G-1}(\mathbf{x}_{\max} - \mathbf{x}_{\min})$ is the grid spacing along each dimension. Each point $\mathbf{x}_{\text{query}} \in \mathbf{X}_{\text{grid}}$ is a query input, and the model π is prompted with the sequence $(\mathbf{x}_1, y_1, \dots, \mathbf{x}_k, y_k, \mathbf{x}_{\text{query}})$ to predict the corresponding class label \hat{y} . The decision boundary is then visualized by plotting the predicted labels \hat{y} over the grid \mathbf{X}_{grid} .

4 Experiments

In this section, we examine existing LLMs through the lens of decision boundaries by conducting a series of binary classification tasks under varying conditions. Our experiments aim to address the following key questions:

- How do existing pretrained LLMs perform on binary classification tasks? §4.1
- How do different factors influence the decision boundaries of these models? §4.2
- How can we improve the smoothness of decision boundaries? §4.3

Classification Tasks. We investigate the decision boundary of LLMs by prompting them with n in-context examples of binary classification tasks, with an equal number of examples for each class. We generate classification datasets using scikit-learn [Pedregosa et al., 2011], creating three types of linear and non-linear classification tasks: linear, circle, and moon, each describing different shapes of ground-truth decision boundaries. Detailed information on the dataset generation can be found in Appendix G. In addition to the in-context examples, we calculate the in-context learning accuracy on a held-out test set of size 100. We sample in-context examples and test points from classification task and convert them into prompt, with an example shown in Appendix F.

Obtaining Decision Boundaries of Language Models. We study an extensive range of models, with sizes ranging from 1.3B to 13B parameters, including open-source models such as Llama2-7B, Llama3-8B, Llama2-13B, Mistral-7B-v0.1, and sheared-Llama-1.3B. We also extend our analysis to state-of-the-art closed-source LLMs, including GPT-40 and GPT-3.5-turbo. We generate the decision boundaries of the open-source models with 8-bit quantization due to computational constraints. We choose a grid size scale of 50 x 50, resulting in 2500 queries for each decision boundary. For the open-source models, we use the approach described in 3.2 to get predictions. For the closed-source models, we use the next token generation as the prediction.

4.1 Non-Smooth Decision Boundaries of LLMs.

Figure 2 compares the decision boundaries of 6 LLMs when provided with 128 in-context examples. Even on simple linearly separable classification problems, all of these models exhibit non-smooth decision boundaries. The decision boundaries vary significantly across models, indicating that these models have different reasoning abilities to interpret the same in-context data. All models show fragmented decision regions, which means small changes in the input features can result in different classifications. This raises concerns about the reliability of LLMs and their practical deployment, as even when test accuracy for classification is high (shown in Figure 3, where test accuracy increases with the number of context examples), the underlying decision boundary lacks generalization. We further demonstrate nonsmoothness in NLP text classification tasks by projecting text input into 2D space, as detailed in Appendix H. In the following sections, we will explore factors that affect decision boundary smoothness and investigate methods to improve smoothness.

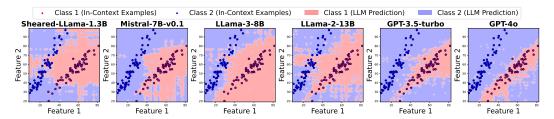


Figure 2: Visualizations of decision boundaries for various LLMs, ranging in size from 1.3B to 13B, on a linearly seperable binary classification task. The in-context data points are shown as scatter points and the colors indicate the label determined by each model. These decision boundaries are obtained using 128 in-context examples. The visualization highlights that the decision boundaries of these language models are not smooth.

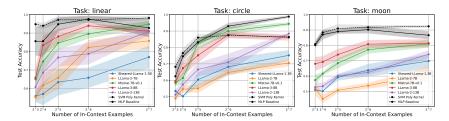


Figure 3: Test accuracy for LLMs and baselines across three classification tasks (linear, circle, and moon), with each subplot illustrating the test accuracy as the number of in-context examples increases. The baselines are the SVM with a polynomial kernel and the MLP with two hidden layers. Shaded regions represent the standard error of the mean accuracy across 5 seeds.

4.2 How Do Different Factors Influence the Decision Boundaries?

Impact of Model Size on Decision Boundary and Accuracy From Figure 2, model sizes increase from left to right, yet there is no clear correlation between model size and the smoothness of the decision boundary. Even the most powerful model, GPT-40, demonstrates fragmented decision regions. This suggests that increasing model size alone is insufficient for improving decision boundary smoothness. However, as shown in Figure 3, larger models tend to perform better in terms of test accuracy compared to smaller models, with Llama-1.3B often performing the worst.

