
Your Finetuned Large Language Model is Already a Powerful Out-of-distribution Detector

Andi Zhang^{1,*}
az381@cam.ac.uk

Tim Z. Xiao^{2,3,4}
zhenzhong.xiao@uni-tuebingen.de

Weiyang Liu^{1,4}
wl396@cam.ac.uk

Robert Bamler²
robert.bamler@uni-tuebingen.de

Damon Wischik¹
djw1005@cam.ac.uk

¹University of Cambridge ²University of Tübingen ³IMPRS-IS

⁴Max Planck Institute for Intelligent Systems, Tübingen

Abstract

We revisit the likelihood ratio between a pretrained large language model (LLM) and its finetuned variant as a criterion for out-of-distribution (OOD) detection. The intuition behind such a criterion is that, the pretrained LLM has the prior knowledge about OOD data due to its large amount of training data, and once finetuned with the in-distribution data, the LLM has sufficient knowledge to distinguish their difference. Leveraging the power of LLMs, we show that, for the first time, the likelihood ratio can serve as an effective OOD detector. Moreover, we apply the proposed LLM-based likelihood ratio to detect OOD questions in question-answering (QA) systems, which can be used to improve the performance of specialized LLMs for general questions. Given that likelihood can be easily obtained by the loss functions within contemporary neural network frameworks, it is straightforward to implement this approach in practice. Since both the pretrained LLMs and its various finetuned models are available, our proposed criterion can be effortlessly incorporated for OOD detection without the need for further training. We conduct comprehensive evaluation across on multiple settings, including far OOD, near OOD, spam detection, and QA scenarios, to demonstrate the effectiveness of the method.

1 Introduction

Detecting out-of-distribution (OOD) is crucial for the safety of artificial intelligence systems. OOD detection aims to identify inputs that substantially deviate from the training data the model was trained on, ensuring the system is alerted to these discrepancies. This capability to identify OOD and anomalous data is particularly critical in high-stakes domains such as healthcare and autonomous driving, where the stakes for accuracy and reliability are exceptionally high.

In OOD detection (Hendrycks & Gimpel, 2016), the term “in-distribution” refers specifically to the distribution of the training data. In natural language processing, OOD detection has been studied in *small and task-specific models* for settings such as translation (Xiao et al., 2020) and question answering (Lyu et al., 2020). With the advancement of *large and general models* like large language models (LLMs), the scope of training data has significantly broadened, positioning these large models embodied with general knowledge and intelligence as “base models”. In the era of large foundation models (Bommasani et al., 2021), the prevalent training paradigm has shifted from end-to-end learning towards finetuning a pretrained base model. This shift calls for a revisiting of the definition of “in-distribution” and its relationship to the training distribution of base models.

*Corresponding Author

There is usually no prior information about the OOD data in conventional OOD detection (Hendrycks & Gimpel, 2016). Bishop (1994) introduced the idea that OOD data could be viewed as data coming from a distinct OOD distribution, suggesting the use of the likelihood ratio between this OOD distribution and the in-distribution as a detection criterion. However, given the absence of information about OOD data, Zhang & Wischik (2022) introduced the concept of OODProxy, a conceptual framework where the OOD distribution is represented by a proxy, incorporating what is known or assumed about the OOD data. From this perspective, it is commonly recognized that many OOD detection methods leveraging likelihood ratios are essentially applying their unique OOD proxy distributions, each reflecting different assumptions or prior knowledge concerning the nature of the OOD data (Ren et al., 2019; Serrà et al., 2019; Schirmermeister et al., 2020; Zhang et al., 2021).

In our paper, we propose that pretrained base models can function as a repository of prior knowledge for OOD data relative to in-distribution data, effectively acting as OOD proxy distributions. Guided by this insight, we discover that the likelihood ratio between the base model and its finetuned counterpart serves as an effective criterion for detecting OOD data. Moreover, for LLM-based question-answering (QA) systems, the same likelihood ratio excels in detecting OOD questions. By identifying and rejecting these OOD questions, we can greatly enhance the robustness of current QA systems.

The convenience of obtaining likelihood from loss functions in current neural networks enables a simple and straightforward implementation of this method in practice. Moreover, it is worth mentioning that numerous practitioners likely already have both a pretrained and a finetuned LLM at their disposal. This setup inherently equips them with the capacity for OOD detection without necessitating additional training efforts.

2 Background and Preliminaries

OOD Detection. We start with an in-distribution dataset, denoted as \mathcal{D}_{in} , assuming that the data within \mathcal{D}_{in} is sampled from an in-distribution probability distribution p_{in} . The objective of OOD detection is to determine whether a given input data x originates from p_{in} . Hendrycks & Gimpel (2016) were the first to highlight the significance of this problem in deep learning era and introduced a practical benchmark for evaluation. This involves training a model on \mathcal{D}_{in}^{train} , the training subset of the in-distribution dataset, with the model providing a detection criterion S for identifying OOD data. We then gather a dataset from domains different from \mathcal{D}_{in} , labeled as \mathcal{D}_{out} . The effectiveness of the criterion S is assessed by applying it to data from $\mathcal{D}_{in}^{test} \cup \mathcal{D}_{out}^{test}$ and evaluating the performance using metrics such as AUROC, AUPR, and FPR95 (Yang et al., 2022). These metrics help determine how well S can differentiate between the in-distribution dataset and the OOD dataset. A high performance across these metrics signifies a robust OOD detection capability.

In Hendrycks & Gimpel (2016)’s foundational study, it is assumed that for the in-distribution dataset \mathcal{D}_{in} , each data point is accompanied by a classification label. Our approach diverges slightly from this premise, as **we operate under the assumption that no labels are available**. This scenario is often referred to as “unsupervised OOD detection” or “OOD detection without in-distribution labels”.

The Paradox in Unsupervised OOD Detection. In the context of unsupervised OOD detection, Nalisnick et al. (2018) revisit Bishop (1994)’s suggestion that a probabilistic generative model p_{θ} could model the in-distribution p_{in} , proposing to use the model output $S(x) = p_{\theta}(x)$ as a criterion for evaluating a given input x . Surprisingly, Nalisnick et al. discovered that in certain scenario - such as when \mathcal{D}_{in} is CIFAR10 and \mathcal{D}_{out} is SVHN, or \mathcal{D}_{in} is FashionMNIST and \mathcal{D}_{out} is MNIST - the OOD data received higher $p_{\theta}(x)$ scores than the in-distribution data. This counterintuitive finding has been labeled as a “paradox”.

