# Firewall Routing: Blocking Leads to Better Hybrid Inference for LLMs

**Anonymous ACL submission**

## Abstract

The rapid advancement of Large Language Models (LLMs) has significantly enhanced performance across various natural language processing (NLP) tasks, yet the high computational costs and latency associated with deploying such models continue to pose critical bottlenecks, limiting their broader applicability. To mitigate these challenges, we propose a dynamic hybrid inference framework, **Firewall Routing**, which efficiently selects between a strong and a weak LLMs based on the complexity of the query. A lightweight routing model is trained to optimize resource allocation by learning from response quality and preventing long-tail queries, which are often unsolvable for both LLMs, from being routed to the stronger model. Moreover, our method incorporates multiple sampling to enhance query evaluation reliability while leveraging **Hard Blocking** and **Soft Blocking** to handle long-tail queries along with refining labels for model selection. Extensive experiments show our method outperforms existing routing strategies by up to 5.29% in APGR, demonstrating state-of-the-art performance across multiple benchmarks.

## 1 Introduction

In recent years, we have witnessed the rapid advancement of artificial intelligence technologies, particularly the rise of large language models (LLMs) such as ChatGPT, which are reshaping the paradigms of our daily work. These models, often containing billions or even trillions of parameters, generate fluent and contextually appropriate responses, enabling natural interactions without requiring specialized user knowledge (OpenAI et al., 2024; Touvron et al., 2023; Grattafiori et al., 2024). However, such remarkable capabilities come at a significant cost: deploying LLMs demands expensive infrastructure, such as multi-GPU systems with high memory capacity, or incurs higher per-token charges in cloud-based LLM services for more capable models (Yu et al., 2022). Moreover, larger models often introduce higher latency, making them less suitable for real-time or resource-constrained applications. Striking a balance among strong model performance, high efficiency, and economical costs remains an "impossible triangle," yet it is precisely this challenge that drives ongoing research efforts in the field.

Making the "impossible triangle" possible requires a paradigm shift in how we allocate computational resources for language model inference. Extensive experiments have demonstrated that not all tasks require the full power of the largest models (Grattafiori et al., 2024). Simpler queries can often be handled effectively by smaller, lower-cost models without compromising quality, whereas more complex queries leverage the advanced capabilities of larger models. This principle forms the foundation of **Hybrid Inference**.

Given the promising potential, **Hybrid Inference** has garnered significant attention from both academia and industry. Existing strategies can be broadly categorized into two main types: **Cascade** methods (Chen et al., 2023; Gupta et al., 2024; Ramírez et al., 2024), and **Route** methods (Shnitzer et al., 2023; Šakota et al., 2024; Lu et al., 2023; Ong et al., 2024; Ding et al., 2024).

**Cascade** methods first process all queries using a weaker model. If the weaker model's confidence in its response is low, typically determined through an internal evaluation mechanism, the query is escalated to a stronger model for reprocessing. Although this approach is conceptually straightforward, it has several inherent limitations. On the one hand, evaluating response quality before completion in generative tasks is inherently difficult, leading to unreliable decision-making(Gupta et al., 2024). On the other hand, evaluating response quality after completion brings greatly increased latency. These factors make **Cascade** methods less
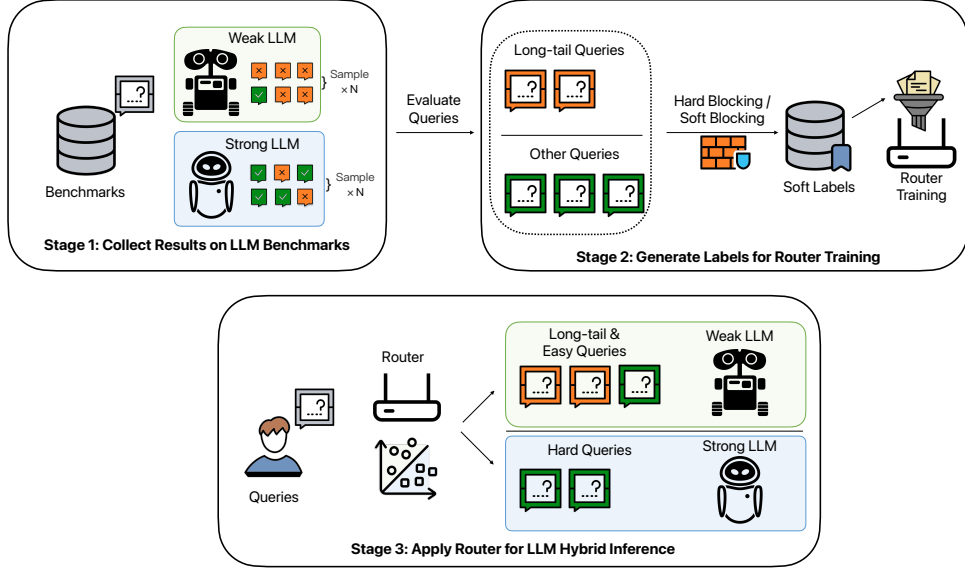
Figure 1: **Firewall Routing** framework for dual-model hybrid inference, comprising a strong model, a weak model, and a router model to balance performance and cost for LLM inference. By blocking unsolvable long-tail queries from being routed to the strong model, the framework achieves state-of-the-art performance.

efficient in real-world applications.