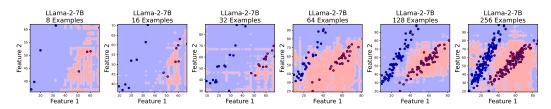


Figure 4: Decision boundary of Llama2-7b with increasing in-context examples from 8 to 256.

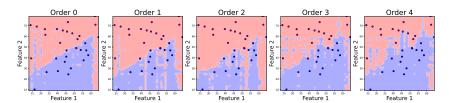
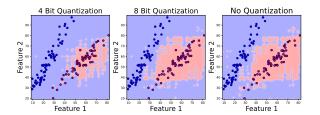
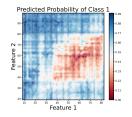


Figure 5: The sensitivity of the Llama3-8b model's decision boundary to the order of in-context examples. Each subplot (Order 0 to Order 4) shows the model's decision boundary with the same 32 examples shuffled differently.

Increasing In-Context Examples Does Not Guarantee Smoother Decision Boundaries While classification accuracies tend to improve with more in-context examples—and it's worth noting that Llama-3-8B and Mistral-7B's accuracy scales similarly to the SVM and MLP baselines—Figure 4 reveals that this does not translate to smoother decision boundaries. Despite the increase in accuracy, the decision boundaries remain fragmented, indicating that merely providing more in-context examples is not sufficient for achieving smoother decision regions.

How Does Quantization Influence Decision Boundaries? Figure 6a illustrates the decision boundaries of the LLaMA-2-7B model under different quantization levels [Dettmers et al., 2022]. When transitioning from 8-bit to 4-bit quantization, the red regions around the red in-context learning examples turn blue. This indicates that the reduced precision from 4-bit quantization significantly affects points near the decision boundary or areas where the model is most uncertain. For further investigation, we plot the probability prediction for class 1 (Figure 6b). The white regions, indicating a 50% probability for both classes, highlight the areas most impacted by quantization. Hence, varying quantization levels can flip the LLM's decisions in the regions of highest uncertainty.





- (a) Decision boundaries of Llama-2-7b with different quantization choices on a linearly seperable tsak.
- (b) Prediction of probability of class 1 with 8-bit quantization.

Figure 6: Impact of quantization on Llama2-7-8b's decision boundaries and probability predictions.

Are Decision Boundaries Sensitive to the Prompt Format? Yes, decision boundaries are sensitive to the labels' names, as shown in Figure 7. Using semantically unrelated labels, such as "Foo" and "Bar" as suggested in [Wei et al., 2023], results in flipped predictions compared to using reversed class names like "Bar" and "Foo". This suggests that the LLM's prediction still depend on its semantic prior knowledge of the labels.

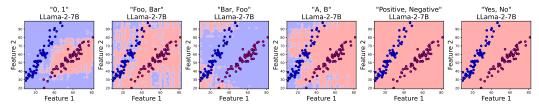


Figure 7: The decision boundaries of LLama-2-7B and LLama-3-8B, across various class labels. Each row corresponds to a model, and each column represents a different class label, shown in quotes.

Are Decision Boundaries Sensitive to the Order of In-Context Learning Examples? Recent works have shown that LLMs are sensitive to the order of in-context examples [Chen et al., 2024], which can significantly influence downstream performance. Similarly, as illustrated in Figure 5, we demonstrate that the model's decision boundaries vary with different shuffles of the in-context examples, highlighting the sensitivity of the decision boundaries to the order of the examples.

4.3 How to Improve the Decision Boundary Smoothness?

Can We Finetune LLMs on the In-Context Examples to Achieve Smoother Decision Boundaries? Our experiments indicate that finetuning LLMs on in-context examples does not result in smoother decision boundaries. Specifically, we finetuned Llama3-8B on 128 in-context learning examples and found that the resulting decision boundaries remained non-smooth. Examples of the decision boundaries after finetuning are provided in Appendix B.

Can We Finetune LLMs on a Dataset of Classification Tasks to Achieve Smoother Decision Boundaries? Previous works have shown that finetuning a pretrained LLM on a large collection of tasks improves its in-context learning performance on unseen tasks [Min et al., 2022a]. In this section, we investigate if the same paradigm helps improve the decision boundary smoothness of LLMs. To do this, we finetune a pretrained Llama model [Touvron et al., 2023] on a set of 1000 binary classification tasks generated from scikit-learn [Pedregosa et al., 2011], where the ground-truth decision boundary is either linear, circle-shaped, or moon-shaped, with equal probabilities. For each task, we sample randomly N=256 data points $x\sim \mathbf{X}_{\rm grid}$ and their corresponding label y's. We then sample the number of context points $m\sim \mathcal{U}[8,128]$, and finetune the LLM to predict $y_{i>m}$ given $x_{i>m}$ and the preceding examples:

$$\mathcal{L}(\pi) = \mathbb{E}\left[\sum_{i=m+1}^{N} \log p(y_i \mid x_i, x_{1:i-1}, y_{1:i-1})\right],$$
(2)

where the expectation is with respect to task, data points $\{(x_i, y_i)\}_{i=1}^N$, and the number of context points m. After training, we evaluate the same finetuned model on various binary classification tasks

with varying numbers of context points. To ensure the test tasks are unseen during training, we use different parameters in creating the datasets, such as the separateness between two classes and the scale between the inner and outer circles in the circle task. See Appendix G for more details.