OOD Proxy. To address the paradox, the studies (Ren et al., 2019; Serrà et al., 2019; Schirmermeister et al., 2020; Zhang et al., 2021) have put forward the idea of utilizing the likelihood ratio as the criterion for identifying OOD data. Zhang & Wischik (2022) integrates these techniques into a comprehensive structure termed the OOD proxy framework. Within this framework, it is posited that in-distribution data can be characterized as samples from

a distribution p_{in} , whereas OOD data are samples from a distribution p_{out} . According to the Neyman-Pearson lemma, the likelihood ratio represents the optimal criterion for OOD detection in theory, which is mathematically expressed as:

$$S(x) = \frac{p_{\text{out}}(x)}{p_{\text{in}}(x)}$$

where obtaining p_{out} is challenging. To address this, Zhang & Wischik (2022) introduced the concept of utilizing a proxy distribution, $p_{\text{out}}^{\text{proxy}}$, which incorporates human subjective understanding of the OOD distribution. For instance, Ren et al. (2019) define $p_{\text{out}}^{\text{proxy}}$ as the distribution representing background statistics; Serrà et al. (2019) consider it as the distribution of data compression; Schirrmeyer et al. (2020) describe $p_{\text{out}}^{\text{proxy}}$ as a general distribution; and for Zhang et al. (2021), $p_{\text{out}}^{\text{proxy}}$ corresponds to the distribution of a local autoregressive model.

In the OOD proxy framework, the use of likelihood as a detection criterion, following the approach of Bishop (1994) and Nalisnick et al. (2018), essentially assumes a uniform distribution for $p_{\text{out}}^{\text{proxy}}$. The suboptimal performance associated with this method is not surprising; it highlights the limitations of an improper prior assumption.

Likelihood of Autoregressive Language Models Autoregressive language models (Brown et al., 2020) are types of probabilistic models that compute the likelihood of a sentence $x = x_1, \dots, x_T$, where T denotes the sentence length and x_t represents each word at position t . By definition of conditional probability, the likelihood $p(x) = p(x_1, \dots, x_T)$ can be calculated as the product of conditional probabilities: $p(x_1)p(x_2|x_1)p(x_3|x_1, x_2) \dots p(x_T|x_1, \dots, x_{T-1})$. Each term $p(x_t|x_{<t})$ represents the probability of predicting the subsequent word x_t given all the previous words x_1, \dots, x_{t-1} . In neural language models, these conditional probabilities are typically modeled using a softmax function over the output vocabulary.

When training or finetuning a language model on a dataset, the objective is to maximize the likelihood of the training data. This is equivalent to minimizing the negative log-likelihood, which is the cross-entropy loss between the predicted word probabilities and the true word labels at each position. Most modern neural network libraries provide built-in functions to compute this cross-entropy loss, making it straightforward to optimize the model parameters to minimize the negative log-likelihood and thus maximize the likelihood of the training data.

3 Pretrained Large Language Model as a OOD Proxy

Considering an autoregressive language model as a distribution p , denote p_θ as the pretrained large language model with parameters θ . Given an in-distribution dataset \mathcal{D}_{in} , the model finetuned on \mathcal{D}_{in} is represented with parameters θ' . For an input x , we introduce the out-of-distribution detection criterion S as follows:

$$S(x) = \frac{p_\theta(x)}{p_{\theta'}(x)}. \quad (1)$$

This criterion essentially employs the pretrained large language model as an OOD proxy introduced in Section 2. This strategy is particularly practical given the widespread availability of pretrained LLMs. Finetuning these models to adapt their distribution for specific domain contexts is a standard practice, meaning many practitioners may already possess finetuned LLMs that represent the distribution of their specific datasets. With both pretrained and finetuned models at hand, calculating the likelihood ratio becomes straightforward, eliminating the need for additional training. For instance, suppose we possess a LLM that has been finetuned on legal documents. Given a new document x , we can determine whether it is a legal document by utilizing the likelihood ratio $S(x)$.

Utilizing a pretrained LLM as an OOD proxy stems from the rationale that the OOD proxy distribution should encapsulate general characteristics of OOD data, including prior knowledge or subjective insights about the OOD data. While it's conceivable to presume

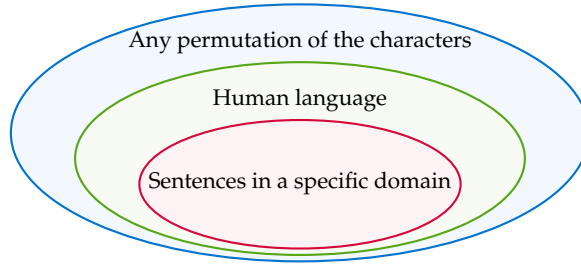


Figure 1: Relationship among sentences within a specific domain, the comprehensive set of human language, and all conceivable character permutations.

OOD data emanate from a uniform distribution devoid of any prior knowledge, evidence presented by Nalisnick et al. (2018) indicates that relying on a uniform distribution as the OOD proxy can be ineffective in practice. In the context of language models, the well-known infinite monkey theorem¹ (Borel, 1913; Eddington, 2019) suggests that any text could theoretically be generated from a uniform distribution; however, this does not serve as an effective representation of any coherent language. Figure 1 depicts the relationship among sentences within a specific domain, the comprehensive set of human language, and all conceivable character permutations. In fact, coherent human language forms a minor subset of all potential character arrangements. Consequently, the assumption that OOD data ought to represent meaningful human language constitutes a strong prior.

Therefore, we advocate for the use of a pretrained LLM as a more suitable OOD proxy. Given their extensive parameters and training on vast corpora (for example, the Llama-2 model, as mentioned by Touvron et al. (2023), is trained on 7 trillion tokens), it is plausible to consider that an LLM encompasses the breadth of human language. Assuming the LLM estimates the distribution of all human language, it is logical to designate the pretrained LLM as the OOD proxy.

4 Likelihood Ratio OOD Detection for QA Systems

In question-answering (QA) systems, identifying OOD questions is crucial for enhancing system robustness through their rejection. However, detecting OOD questions is challenging due to the often brief and uninformative nature of the questions submitted to QA systems, rendering the direct application of the likelihood ratio on the questions themselves ineffective for OOD detection. To overcome this issue, we leverage the observation that while a finetuned LLM generates pertinent answers to in-distribution questions, it tends to produce unreasonable sentences in response to OOD questions (Figure 2). Therefore, we propose a novel approach: for each question, we have the finetuned LLM generate an answer, and then we apply an OOD detection criterion specifically designed for the question-answer pair.

Formally, in the context of autoregressive large language models, consider a question $q = q_1, \dots, q_{T_q}$, from which we generate an answer $a = a_1, \dots, a_{T_a}$ by sampling from the conditional distribution $p(\cdot|q)$. We define the following criterion:

$$S_q(q, a) = \frac{p_\theta(q)}{p_{\theta'}(q)}, \quad S_a(q, a) = \frac{p_\theta(a)}{p_{\theta'}(a)}, \quad S_{q,a}(q, a) = \frac{p_\theta(q, a)}{p_{\theta'}(q, a)},$$

$$S_{a|q}(q, a) = \frac{p_\theta(a|q)}{p_{\theta'}(a|q)} = \frac{p_\theta(q, a)p_{\theta'}(q)}{p_{\theta'}(q, a)p_\theta(q)} = \frac{S_{q,a}(q, a)}{S_q(q, a)},$$

¹The infinite monkey theorem states that a monkey randomly typing on a typewriter for an infinitely long period will almost surely produce any given text, such as the complete works of William Shakespeare.

where S_q , S_a , $S_{q,a}$, and $S_{a|q}$ are defined as follows: S_q is the likelihood ratio for the question, S_a for the answer, $S_{q,a}$ for the question-answer pair, and $S_{a|q}$ for the answer given the question. All these criterion are the ratios between the likelihoods assigned by the finetuned model to those assigned by the base model.