Motivated by these considerations, we focus on **Route** methods, which leverage a lightweight router model to dynamically allocate queries to the most appropriate LLM under a given configuration. However, existing **Route** methods predominantly rely on collected preference data, which are often limited by strict domain-specific constraints (Shnitzer et al., 2023; Šakota et al., 2024; Lu et al., 2023), or heavily depend on model-generated scores (Ong et al., 2024; Ding et al., 2024). Moreover, these methods often depend on preference data or artificially generated labels based on model scoring. In the context of dual-model hybrid inference, where the strong model generally outperforms the weak model, they fail to address long-tail queries that challenge both models, highlighting opportunities for further optimization.

To address these challenges, we propose **Firewall Routing**, a dual-model hybrid inference system that builds on reliable benchmark results and manages to block long-tail queries, enhancing both performance and efficiency.

Specifically, we propose a novel paradigm for training the router model. Unlike existing methods, our approach utilizes multiple sampling during benchmark evaluations to obtain more accurate estimations of the capabilities of both the strong and weak models. These estimations are then used to construct soft labels for router training. Through mathematical derivations, this paradigm highlights the generality of soft label training in the domain of router optimization and demonstrates that the hard label approach is a specific instance of this broader framework.

To further address the challenge of long-tail queries, we propose two novel approaches—**Hard Blocking** and **Soft Blocking**—designed to effectively manage unsolvable cases. **Hard Blocking** utilizes statistical information to identify unsolvable long-tail queries and assigns them the label "route to the weak model," minimizing unnecessary computational overhead. In contrast, **Soft Blocking** leverages the **Pass Rate** (pass@1) to generate refined soft labels with more precise routing conditions, further reducing computational inefficiencies.

To summarize, we make the following contributions:

1. We propose a novel router training paradigm leveraging multiple sampling to generate soft labels, which generalizes router optimization and demonstrates hard label training as a specific case within this framework.

2. We propose **Hard Blocking** and **Soft Blocking** as automated mechanisms to enable our approach to overcome the challenges associated with long-tail queries.

3. We validate our approach through extensive experiments across diverse configurations.

2

## 2 Related Works

**Hybrid Inference** balances response quality and inference cost by dynamically selecting models based on task complexity. For image classification, Kag et al. (2023) explored joint training of a small model, a large model, and a router, while in NLP tasks, the Tryage architecture (Hari and Thomson, 2023) employed a joint-trained router to optimize performance across domains. However, for LLMs, joint training is computationally expensive and deviates from the pre-training paradigm, leading to two main approaches: **Cascade Methods** and **Route Methods**.

**Cascade Methods** first query a weaker model and escalate the request to a stronger model only when necessary. FrugalGPT (Chen et al., 2023) estimates response confidence using an LLM-based heuristic to decide whether a query should be forwarded to a larger model. Similarly, Gupta et al. (2024) proposed a confidence estimation method based on the conditional probability of the generated response, serving as a reliability metric. By assessing the correctness of the weaker model's responses, these methods effectively reduce the number of strong model invocations while maintaining high response quality. However, this approach introduces significant response time overhead, as the weaker model must first generate an output before determining whether escalation is required.

Margin Sampling (Ramírez et al., 2024) is a different cascade approach without introducing extra response time. Only when the probability difference between the top two predicted tokens is small at the beginning of generation, indicating uncertainty, is the query escalated to the strong model.

While both cascade hybrid inference and speculative decoding involve a weak and strong model processing the same query, their goals differ. Speculative decoding (Kim et al., 2024; Leviathan et al., 2023) speeds up text generation by having a smaller model propose tokens, which are then verified by a larger model, but this frequent validation incurs high computational costs. In contrast, cascade methods prioritize reducing reliance on the strong model, balancing performance and efficiency by minimizing its usage.

**Route Methods** introduce a router model to determine which model should handle a given query. For example, Shnitzer et al. (2023) frame this as an out-of-distribution (OOD) detection problem, where they predict a model's response correctness and confidence using k-nearest embedded queries. Similarly, Šakota et al. (2024) train a model to predict whether a query can be correctly answered, incorporating a special token to indicate which LLM should be used. Lu et al. (2023) distilled a reward model to predict which LLM serves as the optimal expert for a given query.