We consider several finetuning settings for ablation studies. 1) In the first setting, we finetune the pretrained LLM using LoRA [Hu et al., 2021] and finetune the attention layers. 2) We finetune only the token embedding layer of LLM. 3) We finetune only the linear head layer of LLM. Then we consider modifying the architecture of the LLM: In this setting, we keep the core transformer backbone of the LLM frozen, attach randomly initialized embedding layers and prediction head to the model, and train the entire model using objective (2). This stems from the intuition that task-specific embedding and prediction layers allow the model to maximally utlize the general pattern-matching capabilities of the transformer backbone for the new task. We refer to this model as CustomLLM, and consider its three variants, which add 1) a new embedding layer for x, 2) a new prediction head for y, and 3) new embedding layers for x, y, and a new prediction head for y. The embedding layers and prediction head are MLPs with one hidden layer. We embed the raw numerical values instead of the text representation of x whenever a new embedding layer for x is used (same for y), and predict directly the binary class values instead of text labels whenever the new prediction head is used. Results of Finetuning LLM and CustomLLM in Figure 8 and Figure 9 show that finetuning the intermediate and earlier embedding layers leads to smoother decision boundary compared to finetuning the top prediction head.

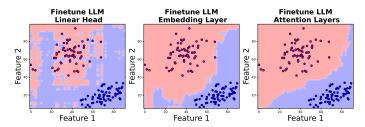


Figure 8: Decision boundary of Llama3-8B post finetuning the linear head, embedding layer and the attention layers. Finetuning the latter two layers improves the smoothness.

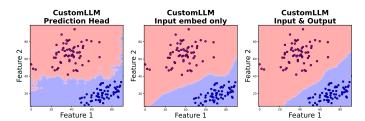
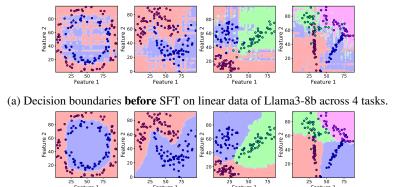


Figure 9: CustomLLM finetuning ablations. Decision boundary after finetuning the prediction head, input embedding layer and both layers for the CustomLLM.

Can LLMs finetuned on one in-context learning task generalize to more complex in-context learning tasks? In this section, we further explore whether a LLM fine-tuned only on a linear task can achiever smoother decision boundaries on unseen and more complex tasks. As shown in Figure 10, we compare the decision boundaries of Llama3-8b before and after SFT on the linear task only. Unexpectedly, we found it generalizes to unseen non-linear tasks as well as 3-class and 4-class classification tasks, despite only being trained on a binary linear task. The smoother decision boundaries observed in these unseen tasks suggest that fine-tuning on a synthetic in-context learning task can have downstream benefits for other tasks, enabling the model to be more robust in in-context learning.

Can we train a transformer from scratch to learn smooth decision boundary in-context? One may wonder whether a small transformer trained from scratch can provide smooth decision boundaries. To answer this, we train TNPs [Nguyen and Grover, 2022], a transformer-based model specifically designed for in-context learning. For each sequence of data points $\{(x_i,y_i)\}_{i=1}^N$ from a task C, TNPs learn to predict the query labels $y_{i>m}$ given the query inputs $x_{i>m}$ and the context pairs, assuming



(b) Decision boundaries after SFT on linear data of Llama3-8b across 4 unseen tasks.

Figure 10: Generalization ability of Llama-3-8B after supervised fine-tuning on a single binary linear classification task. The first two columns show the model's performance on non-linear classification tasks before and after fine-tuning, while the last two columns demonstrate its ability to generalize to 3-class and 4-class classification tasks.

conditional independence among the queries given the context:

$$\mathcal{L}(\theta) = \mathbb{E}\left[\sum_{i=m+1}^{N} \log p(y_i \mid x_i, x_{1:m}, y_{1:m})\right],\tag{3}$$

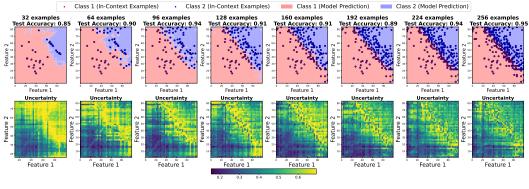
where the expectation is with respect to task C, data points $\{(x_i,y_i)\}_{i=1}^N$, and the number of context points m. TNPs employ a specialized mask to ensure the conditional independence assumption. We showed in Appendix D that transformers trained from sctrach can learn to in-context learn smooth decision boundary. Details are in Appendix D.