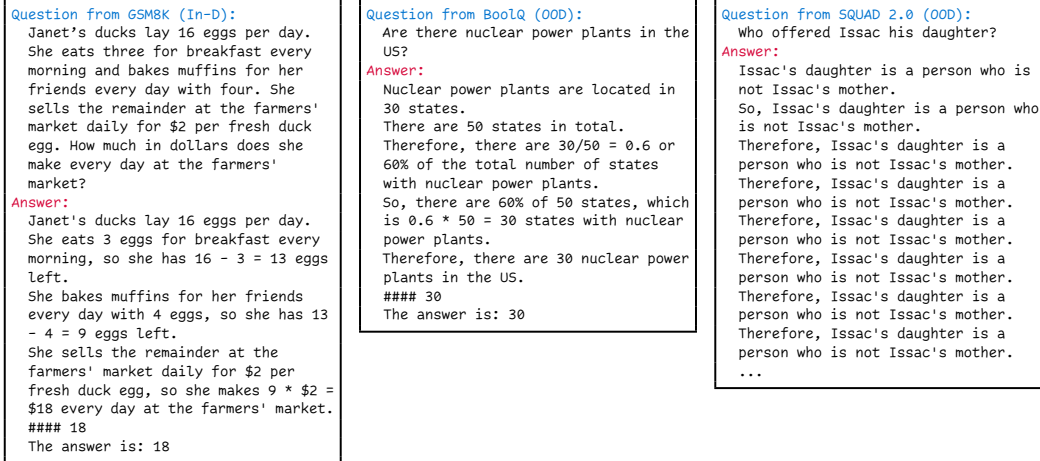


Figure 2: Example question-answer sets produced by MetaMath-7B. The responses to In-D questions are accurate and logical. However, for OOD questions, MetaMath-7B generates unreasonable answers, responding to a straightforward query with unnecessary mathematical calculations or producing repetitive sentences with no useful information.

5 Experiments and Results

In this section, we conduct a comprehensive evaluation across various scenarios, including far OOD, near OOD, spam detection, and QA, to demonstrate the effectiveness of our approach.

We adhere to the definitions of far OOD and near OOD as outlined by Yang et al. (2022) in their work. Near OOD datasets exhibit only a semantic shift from the In-D datasets, whereas far OOD also encompasses a significant covariate (domain) shift. For far OOD evaluations, we designate two distinct datasets as In-D and OOD. For near OOD, we divide a single dataset into two groups: one serving as the In-D with certain classes and the other as OOD with a different set of classes.

Additionally, we demonstrate the capability of our proposed method in detecting OOD instances within the context of spam detection (Labonne & Moran, 2023) - a practical application for our unsupervised OOD detection technique, especially where In-D labels are absent. We show that our method achieves commendable results in spam detection even without access to any spam data. Moreover, when spam data is available and we further finetune the OOD proxy distribution using this data, our results are competitive with state-of-the-art (SOTA) spam detection algorithms.

Finally, we evaluate our approach within a real question-answering (QA) context, utilizing MetaMath - a Llama-2 model finetuned for math problem-solving, as described by Yu et al. (2023). By implementing the likelihood-ratio-based criteria outlined in Section 4, we find that for specific short questions, having the LLM provide an answer and subsequently applying a criterion that analyzes both the question and answer leads to consistently improved outcomes in identifying OOD questions.

Evaluation Metrics We employ AUROC (Area Under the Receiver Operating Characteristic curve), AUPR (Area Under the Precision-Recall curve), and FPR95 (False Positive Rate at 95% True Positive Rate) as our evaluation metrics. These metrics are commonly utilized in

In-D	OOD	Method	AUROC \uparrow	AUPR (OOD) \uparrow	FPR95 \downarrow	
	SST-2	Zhou et al. (2021)	0.978	0.865	0.015	
		CE (Hendrycks & Gimpel, 2016)	0.981	0.942	0.087	
		TAPT (Gururangan et al., 2020)	0.981	0.939	0.088	
		SupCon (Khosla et al., 2020)	0.980	0.943	0.094	
		Uppaal et al. (2023)	1.000	0.999	0.000	
		Llama-7B LH	0.008	0.541	0.999	
		Llama-7B LR	1.000	1.000	0.000	
		Mistral-7B LH	0.008	0.541	1.000	
		Mistral-7B LR	0.995	0.999	0.009	
		Llama-13B LH	0.009	0.541	1.000	
		Llama-13B LR	1.000	1.000	0.000	
		20NG	RTE	Zhou et al. (2021)	0.956	0.860
	CE (Hendrycks & Gimpel, 2016)			0.945	0.902	0.285
	TAPT (Gururangan et al., 2020)			0.919	0.869	0.352
	SupCon (Khosla et al., 2020)			0.952	0.914	0.248
	Uppaal et al. (2023)			1.000	0.999	0.000
	Llama-7B LH			0.063	0.443	0.998
	Llama-7B LR			1.000	1.000	0.001
	Mistral-7B LH			0.074	0.446	0.998
	Mistral-7B LR			0.997	0.999	0.006
	Llama-13B LH			0.070	0.445	0.997
	Llama-13B LR			1.000	1.000	0.000
	IMDB				Zhou et al. (2021)	0.969
		CE (Hendrycks & Gimpel, 2016)	0.961		0.995	0.206
		TAPT (Gururangan et al., 2020)	0.965		0.995	0.159
		SupCon (Khosla et al., 2020)	0.970		0.996	0.150
		Uppaal et al. (2023)	0.990		0.998	0.012
		Llama-7B LH	0.755		0.311	0.932
		Llama-7B LR	1.000		1.000	0.001
		Mistral-7B LH	0.767		0.943	0.926
Mistral-7B LR		0.999	0.998		0.003	
Llama-13B LH		0.773	0.332		0.919	
Llama-13B LR		1.000	1.000		0.000	

Table 1: Results of far OOD detection, utilizing the same experimental setup as described by Uppaal et al. (2023). Results for methods not originating from our work are cited directly from Uppaal et al. (2023). This table presents a condensed version; for the complete details, please see Table 5 in Appendix.

evaluating the performance of OOD detection methods (Hendrycks & Gimpel, 2016; Yang et al., 2022).

5.1 Far OOD Detection

For text data, detecting far OOD instances is relatively less challenging. Uppaal et al. (2023) have illustrated that utilizing the latent distance from a pretrained RoBERTa model significantly addresses far OOD detection, especially when the 20 Newsgroups (20NG) dataset serves as the in-distribution. In this context, we present that our proposed method also attains nearly perfect performance under the same experimental conditions, as detailed in Table 1.