Many recent works focus on dual-model hybrid inference systems. For instance, RouteLLM (Ong et al., 2024) uses preference pairs from multiple LLMs in Chatbot Arena to train a Bradley-Terry model (Bradley and Terry, 1952) as the router. Hybrid LLM (Ding et al., 2024) derives Win Rates for queries through a biased comparison of response BARTScores, creating a desired label distribution to train the router. These approaches highlight the potential for training routers with more reliable evidence, such as pass@k (Chen et al., 2021), to improve model selection.

## 3 Method

### 3.1 Router Training Criterions

#### 3.1.1 Train with Hard Label

Early works on building up hybrid inference systems usually train a system with the router model as a whole, where the router model learns how to route under a fixed configuration(Kag et al., 2023). Due to the high training costs associated with large-scale models, most works in LLM hybrid inference only train the router model.

In existing evaluation frameworks for large language models, generative tasks typically follow a greedy decoding paradigm, where the model outputs the token with the highest probability while disregarding alternative token possibilities. Based on this setting, existing methods (Ding et al., 2024) adopt a "Hard Label" approach for router training.

Specifically, for a single query $x_i \in Q$, let $S(x_i)$ and $W(x_i)$ represent the responses generated by the strong model $S$ and the weak model $W$, respectively, using greedy decoding. The correctness of these responses is denoted as $\delta(S(x_i))$ and $\delta(W(x_i))$, where $\delta(\cdot) \in \{0, 1\}$, with 1 indicating a correct response and 0 indicating an incorrect one. The decision on *whether to route the query to the weak model* is determined by the label $y_i$, defined as $y_i := \mathbb{I}[\delta(S(x_i)) \leq \delta(W(x_i)]$. Here, $y_i = 1$ implies the weak model is capable of performing at least as well as the strong model for query $x_i$, and thus the query should be routed to the weak model.

The hard-label router is trained by minimizing the binary cross-entropy loss:

$$\mathcal{L}(\theta) = -\frac{1}{|Q|}\sum_{i=1}^{|Q|}((1-y_i)log(1-p_\theta(x_i)) \\ + y_i log(p_\theta(x_i))), \qquad (1)$$

where $p_\theta(x)$ is output of router $\theta$ toward query $x$, where larger $p_\theta(x)$ indicates that the queries should more likely to be routed to the weak model.

The hard label approach is limited by its inability to account for the inherent variability in the responses of large models, thereby restricting the router's ability to make fine-grained decisions. This limitation becomes particularly apparent in scenarios where the smaller model's performance is often comparable to that of the larger model. An ideal hybrid inference system should incorporate the inherent variability of model responses into the training labels for the router, enabling it to make more reliable and cost-effective routing decisions, thereby enhancing overall system efficiency.

### 3.1.2 Train with Soft Label

To more objectively reflect the performance of large models, existing evaluations often involve multiple sampling of model outputs. Inspired by this approach, we extend our approach by incorporating multiple sampling, which allows us to evaluate the models more thoroughly and account for response variability. This enhancement aims to improve the robustness and efficiency of the routing decisions in our hybrid inference framework.

Specifically, for a single query $x_i$ , let $S^1(x_i), \ldots, S^n(x_i)$ and $W^1(x_i), \ldots, W^n(x_i)$ denote the responses generated by the strong model $S$ and the weak model $W$ over $n$ sampling iterations. The correctness of these responses is represented by $\delta(S^j(x_i))$ and $\delta(W^j(x_i))$, where $\delta(\cdot) \in \{0,1\}$, with 1 indicating a correct response and 0 indicating an incorrect one. Each sampling iteration produces a noisy observation of $y_i$, denoted as $y_i^j = \mathbb{I}[\delta(S^j(x_i)) \leq \delta(W^j(x_i))]$. In this setting, $x_i$ is associated with $n$ data pairs in the training set, denoted as $(x_i, y_i^1), (x_i, y_i^2), \ldots, (x_i, y_i^n)$.

Using this data, the router can still be trained with a hard label-based objective. However, this approach presents two significant challenges: first, the training cost scales proportionally with the number of sampling attempts $n$; second, a single input

can correspond to varying labels, potentially misleading the router's behavior.