How to Use Uncertainty-aware Active Learning to Learn Decision Boundaries We investigate whether the decision boundary can be smoothed by providing the LLM with labels of the most uncertain points on the grid as additional in-context examples. Uncertainty is measured as the entropy of the probability distribution of the two classes after softmax normalization of the logits. Our study focuses on an active learning scheme where new in-context examples are incrementally added based on the LLM's current uncertainty. Initially, we obtain the decision boundary conditioned on the existing in-context examples. To refine this boundary, we query the LLM over a grid and select the top-k most uncertain points, ensuring they are spatially distant from each other using a greedy sampling approach. For labeling these uncertain points, we use a logistic regression model well-trained on a larger dataset with perfect accuracy as the ground truth decision boundary. As shown in Figure 11, this uncertainty-aware active sampling method results in a smoother decision boundary over iterations compared to random sampling. The iterative refinement enhances the model's generalization capabilities, leading to higher test set accuracies and greater sample efficiency, requiring fewer additional in-context examples to achieve performance gains. These findings indicate that leveraging the LLM's uncertainty measurements is valuable for selecting new in-context examples in resource-constrained settings where labeled data is scarce. We show more examples in Appendix I.

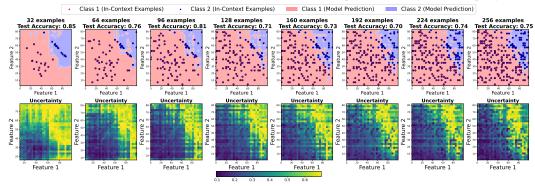
5 Related Works

Understanding in-context learning in transformers and LLMs is an active area of research, with existing works approaching this problem from both theoretical and practical perspectives.

Theoretical understanding of in-context learning Recent works aim to establish a theoretical connection between in-context learning and gradient descent (GD). The pioneering work by Akyürek et al. proves transformers can implement learning algorithms for linear models based on GD and closed-form ridge regression by construction. Von Oswald et al. [2023] proves the equivalence between linear self-attention and GD on linear regression by construction. Similarly, Dai et al.



(a) Decision boundaries with different numbers of context examples when using active sampling.



(b) Decision boundaries with different numbers of context examples when using random sampling.

Figure 11: Comparison of active and random sampling methods. We plot the decision boundaries and uncertainty plot across different number of in-context examples from 32 to 256, where the in-context examples are gradually added to the prompt using active or random methods. Active sampling gives smoother decision boundary and the uncertain points lie on it. The test set accuracies is plotted in the titles.

[2023] shows that attention in transformers has a dual form of GD and views transformers as meta-optimizers. Subsequent works extend these ideas to characterize the global optimum of single-layer linear transformers. Ahn et al. [2024] observe that with the optimal parameters, the transformer implements a single step of preconditioned gradient descent, while Zhang et al. [2023] shows that at the global optimum, the transformer achieves competitive prediction error with the best linear predictor on a new prediction task. In addition to theoretical connections to GD, a complementary direction aims to establish statistical complexity and generalization bounds of in-context learning in transformers [Bai et al., 2024, Li et al., 2023b, Wies et al., 2024, Wu et al., 2023]. The common limitation of these existing theoretical frameworks is the reliance on strong assumptions about the transformer architecture or the functional form of the in-context learning tasks which may not reflect real-world practices.

Practical understanding of in-context learning More relevant to our paper is a line of works focusing on understanding the practical aspects of in-context learning in LLMs. Many existing works investigate the roles of in-context examples and prompts. Min et al. [2022b] show a surprising result that ground-truth demonstrations are not required for in-context learning, while other factors such as the label space, input text distribution, and overall sequence format play an important role. Shi et al. [2023] investigate the distractibility of LLMs and shows that their performance dramatically drops when irrelevant context is included. Subsequently, Wei et al. [2023] characterize these behaviors of LLMs with respect to model size, and show that larger language models perform in-context learning differently in the presence of flipped or semantically unrelated labels. Webson and Pavlick [2022] argue against the current practice of prompt engineering, showing that intentionally irrelevant or even pathologically misleading prompts achieve similar downstream performance to instructively good prompts. Orthogonally, Lampinen et al. [2022] find that including explanations in the in-context

examples significantly improves the few-shot performance of LLMs. Finally, given the expanded context windows of modern LLMs, recent works have explored in-context learning in the many-shot setting with hundreds or thousands of examples [Agarwal et al., 2024, Li et al., 2023a, Bertsch et al., 2024].