Note that, the notation ‘Model-XB LH/LR’ is used, where ‘Model’ can be either Llama (Touvron et al., 2023) or Mistral (Jiang et al., 2023), denoting two popular open-source large language models (LLMs). Here, ‘Llama’ specifically refers to the Llama 2 version. The ‘XB’ indicates the model’s parameter count, either 7B or 13B. ‘LH’ stands for likelihood, signifying the use of $p_\theta(x)$ as the criterion for OOD detection; ‘LR’ denotes likelihood-ratio, referring to the employment of $S(x)$ as outlined in Equation (1) for OOD identification.

Examining Table 1, we observe that for far OOD detection, the Llama-13B LR model nearly perfectly addresses the challenge across all the In-D OOD pairs for far OOD detection.

5.2 Near OOD Detection

We select the ROSTD (Gangal et al., 2020), SNIPS (Coucke et al., 2018), and CLINC150 (Larson et al., 2019) datasets for our near OOD detection experiments. The ROSTD and

Dataset	In-D Label	Model	AUROC \uparrow	AUPR \uparrow	FPR95 \downarrow
ROSTD	No	Gangal et al. (2020)	0.981	0.958	0.077
		Jin et al. (2022)	0.990	0.973	0.041
		Llama-7B LH	0.960	0.890	0.168
		Llama-7B LR	0.994	0.984	0.023
		Mistral-7B LH	0.964	0.901	0.158
		Mistral-7B LR	0.992	0.978	0.033
		Llama-13B LH	0.965	0.905	0.166
		Llama-13B LR	0.994	0.988	0.018
	Yes	Podolskiy et al. (2021)	0.998	0.994	0.008
	SNIPS	No	Gangal et al. (2020)	0.955	0.903
Jin et al. (2022)			0.963	0.910	0.145
Llama-7B LH			0.912	0.829	0.391
Llama-7B LR			0.993	0.986	0.029
Mistral-7B LH			0.912	0.819	0.417
Mistral-7B LR			0.987	0.968	0.087
Llama-13B LH			0.942	0.872	0.280
Llama-13B LR			0.995	0.988	0.028
Yes		Podolskiy et al. (2021)	0.978	0.933	0.120
CLINC150		No	Gangal et al. (2020)	0.883	0.677
	Jin et al. (2022)		0.902	0.703	0.417
	Llama-7B LH		0.821	0.456	0.538
	Llama-7B LR		0.917	0.766	0.384
	Mistral-7B LH		0.823	0.454	0.540
	Mistral-7B LR		0.913	0.730	0.399
	Llama-13B LH		0.820	0.450	0.546
	Llama-13B LR		0.915	0.742	0.386
	Yes	Podolskiy et al. (2021)	0.982	0.939	0.092

Table 2: Results for near OOD detection. Since the experimental configurations in the studies by Gangal et al. (2020), Podolskiy et al. (2021), and Jin et al. (2022) differ, we have replicated their methods and aligned the dataset splitting for consistency.

CLINC150 datasets are specifically crafted for OOD detection and include designated classes representing OOD data from the same domain. The SNIPS dataset comprises user utterances distributed among seven intent classes, such as GetWeather and RateBook. As it does not inherently provide OOD utterances, we classify the GetWeather and BookRestaurant intents as OOD for the purpose of our experiments. It’s noteworthy that this classification diverges from the one in the study by Jin et al. (2022), which does not explicitly detail their data splitting methodology.

Table 2 demonstrates that the likelihood ratio between the pretrained Llama model and the finetuned Llama model yields the highest performance among unsupervised OOD detection methods. In the case of CLINC150, the supervised OOD detection method introduced by Podolskiy et al. (2021) significantly surpasses our approach, a point that is further discussed in Section 6.

5.3 Spam Detection

Given that the concept of unsupervised OOD detection aligns closely with spam detection, we evaluate our method using the spam detection benchmark introduced by Labonne & Moran (2023). This benchmark includes four specific spam detection datasets: Ling-Spam Dataset (Sakkis et al., 2003), SMS Spam Collection (Almeida et al., 2011), SpamAssassin Public Corpus, and Enron Email Dataset (Metsis et al., 2006). It compares the performance of both traditional and deep learning-based binary classifiers. Our method, being rooted in OOD detection, requires only the non-spam (ham) data for finetuning the LLM. Table 3 indicates that our method, without any spam data, can still reliably identify spam. Furthermore, when spam data is available, we can finetune the OOD proxy using this data and apply the likelihood ratio between the two finetuned LLMs. This approach achieves performance that is competitive with the SOTA in spam detection.

Dataset	Spam Data	Model	AUROC \uparrow	AUPR \uparrow	FPR95 \downarrow
SMS	No	Llama-7B LH	0.960	0.699	0.088
		Llama-7B LR	0.866	0.582	0.487
		Llama-13B LH	0.957	0.689	0.093
		Llama-13B LR	0.810	0.518	0.761
	Yes	NB	0.988	0.949	0.113
		Logistic	0.985	0.946	0.124
		KNN	0.863	0.830	0.811
		SVM	0.997	0.980	0.024
		XGBoost	0.918	0.873	0.676
		LightGBM	0.978	0.921	0.103
		RoBERTa	0.997	0.988	0.004
		Spam-T5	0.985	0.959	0.082
		Llama-7B LR	1.000	1.000	0.000
		Llama-13B LR	0.999	0.995	0.000
SpamAssassin	No	Llama-7B LH	0.964	0.884	0.096
		Llama-7B LR	0.960	0.935	0.296
		Llama-13B LH	0.956	0.897	0.169
		Llama-13B LR	0.941	0.917	0.398
	Yes	NB	0.971	0.917	0.070
		Logistic	0.992	0.986	0.029
		KNN	0.931	0.935	0.578
		SVM	0.990	0.983	0.046
		XGBoost	0.994	0.989	0.019
		LightGBM	1.000	0.999	0.000
		RoBERTa	0.999	0.997	0.000
		Spam-T5	0.996	0.994	0.012
		Llama-7B LR	0.998	0.996	0.005
		Llama-13B LR	0.994	0.989	0.019
Enron	No	Llama-7B LH	0.721	0.728	0.798
		Llama-7B LR	0.991	0.989	0.043
		Llama-13B LH	0.719	0.723	0.798
		Llama-13B LR	0.992	0.990	0.035
	Yes	NB	0.992	0.991	0.035
		Logistic	0.994	0.992	0.025
		KNN	0.915	0.927	0.239
		SVM	0.998	0.998	0.008
		XGBoost	0.975	0.967	0.111
		LightGBM	0.997	0.997	0.013
		RoBERTa	1.000	1.000	0.001
		Spam-T5	1.000	1.000	0.001
		Llama-7B LR	0.999	0.999	0.001
		Llama-13B LR	1.000	1.000	0.000

Table 3: Results of spam detection. This table presents a condensed version; for the complete details, please see Table 6 in Appendix.