Learning to learn in-context In contrast to the emergent in-context capabilities of LLMs, existing works have also studied methods that learn to perform in-context learning explicitly. Min et al. [2022a] propose MetaICL, a meta-training framework for finetuning pretrained LLMs to perform in-context learning on a large and diverse collection of tasks. MetaICL outperforms several baselines including emergent in-context learning and multi-task learning followed by zero-shot transfer. Beyond text, TNP [Nguyen and Grover, 2022, Nguyen et al., 2023, Nguyen and Grover, 2024] and PFNs [Müller et al., 2021] propose to train transformer models to perform in-context prediction for a family of functions, which allows in-context generalization to unseen functions after training. Similarly, Garg et al. [2022] show that autoregressive transformers can be trained from scratch to learn function classes such as linear functions and 2-layer ReLU networks. Other work also shows that alignment can be done in-context [Zhao et al., 2023], where in-context learned reward model can be used for inference-time preference alignment. These works present an interesting set of baselines for our work to examine the in-context learning ability of LLMs.

6 Conclusion

We propose a novel approach to understanding in-context learning in LLMs by probing their decision boundaries in in-context learning in binary classification tasks. Despite achieving high test accuracy, we observe that the decision boundaries of LLMs are often irregularly non-smooth. Through extensive experiments, we identify factors that affect this decision boundary. We also explore fine-tuning and adaptive sampling methods, finding them effective in improving boundary smoothness. Our findings provide new insights into the mechanics of in-context learning and suggest pathways for further research and optimization.

Acknowledgments

This research is supported by NSF CAREER Award 2341040, Schmidt Sciences AI2050 Fellowship, Samsung, and Cisco.

References

- J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- R. Agarwal, A. Singh, L. M. Zhang, B. Bohnet, S. Chan, A. Anand, Z. Abbas, A. Nova, J. D. Co-Reyes, E. Chu, et al. Many-shot in-context learning. arXiv preprint arXiv:2404.11018, 2024.
- K. Ahn, X. Cheng, H. Daneshmand, and S. Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. Advances in Neural Information Processing Systems, 36, 2024.
- E. Akyürek, D. Schuurmans, J. Andreas, T. Ma, and D. Zhou. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*.
- Y. Bai, F. Chen, H. Wang, C. Xiong, and S. Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *Advances in neural information processing systems*, 36, 2024.
- A. Bertsch, M. Ivgi, U. Alon, J. Berant, M. R. Gormley, and G. Neubig. In-context learning with long-context models: An in-depth exploration. *arXiv preprint arXiv:2405.00200*, 2024.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- X. Chen, R. A. Chi, X. Wang, and D. Zhou. Premise order matters in reasoning with large language models. *arXiv preprint arXiv:2402.08939*, 2024.
- A. Conneau and D. Kiela. Senteval: An evaluation toolkit for universal sentence representations. *arXiv* preprint arXiv:1803.05449, 2018.
- D. Dai, Y. Sun, L. Dong, Y. Hao, S. Ma, Z. Sui, and F. Wei. Why can gpt learn in-context? language models secretly perform gradient descent as meta-optimizers. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4005–4019, 2023.
- T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale, 2022.
- S. Garg, D. Tsipras, P. S. Liang, and G. Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- E. J. Hu, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. d. l. Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- A. Lampinen, I. Dasgupta, S. Chan, K. Mathewson, M. Tessler, A. Creswell, J. McClelland, J. Wang, and F. Hill. Can language models learn from explanations in context? In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 537–563, 2022.
- M. Li, S. Gong, J. Feng, Y. Xu, J. Zhang, Z. Wu, and L. Kong. In-context learning with many demonstration examples. *arXiv preprint arXiv:2302.04931*, 2023a.
- Y. Li, M. E. Ildiz, D. Papailiopoulos, and S. Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*, pages 19565–19594. PMLR, 2023b.