5.4 OOD Question Detection in QA Systems

We test the effectiveness of the criterion we introduced in Section 4 within a QA scenario. Here, we employ MetaMath (Yu et al., 2023), which is a Llama2 model finetuned on mathematics data for solving mathematical problems. The objective in this QA context is to identify OOD questions that originate from domains outside the model’s expertise. Filtering out OOD questions enhances the system’s robustness. For this evaluation, we designate the test sets of the GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021) datasets as In-D and use SQUAD (Rajpurkar et al., 2018), BoolQ (Clark et al., 2019), and PIQA (Bisk et al., 2020) as OOD datasets. Table 4 indicates that the criterion $S_{a|q}$ consistently outperforms the other evaluated criteria, with all its AUROC values exceeding 0.5. This demonstrates its effectiveness in detecting OOD. Notably, S_q exhibits subpar performance in most scenarios, underscoring the necessity of our proposed approach that generates an answer and formulates the criterion based on the question-and-answer pair.

In-D	OOD	Model	Criterion	AUROC \uparrow	AUPR (OOD) \uparrow	FPR95 \downarrow
GSM8K	SQUAD 2.0	7B	S_q	0.1116	0.0546	0.9894
			S_a	0.5463	0.0979	0.9947
			$S_{q,a}$	0.5363	0.0959	0.9924
			$S_{a q}$	0.6877	0.1376	0.9704
	BoolQ	7B	S_q	0.0538	0.1618	0.9955
			S_a	0.5045	0.2616	0.9932
			$S_{q,a}$	0.4797	0.2536	0.9879
			$S_{a q}$	0.7156	0.4041	0.9310
	PIQA	7B	S_q	0.9762	0.9779	0.0735
			S_a	0.9612	0.9746	0.0569
			$S_{q,a}$	0.9975	0.9983	0.0038
			$S_{a q}$	0.9944	0.9944	0.0099
MATH	SQUAD 2.0	7B	S_q	0.2139	0.2304	0.9474
			S_a	0.6384	0.3963	0.8916
			$S_{q,a}$	0.6527	0.4477	0.8436
			$S_{a q}$	0.7385	0.5305	0.7914
	BoolQ	7B	S_q	0.1303	0.4472	0.9658
			S_a	0.6135	0.7111	0.8580
			$S_{q,a}$	0.6361	0.7474	0.8008
			$S_{a q}$	0.7507	0.8266	0.6870
	PIQA	7B	S_q	0.9681	0.9902	0.0732
			S_a	0.9206	0.9775	0.1133
			$S_{q,a}$	0.9873	0.9962	0.0242
			$S_{a q}$	0.9876	0.9956	0.0812

Table 4: Outcomes of OOD question detection in QA settings. This table presents a condensed version; for the complete details, please see Table 7 in Appendix.

6 Discussions

Nalisnick’s Paradox in Language OOD Detection Our far OOD detection experiments in Table 1 show that the AUROC for likelihood-based OOD detection methods can be exceptionally low. This finding echoes the phenomenon highlighted by Nalisnick et al. (2018), where language OOD data may unexpectedly exhibit higher likelihood values. Upon examining the characteristics of the datasets involved, it becomes apparent that the texts from the 20 Newsgroups (20NG) (Lang, 1995) dataset are significantly longer than those from the comparative OOD datasets, such as SST-2 and RTE, especially in instances where the AUROC is notably low. This observation reveals a tendency among language models to assign higher likelihoods to shorter sentences, irrespective of their actual semantic content. It highlights the crucial importance of adopting the likelihood ratio rather than solely relying on raw likelihood values to enhance the accuracy of OOD detection.

The Effectiveness of In-D Labels In the near OOD detection experiments presented in Table 2, Podolskiy et al. (2021)’s approach outperforms competing methods significantly. This superior performance may be attributed to the distinct class distribution across datasets. Specifically, the CLINC150 dataset comprises 150 In-D classes, in stark contrast to the SNIPS and ROSTD datasets, which offer a mere 7 (with only 5 for In-D) and 13 In-D classes, respectively. The substantially greater number of In-D classes in CLINC150 compared to the other datasets likely enhances Podolskiy et al. (2021)’s method’s ability to leverage the extensive class label information, thus yielding improved OOD detection results.

7 Related Work

Building upon the OOD detection method using likelihood ratios introduced by Ren et al. (2019), Gangal et al. (2020) suggested employing the likelihood ratio between two LSTM language models. In their approach, one model functions as a “background model” representing OOD data and is trained on random combinations from the vocabulary.

Jin et al. (2022) used the likelihood ratio between a pretrained GPT-2 and a finetuned version of GPT-2 as a baseline to compare with their proposed contrastive learning-based OOD

detection method. They showed that their proposed method can outperform the likelihood ratio based method. However, they only used GPT-2 in their likelihood ratio baseline, which is very small both in terms of training data and the model size by today standard.

In comparison, our study provides a more comprehensive analysis of likelihood ratio between base models and fine-tuned models using much larger LLMs. We show that leveraging these more advanced LLMs, likelihood ratio significantly outperforms results from Jin et al. (2022) in OOD detection. Additionally, in the current era, accessing and sharing both pretrained and finetuned LLMs has become considerably easier. Using likelihood ratio in LLMs for OOD detection is easy, accessible, and very effective.

8 Concluding Remarks

We revisit and validate the likelihood ratio between a pretrained LLM and its finetuned variant as a criterion for OOD detection across various scenarios, without the need for additional training. This LLM-based likelihood ratio, despite being very easy to implement, shows surprising empirical effectiveness in OOD detection, and more importantly, it enables us to build robust QA systems that are able to answer both general and domain-specific questions (i.e., using the likelihood ratio to determine which LLM should be used to answer the input question). We expect that our LLM-based likelihood ratio can benefit many other applications in the future.

References

- Tiago A Almeida, José María G Hidalgo, and Akebo Yamakami. Contributions to the study of sms spam filtering: new collection and results. In *Proceedings of the 11th ACM symposium on Document engineering*, pp. 259–262, 2011. 7
- Christopher M Bishop. Novelty detection and neural network validation. *IEE Proceedings-Vision, Image and Signal processing*, 141(4):217–222, 1994. 2, 3
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020. 8
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021. 1
- Émile Borel. La mécanique statique et l’irréversibilité. *J. Phys. Theor. Appl.*, 3(1):189–196, 1913. 4
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 3
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*, 2019. 8
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021. 8
- Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, et al. Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. *arXiv preprint arXiv:1805.10190*, 2018. 6
- Arthur Eddington. *The nature of the physical world: THE GIFFORD LECTURES 1927*, volume 23. BoD–Books on Demand, 2019. 4