- P. Malo, A. Sinha, P. Korhonen, J. Wallenius, and P. Takala. Good debt or bad debt: Detecting semantic orientations in economic texts. *Journal of the Association for Information Science and Technology*, 65(4):782–796, 2014.
- S. Min, M. Lewis, L. Zettlemoyer, and H. Hajishirzi. Metaicl: Learning to learn in context. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2791–2809, 2022a.
- S. Min, X. Lyu, A. Holtzman, M. Artetxe, M. Lewis, H. Hajishirzi, and L. Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11048–11064, 2022b.
- I. Mollas, Z. Chrysopoulou, S. Karlos, and G. Tsoumakas. Ethos: an online hate speech detection dataset. arXiv preprint arXiv:2006.08328, 2020.
- S. Müller, N. Hollmann, S. P. Arango, J. Grabocka, and F. Hutter. Transformers can do bayesian inference. In *International Conference on Learning Representations*, 2021.
- T. Nguyen and A. Grover. Transformer neural processes: Uncertainty-aware meta learning via sequence modeling. In *International Conference on Machine Learning*, pages 16569–16594. PMLR, 2022.
- T. Nguyen and A. Grover. Lico: Large language models for in-context molecular optimization. *arXiv* preprint arXiv:2406.18851, 2024.
- T. Nguyen, S. Agrawal, and A. Grover. Expt: Synthetic pretraining for few-shot experimental design. *Advances in Neural Information Processing Systems*, 36:45856–45869, 2023.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- M. Reid, N. Savinov, D. Teplyashin, D. Lepikhin, T. Lillicrap, J.-b. Alayrac, R. Soricut, A. Lazaridou, O. Firat, J. Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- F. Shi, X. Chen, K. Misra, N. Scales, D. Dohan, E. H. Chi, N. Schärli, and D. Zhou. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*, pages 31210–31227. PMLR, 2023.
- R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642, 2013.
- H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023.
- J. Von Oswald, E. Niklasson, E. Randazzo, J. Sacramento, A. Mordvintsev, A. Zhmoginov, and M. Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference* on *Machine Learning*, pages 35151–35174. PMLR, 2023.
- A. Wang, Y. Pruksachatkun, N. Nangia, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32, 2019.
- A. Webson and E. Pavlick. Do prompt-based models really understand the meaning of their prompts? In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2300–2344, 2022.
- J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, et al. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022.

- J. Wei, J. Wei, Y. Tay, D. Tran, A. Webson, Y. Lu, X. Chen, H. Liu, D. Huang, D. Zhou, et al. Larger language models do in-context learning differently. *arXiv preprint arXiv:2303.03846*, 2023.
- N. Wies, Y. Levine, and A. Shashua. The learnability of in-context learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- J. Wu, D. Zou, Z. Chen, V. Braverman, Q. Gu, and P. Bartlett. How many pretraining tasks are needed for in-context learning of linear regression? In *The Twelfth International Conference on Learning Representations*, 2023.
- M. Xia, T. Gao, Z. Zeng, and D. Chen. Sheared llama: Accelerating language model pre-training via structured pruning. *arXiv preprint arXiv:2310.06694*, 2023.
- R. Zhang, S. Frei, and P. L. Bartlett. Trained transformers learn linear models in-context. *arXiv* preprint arXiv:2306.09927, 2023.
- X. Zhang, J. Zhao, and Y. LeCun. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28, 2015.
- S. Zhao, J. Dang, and A. Grover. Group preference optimization: Few-shot alignment of large language models. *arXiv preprint arXiv:2310.11523*, 2023.

A Pretrained LLMs decision boundary on linear and non-linear classification tasks

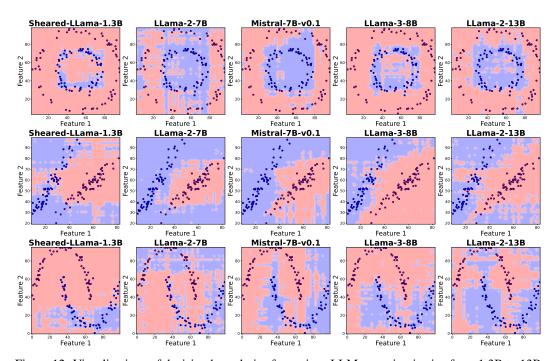


Figure 12: Visualizations of decision boundaries for various LLMs, ranging in size from 1.3B to 13B, on three classification tasks. The tasks are, from top to bottom, circle, linear, and moon classifications. Note that the circle and moon tasks are not linearly separable. The in-context data points are shown as scatter points and the colors indicate the label determined by each model. These decision boundaries are obtained using 128 in-context examples. The visualization highlights that the decision boundaries of these language models are not smooth.

B Finetune on in-context examples only

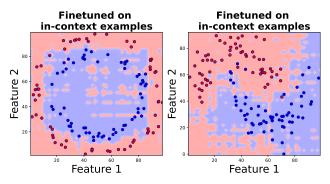


Figure 13: Two examples of Llama2-7B finetuned on the in-context examples points, which are scattered points in the plot.

C SFT LLMs for in-context classification

We used LoRA [Hu et al., 2021] to supervise fine-tune the Llama series models on both non-linear and linear classification tasks, including circle, linear, and moon datasets. The models fine-tuned are Sheared-Llama-1.3B, Llama2-7B, Llama2-13B, and Llama3-8B. Visualization in Figure 14 demonstrates that these language models produce smoother decision boundaries after training on the classification datasets using SFT.