-
- Varun Gangal, Abhinav Arora, Arash Einolghozati, and Sonal Gupta. Likelihood ratios and generative classifiers for unsupervised out-of-domain detection in task oriented dialog. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 7764–7771, 2020. 6, 7, 9
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. Don’t stop pretraining: Adapt language models to domains and tasks. *arXiv preprint arXiv:2004.10964*, 2020. 6, 16
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016. 1, 2, 6, 16
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021. 8
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 13
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023. 6
- Di Jin, Shuyang Gao, Seokhwan Kim, Yang Liu, and Dilek Hakkani-Tür. Towards textual out-of-domain detection without in-domain labels. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:1386–1395, 2022. 7, 9, 10
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020. 6, 16
- Maxime Labonne and Sean Moran. Spam-t5: Benchmarking large language models for few-shot email spam detection. *arXiv preprint arXiv:2304.01238*, 2023. 5, 7
- Ken Lang. Newsweeder: Learning to filter netnews. In *Machine learning proceedings 1995*, pp. 331–339. Elsevier, 1995. 9
- Stefan Larson, Anish Mahendran, Joseph J Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K Kummerfeld, Kevin Leach, Michael A Laurenzano, Lingjia Tang, et al. An evaluation dataset for intent classification and out-of-scope prediction. *arXiv preprint arXiv:1909.02027*, 2019. 6
- Zhihao Lyu, Danier Duolikun, Bowei Dai, Yuan Yao, Pasquale Minervini, Tim Z. Xiao, and Yarin Gal. You need only uncertain answers: Data efficient multilingual question answering. *ICML 2020 Workshop on Uncertainty and Robustness in Deep Learning*, 2020. 1
- Vangelis Metsis, Ion Androutsopoulos, and Georgios Paliouras. Spam filtering with naive bayes-which naive bayes? In *CEAS*, volume 17, pp. 28–69. Mountain View, CA, 2006. 7
- Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do deep generative models know what they don’t know? *arXiv preprint arXiv:1810.09136*, 2018. 2, 3, 4, 9
- Alexander Podolskiy, Dmitry Lipin, Andrey Bout, Ekaterina Artemova, and Irina Piontkovskaya. Revisiting mahalanobis distance for transformer-based out-of-domain detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 13675–13682, 2021. 7, 9
- Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions for squad. *arXiv preprint arXiv:1806.03822*, 2018. 8

-
- Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. *Advances in neural information processing systems*, 32, 2019. 2, 3, 9
- Georgios Sakkis, Ion Androutsopoulos, Georgios Paliouras, Vangelis Karkaletsis, Constantine D Spyropoulos, and Panagiotis Stamatopoulos. A memory-based approach to anti-spam filtering for mailing lists. *Information retrieval*, 6:49–73, 2003. 7
- Robin Schirmer, Yuxuan Zhou, Tonio Ball, and Dan Zhang. Understanding anomaly detection with deep invertible networks through hierarchies of distributions and features. *Advances in Neural Information Processing Systems*, 33:21038–21049, 2020. 2, 3
- Joan Serra, David Álvarez, Vicenç Gómez, Olga Slizovskaia, José F Núñez, and Jordi Luque. Input complexity and out-of-distribution detection with likelihood-based generative models. *arXiv preprint arXiv:1909.11480*, 2019. 2, 3
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 4, 6
- Rheeya Uppaal, Junjie Hu, and Yixuan Li. Is fine-tuning needed? pre-trained language models are near perfect for out-of-domain detection. *arXiv preprint arXiv:2305.13282*, 2023. 6, 16
- Tim Z. Xiao, Aidan N. Gomez, and Yarin Gal. Wat zei je? detecting out-of-distribution translations with variational transformers. *arXiv preprint arXiv:2006.08344*, 2020. 1
- Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi Wang, Guangyao Chen, Bo Li, Yiyu Sun, et al. Openood: Benchmarking generalized out-of-distribution detection. *Advances in Neural Information Processing Systems*, 35: 32598–32611, 2022. 2, 5, 6
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023. 5, 8, 13
- Andi Zhang and Damon Wischik. Falsehoods that ml researchers believe about ood detection. *arXiv preprint arXiv:2210.12767*, 2022. 2, 3
- Mingtian Zhang, Andi Zhang, and Steven McDonagh. On the out-of-distribution generalization of probabilistic image modelling. *Advances in Neural Information Processing Systems*, 34:3811–3823, 2021. 2, 3
- Wenxuan Zhou, Fangyu Liu, and Muhao Chen. Contrastive out-of-distribution detection for pretrained transformers. *arXiv preprint arXiv:2104.08812*, 2021. 6, 16

A Discussion (Continued)

A.1 The cases that LH is better than LR

In the spam detection experiments detailed in Table 3, particularly with data from the SMS and SpamAssassin datasets, we observe that without spam data, the likelihood (LH) method outperforms the likelihood ratio (LR). While LR generally shows superior and more consistent performance across most experiments - as LH can exhibit extremely poor performance in certain cases, a point elaborated in Section 6 - there are specific instances where LH is more effective.

The rationale behind using large language models (LLMs) as an OOD proxy, as introduced in Section 2, is based on the assumption that OOD data deviates from domain-specific natural language content. However, spam messages in the SMS dataset often include intentionally misspelled words to circumvent detection mechanisms, thereby violating our natural language assumption. Similarly, the content from the SpamAssassin dataset, being highly structured in email and data transaction formats (header information), also diverges from typical natural language patterns.

Given these deviations from the natural language assumption, it is understandable why LH might outperform LR in these unique scenarios.

A.2 A Qualitative Analysis on QA OOD Detection

Figure 3 demonstrates sample QA pairs generated by MetaMath, accompanied by their associated criteria: S_q , S_a , $S_{q,a}$, and $S_{a|q}$, which are introduced in Section 4. A higher score for a given criterion suggests that the QA pair is more likely to be OOD. Example (a) is a classic In-D instance, presenting a mathematical problem paired with a logical solution that culminates in a quantitative answer. Example (b) poses a conventional question paired with an irrelevant quantitative answer under the guise of mathematical logic. Example (c) is an undesirable case where MetaMath repeats sentences without providing a substantial answer. Examples (d) and (f) display standard questions with appropriate responses from MetaMath. In example (e), MetaMath attempts to address a general question through mathematical reasoning but does not succeed.

The criterion S_q for example (a) exhibits the highest value among all S_q scores, indicating that S_q may not effectively detect OOD questions. This observation is consistent with our analysis in Section 4 and our discussion in Section 6. As for the most effective score, $S_{a|q}$, we observe that when MetaMath provides an appropriate response to a general question, $S_{a|q}$ is able to effectively flag it as an OOD question, as demonstrated in examples (d) and (f). This supports the hypothesis that a typical answer generated from a conventional question is more likely to receive a higher likelihood in Llama - a model oriented towards natural language - compared to MetaMath, which is specialized in mathematical problems.

B Finetuning Details

For finetuning, with the exception of MetaMath (Yu et al., 2023), which is comprehensively finetuned from the Llama-2 model, we employ LoRA (Low-Rank Adaptation) as described by Hu et al. (2021) for adjusting the models. Our aim is to demonstrate the general applicability of our method; therefore, all our LoRA finetuning follows a uniform parameter configuration. In accordance with prevalent practices, we finetune the parameters in the Q (query) and V (value) projections. We employ a LoRA rank of 16, a LoRA alpha of 32, and a LoRA dropout rate of 0.05. Given the varying sizes of datasets, we adjust the learning rate between 10^{-3} and 10^{-4} , and the training typically spans 10 to 20 epochs. We select the model checkpoint with the best evaluation loss as our finetuned model. For detailed information, please refer to the code provided.

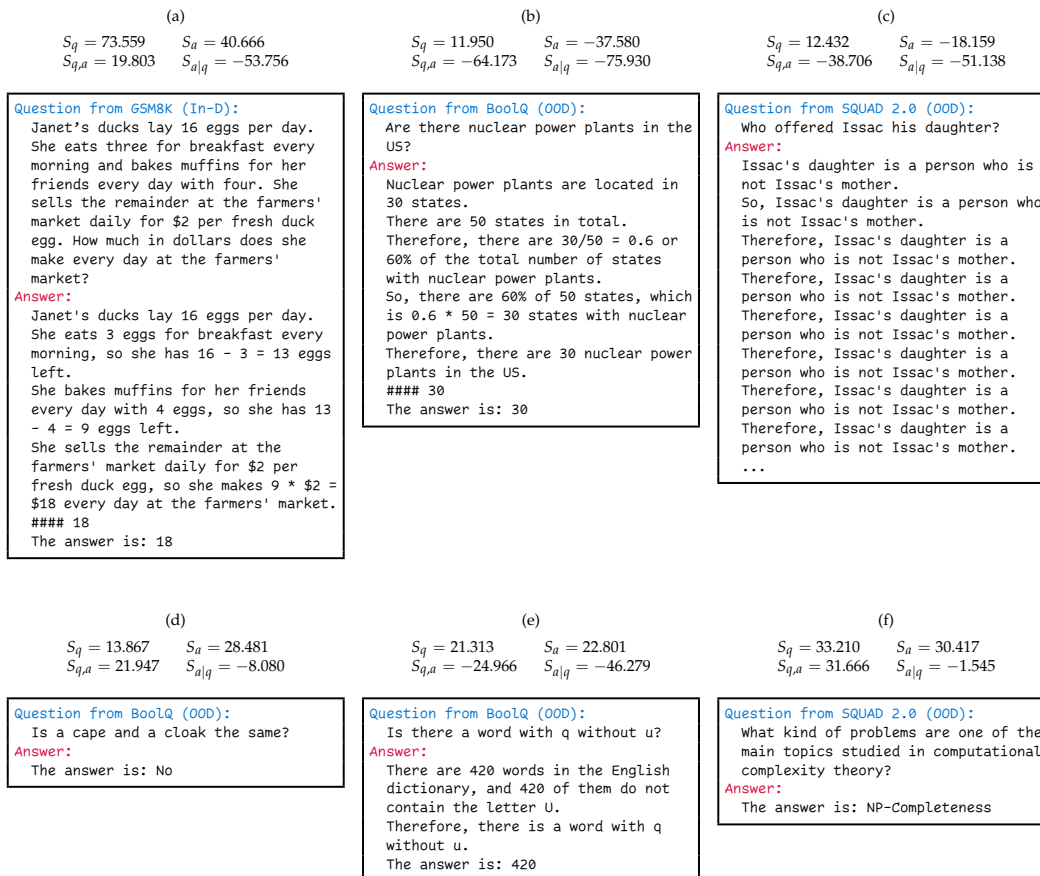


Figure 3: Example question-answer sets produced by MetaMath-7B, with the applicable criteria displayed above each pair. The criteria are represented on a logarithmic scale.

C Full Experiment Results

Full experiment results are detailed below.

In-D	OOD	Method	AUROC \uparrow	AUPR (OOD) \uparrow	FPR95 \downarrow
	SST-2	Zhou et al. (2021)	0.978	0.865	0.015
		CE (Hendrycks & Gimpel, 2016)	0.981	0.942	0.087
		TAPT (Gururangan et al., 2020)	0.981	0.939	0.088
		SupCon (Khosla et al., 2020)	0.980	0.943	0.094
		Uppaal et al. (2023)	1.000	0.999	0.000
		Llama-7B LH	0.008	0.541	0.999
		Llama-7B LR	1.000	1.000	0.000
		Mistral-7B LH	0.008	0.541	1.000
		Mistral-7B LR	0.995	0.999	0.009
		Llama-13B LH	0.009	0.541	1.000
		Llama-13B LR	1.000	1.000	0.000
			MNLI	Zhou et al. (2021)	0.964
CE (Hendrycks & Gimpel, 2016)	0.968			0.989	0.166
TAPT (Gururangan et al., 2020)	0.964			0.988	0.175
SupCon (Khosla et al., 2020)	0.970			0.990	0.156
Uppaal et al. (2023)	1.000			1.000	0.000
Llama-7B LH	0.020			0.119	0.999
Llama-7B LR	1.000			1.000	0.001
Mistral-7B LH	0.024			0.119	0.999
Mistral-7B LR	0.996			0.996	0.008
Llama-13B LH	0.024			0.119	0.998
Llama-13B LR	1.000			1.000	0.000
	RTE			Zhou et al. (2021)	0.956
		CE (Hendrycks & Gimpel, 2016)	0.945	0.902	0.285
		TAPT (Gururangan et al., 2020)	0.919	0.869	0.352
		SupCon (Khosla et al., 2020)	0.952	0.914	0.248
		Uppaal et al. (2023)	1.000	0.999	0.000
		Llama-7B LH	0.063	0.443	0.998
		Llama-7B LR	1.000	1.000	0.001
		Mistral-7B LH	0.074	0.446	0.998
		Mistral-7B LR	0.997	0.999	0.006
		Llama-13B LH	0.070	0.445	0.997
		Llama-13B LR	1.000	1.000	0.000
		20NG	IMDB	Zhou et al. (2021)	0.969
CE (Hendrycks & Gimpel, 2016)	0.961			0.995	0.206
TAPT (Gururangan et al., 2020)	0.965			0.995	0.159
SupCon (Khosla et al., 2020)	0.970			0.996	0.150
Uppaal et al. (2023)	0.990			0.998	0.012
Llama-7B LH	0.755			0.311	0.932
Llama-7B LR	1.000			1.000	0.001
Mistral-7B LH	0.767			0.943	0.926
Mistral-7B LR	0.999			0.998	0.003
Llama-13B LH	0.773			0.332	0.919
Llama-13B LR	1.000			1.000	0.000
	Multi30K			Zhou et al. (2021)	0.980
		CE (Hendrycks & Gimpel, 2016)	0.962	0.920	0.175
		TAPT (Gururangan et al., 2020)	0.956	0.922	0.167
		SupCon (Khosla et al., 2020)	0.955	0.918	0.201
		Uppaal et al. (2023)	1.000	1.000	0.000
		Llama-7B LH	0.002	0.470	1.000
		Llama-7B LR	1.000	1.000	0.000
		Mistral-7B LH	0.002	0.470	1.000
		Mistral-7B LR	0.995	0.998	0.008
		Llama-13B LH	0.002	0.470	1.000
		Llama-13B LR	1.000	1.000	0.000
			NewsCategory	Zhou et al. (2021)	0.955
CE (Hendrycks & Gimpel, 2016)	0.957			0.984	0.234
TAPT (Gururangan et al., 2020)	0.947			0.981	0.243
SupCon (Khosla et al., 2020)	0.962			0.986	0.219
Uppaal et al. (2023)	1.000			1.000	0.000
Llama-7B LH	0.014			0.128	1.000
Llama-7B LR	1.000			1.000	0.001
Mistral-7B LH	0.019			0.129	1.000
Mistral-7B LR	0.997			0.997	0.006
Llama-13B LH	0.017			0.128	0.999
Llama-13B LR	1.000			1.000	0.000
	CLINC150			Zhou et al. (2021)	0.988
		CE (Hendrycks & Gimpel, 2016)	0.964	0.844	0.189
		TAPT (Gururangan et al., 2020)	0.959	0.830	0.213
		SupCon (Khosla et al., 2020)	0.957	0.821	0.230
		Uppaal et al. (2023)	1.000	1.000	0.000
		Llama-7B LH	0.001	0.661	1.000
		Llama-7B LR	1.000	1.000	0.000
		Mistral-7B LH	0.001	0.661	1.000
		Mistral-7B LR	0.995	1.000	0.008
		Llama-13B LH	0.001	0.661	1.000
		Llama-13B LR	1.000	1.000	0.000

Table 5: Results of far OOD detection, utilizing the same experimental setup as described by Uppaal et al. (2023). Results for methods not originating from our work are cited directly from Uppaal et al. (2023).