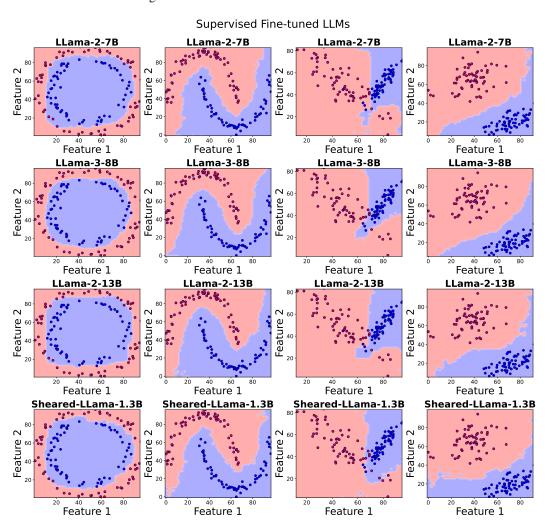


Figure 14: Decision boundary of in-context learning on 128 examples across Llama series models after supervised finetuning with LoRA.

D Training Transformers from Scratch: TNP models decision boundaries

We trained TNP models of four different sizes as shown in the Table 1 below. We plot how does the TNP models decision boundary changes as more in-context examples are added in Figure 15. TNP models learn smooth deicision boundary for this moon-shaped non-linear task. And we did not observe a scaling law of transformer sizes versus the decision boundary smoothness. In contrast the smaller model generalize better than the larger model.

Table 1: TNP	transformers	model sizes	and a	rchitectures

Model	Parameters (M)	Input embed dim	feedforward dim	num heads	num layers
Small	0.1	64	64	2	3
Medium	0.6	128	128	4	6
Large	1.6	128	256	8	12
X-Large	9.7	256	512	16	18

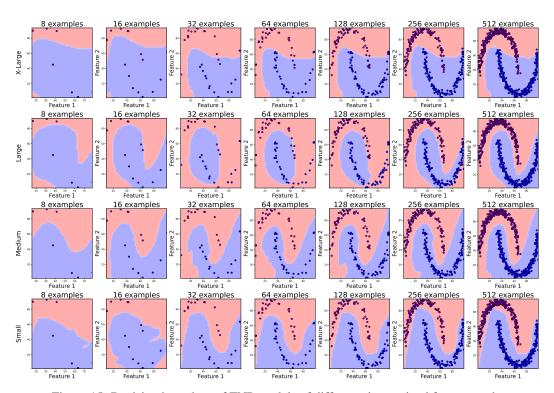


Figure 15: Decision boundary of TNP models of different sizes trained from scratch.

E Traditional Classifiers Model Details

In our experiments, we used several classical machine learning models with the following hyperparameters:

- **Decision Tree Classifier:** We set the maximum depth of the tree to 3.
- **Multi-Layer Perceptron:** The neural network consists of two hidden layers, each with 256 neurons, and the maximum number of iterations is set to 1000.
- **K-Nearest Neighbors:** The number of neighbors is set to 5.
- Support Vector Machine (SVM): We used a radial basis function (RBF) kernel with a gamma value of 0.2.

F Prompt Format for binary classification

```
Given pairs of numbers and their labels, predict the label for a new
input pair of numbers based on the provided data.
Answer with only one of the labels 'Foo' and 'Bar':
Input:
       64 24
Label: Bar
Input: 34 41
Label: Bar
Input: 71 66
Label: Bar
Input: 96 49
Label: Foo
Input: 21 56
Label: Foo
What is the label for this input?
Input: 23
Label:
```

Figure 16: Few-shot in-context prompt with n context questions.

G Classification Datasets

We use three types of classification tasks from scikit-learn [Pedregosa et al., 2011] to probe the decision boundary of LLMs and transformers: linear, circle, and moon classification problems. For linear classification tasks, we utilize the make_classification function, which generates random classification problems by creating clusters of points normally distributed around the vertices of a hypercube with sides of length 2 × class_sep. Circle classification tasks are generated using the make_circles function, creating a binary classification problem with a large circle containing a smaller circle. The factor parameter controls the scale of the inner circle relative to the outer circle. Moon classification tasks are generated using the make_moons function, creating a binary classification problem with two interleaving half circles. The noise parameter controls the standard deviation of Gaussian noise added to the data points.

For training tasks, the class_sep parameter is randomly sampled from the range [1.5,2], and the factor parameter for circular tasks is sampled from [0.1,0.4]. For testing tasks, the class_sep parameter is sampled from [1,1.4], and the factor parameter from [0.5,0.9], ensuring that testing tasks differ from training tasks. The noise parameter for moon-shaped tasks is sampled from [0.05,0.1] for training and [0.1,0.2] for testing, introducing varying levels of complexity in the classification problems.

H Decision Boundary of LLMs on NLP tasks.