Dataset	Spam Data	Model	AUROC \uparrow	AUPR \uparrow	FPR95 \downarrow
Ling	No	Llama-7B LH	0.552	0.215	0.934
		Llama-7B LR	0.967	0.858	0.174
		Llama-13B LH	0.534	0.182	0.929
		Llama-13B LR	0.933	0.746	0.336
	Yes	NB	1.000	1.000	0.000
		Logistic	1.000	1.000	0.000
		KNN	0.968	0.932	0.021
		SVM	1.000	1.000	0.000
		XGBoost	0.995	0.973	0.017
		LightGBM	0.997	0.979	0.008
		RoBERTa	1.000	1.000	0.000
		Spam-T5	1.000	1.000	0.000
		Llama-7B LR	0.998	0.993	0.008
		Llama-13B LR	0.997	0.988	0.012
SMS	No	Llama-7B LH	0.960	0.699	0.088
		Llama-7B LR	0.866	0.582	0.487
		Llama-13B LH	0.957	0.689	0.093
		Llama-13B LR	0.810	0.518	0.761
	Yes	NB	0.988	0.949	0.113
		Logistic	0.985	0.946	0.124
		KNN	0.863	0.830	0.811
		SVM	0.997	0.980	0.024
		XGBoost	0.918	0.873	0.676
		LightGBM	0.978	0.921	0.103
		RoBERTa	0.997	0.988	0.004
		Spam-T5	0.985	0.959	0.082
		Llama-7B LR	1.000	1.000	0.000
		Llama-13B LR	0.999	0.995	0.000
SpamAssassin	No	Llama-7B LH	0.964	0.884	0.096
		Llama-7B LR	0.960	0.935	0.296
		Llama-13B LH	0.956	0.897	0.169
		Llama-13B LR	0.941	0.917	0.398
	Yes	NB	0.971	0.917	0.070
		Logistic	0.992	0.986	0.029
		KNN	0.931	0.935	0.578
		SVM	0.990	0.983	0.046
		XGBoost	0.994	0.989	0.019
		LightGBM	1.000	0.999	0.000
		RoBERTa	0.999	0.997	0.000
		Spam-T5	0.996	0.994	0.012
		Llama-7B LR	0.998	0.996	0.005
		Llama-13B LR	0.994	0.989	0.019
Enron	No	Llama-7B LH	0.721	0.728	0.798
		Llama-7B LR	0.991	0.989	0.043
		Llama-13B LH	0.719	0.723	0.798
		Llama-13B LR	0.992	0.990	0.035
	Yes	NB	0.992	0.991	0.035
		Logistic	0.994	0.992	0.025
		KNN	0.915	0.927	0.239
		SVM	0.998	0.998	0.008
		XGBoost	0.975	0.967	0.111
		LightGBM	0.997	0.997	0.013
		RoBERTa	1.000	1.000	0.001
		Spam-T5	1.000	1.000	0.001
		Llama-7B LR	0.999	0.999	0.001
		Llama-13B LR	1.000	1.000	0.000

Table 6: Results of spam detection.

In-D	OOD	Model	Criterion	AUROC \uparrow	AUPR (OOD) \uparrow	FPR95 \downarrow
GSM8K	SQUAD 2.0	7B	S_q	0.1116	0.0546	0.9894
			S_a	0.5463	0.0979	0.9947
			$S_{q,a}$	0.5363	0.0959	0.9924
			$S_{a q}$	0.6877	0.1376	0.9704
		13B	S_q	0.0519	0.0524	0.9955
			S_a	0.4958	0.0891	0.9955
			$S_{q,a}$	0.4017	0.0764	0.9947
			$S_{a q}$	0.6144	0.1136	0.9765
	BoolQ	7B	S_q	0.0538	0.1618	0.9955
			S_a	0.5045	0.2616	0.9932
			$S_{q,a}$	0.4797	0.2536	0.9879
			$S_{a q}$	0.7156	0.4041	0.9310
		13B	S_q	0.1008	0.1659	0.9924
			S_a	0.5507	0.2861	0.9833
			$S_{q,a}$	0.4778	0.2530	0.9795
			$S_{a q}$	0.6967	0.3886	0.9303
PIQA	7B	S_q	0.9762	0.9779	0.0735	
		S_a	0.9612	0.9746	0.0569	
		$S_{q,a}$	0.9975	0.9983	0.0038	
		$S_{a q}$	0.9944	0.9944	0.0099	
	13B	S_q	0.9900	0.9906	0.0318	
		S_a	0.8879	0.9331	0.1122	
		$S_{q,a}$	1.0000	0.9999	0.0000	
		$S_{a q}$	1.0000	1.0000	0.0000	
MATH	SQUAD 2.0	7B	S_q	0.2139	0.2304	0.9474
			S_a	0.6384	0.3963	0.8916
			$S_{q,a}$	0.6527	0.4477	0.8436
			$S_{a q}$	0.7385	0.5305	0.7914
		13B	S_q	0.1880	0.2232	0.9460
			S_a	0.6098	0.3841	0.8890
			$S_{q,a}$	0.5628	0.3649	0.8882
			$S_{a q}$	0.6786	0.4737	0.8166
	BoolQ	7B	S_q	0.1303	0.4472	0.9658
			S_a	0.6135	0.7111	0.8580
			$S_{q,a}$	0.6361	0.7474	0.8008
			$S_{a q}$	0.7507	0.8266	0.6870
		13B	S_q	0.2612	0.5223	0.9304
			S_a	0.6488	0.7474	0.8148
			$S_{q,a}$	0.6384	0.7521	0.7818
			$S_{a q}$	0.7350	0.8191	0.6854
PIQA	7B	S_q	0.9681	0.9902	0.0732	
		S_a	0.9206	0.9775	0.1133	
		$S_{q,a}$	0.9873	0.9962	0.0242	
		$S_{a q}$	0.9876	0.9956	0.0812	
	13B	S_q	0.9795	0.9938	0.0452	
		S_a	0.8495	0.9572	0.1627	
		$S_{q,a}$	0.9897	0.9969	0.0198	
		$S_{a q}$	0.9626	0.9895	0.0484	

Table 7: Outcomes of OOD question detection in QA settings.