We extend our analysis to multi-class NLP classification tasks using high-dimensional real-world datasets. To address the challenge of visualizing high-dimensional text embeddings, we project them onto a 2D space using t-SNE and send the 2D embeddings as input in the prompt to the LLM. While any dimensionality reduction technique inevitably introduces confounding factors, this approach allows us to extend our analysis to more complex, real-world scenarios. Our experiments encompass six widely-used NLP classification tasks, covering both binary and multi-class settings. These include Subjective/Obejective sentence classification (SUBJ) [Conneau and Kiela, 2018], financial sentiment analysis (FP) [Malo et al., 2014], textual entailment recognition (RTE) [Wang et al., 2019], hate speech detection (ETHOS) [Mollas et al., 2020], sentiment analysis (SST-2) [Socher et al., 2013]

and news topic classification (AG_NEWS) [Zhang et al., 2015]. The results, presented in Figure 17, demonstrate that the non-smooth decision boundary characteristics observed in our synthetic datasets persist in these more complex NLP tasks.

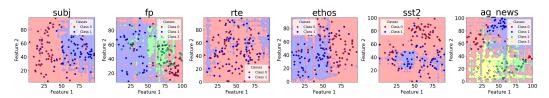


Figure 17: Decision boundaries of Llama-3-8b on six NLP tasks, ranging from binary to multi-class classification. Since text embeddings are natively high-dimensional, we projected text embeddings onto a 2D space using t-SNE. The irregular, non-smooth behaviors are also seen in these tasks.

I Uncertainty Aware Active Sampling For Smoother Decision Boundary and Better Test set Accuracy

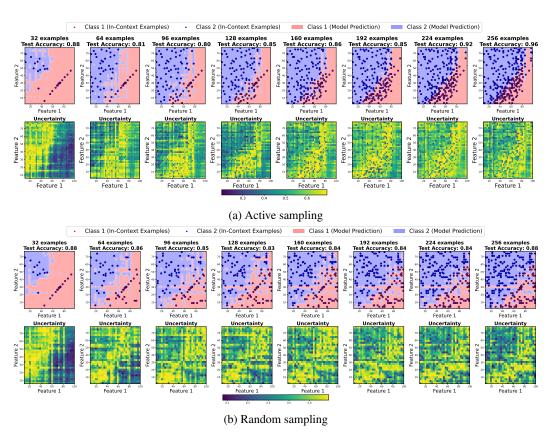


Figure 18: Comparison of active and random sampling methods.

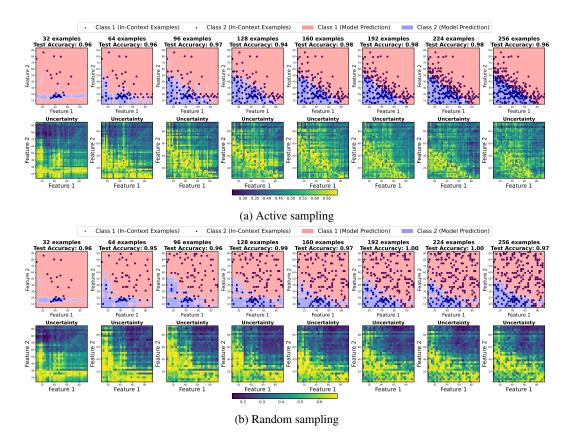


Figure 19: Comparison of active and random sampling methods.

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: It reflects our proposed method and experimental findings.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: We provide limitation in the appendix section.

Guidelines:

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

3. Theory Assumptions and Proofs

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

Answer: [NA]

Justification: The paper does not include theoretical studies.

Guidelines:

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

4. Experimental Result Reproducibility

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

Answer: [Yes]

Justification: We present complete method and experiment details in Section 3 and 4.

Guidelines:

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

5. Open access to data and code

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

Answer: [Yes]

Justification: All code and checkpoints will be released publicly upon acceptance.

Guidelines:

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

6. Experimental Setting/Details

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

Answer: [Yes]

Justification: We present complete method and experiment details in Section 3 and 4.

Guidelines:

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

7. Experiment Statistical Significance

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

Answer: [Yes]

Justification: Although we do not provide error bar for every plot, this is due to the nature of our work, since we are visualizing the decision boundary for qualitative understanding. We justify this with additional plots in various settings in the appendix. Apart from the decision boundary plots, we do plot the accuracy plot with error bar.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

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

Answer: [Yes]

Justification: We provide the compute in the appendix.

Guidelines:

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

9. Code Of Ethics

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

Answer: [Yes]

Justification: The paper conforms to the NeurIPS Code of Ethics.

Guidelines:

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

10. Broader Impacts

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

Answer: [NA]

Justification: Our work does not have societal impact.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [NA]
Justification: NA

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes],

Justification: Yes, we cited every dataset and models we used.

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

 If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

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

Answer: [NA]

Justification: We don't release new dataset.

Guidelines:

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

14. Crowdsourcing and Research with Human Subjects

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

Answer: [NA]

Justification: we don't have any.

Guidelines:

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

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

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

Answer: [NA]

Justification: We don't have this.

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