Thus, we introduce the concept of the weak-to-strong **Win Rate**, defined as $r_i := \frac{1}{n}\sum_{j=1}^{n} y_i^j$, which represents the probability that the weak model matches or exceeds the performance of the strong model. Furthermore, we demonstrate that optimization objectives based on **Win Rate** exhibit greater generality for router training. Notably, hard label training inherently captures the concept of **Win Rate**, which can be expressed in the following form:

$$\mathcal{L}(\theta) = -\frac{1}{n|Q|}\sum_{i=1}^{|Q|}\sum_{j=1}^{n}((1-y_i^j)log(1-p_\theta(x_i)) \\ + y_i^j log(p_\theta(x_i)))$$

$$= -\frac{1}{n|Q|}\sum_{i=1}^{|Q|}((n-\sum_{j=1}^{n} y_i^j)log(1-p_\theta(x_i)) \\ + (\sum_{j=1}^{n} y_i^j)log(p_\theta(x_i)))$$

$$= -\frac{1}{|Q|}\sum_{i=1}^{|Q|}((1-r_i)log(1-p_\theta(x_i)) \\ + r_i log(p_\theta(x_i))). \qquad (2)$$

Here, $p_\theta(x)$ represents the output of the router $\theta$ for the query $x$, where a larger $p_\theta(x)$ indicates a higher likelihood that the query should be routed to the weak model. This formula motivates us to explore more refined soft labels that better characterize the behavior of large models through their win rates.

### 3.2 Blocking Long-tail Queries

Even for large models, there are instances where, despite multiple sampling attempts $n$, the model is still unable to resolve certain long-tail queries. This limitation arises from the inherent complexity and ambiguity in some queries, which even powerful models may struggle to address consistently, regardless of the number of samples taken. Consequently, such cases highlight the need for more sophisticated handling of unsolvable long-tail queries in hybrid inference systems.

### 3.2.1 Hard Blocking

To automatically identify long-tail queries, we introduce multiple sample **Pass Rate** (pass@k when k=1) from Chen et al. (2021)'s work to substitute single sample correctness. For a single query $x_i \in Q$ with $n$ sampled responses $R^1(x_i), ..., R^n(x_i)$
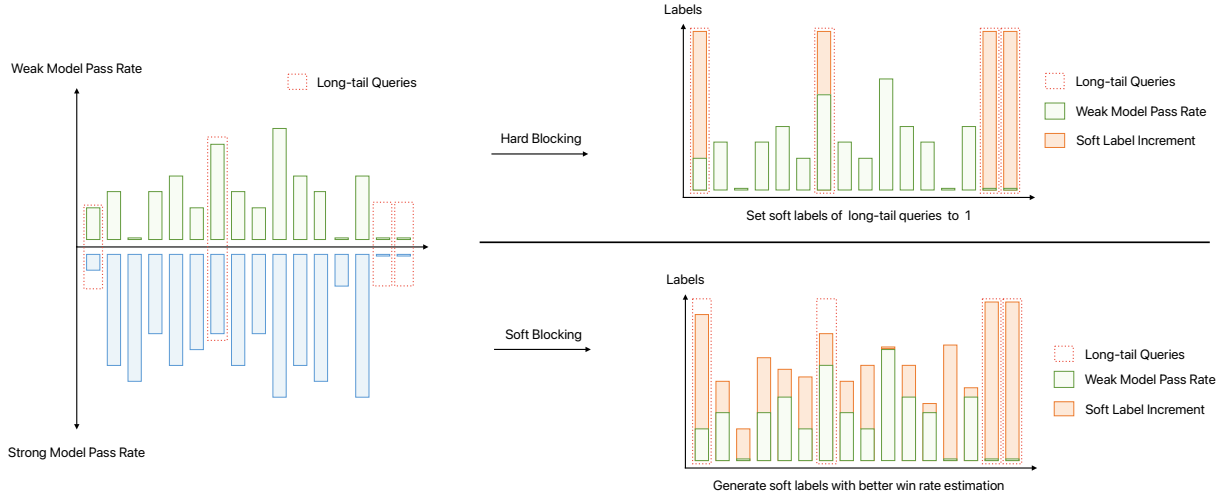
4

Figure 2: Hard Blocking and Soft Blocking facilitate the automatic handling of unsolvable long-tail queries by generating reliable soft labels for router training. Queries assigned larger soft label values are more likely to be routed to the weak model.

from model $R$, **Pass Rate** is defined as the average correctness of these responses:

$$pr(x_i) := \frac{1}{n} \sum_{j=1}^{n} \delta(R^j(x_i)). \quad (3)$$

We are able to split queries into two sets, $Q_u$ and $Q_s = Q - Q_u$, representing unsolvable and solvable queries, satisfying:

$$\forall x^u \in Q_u, pr_s(x^u) \leq pr_w(x^u),$$
$$\forall x^s \in Q_s, pr_s(x^s) > pr_w(x^s). \quad (4)$$

By addressing unsolvable long-tail queries through routing them to the weak model, the decision to route other queries similarly hinges entirely on the weak model's capability to handle these queries effectively:

$$label_i = \begin{cases} pr_w(x_i), x_i \in Q^s, \\ 1, x_i \in Q^u, \end{cases} \quad (5)$$

where $label_i$ is the soft label used in router training to substitute $r_i$ in Eq.2.

To further reduce the cost associated with label collection in this method, it is also possible to split $Q_u$ and $Q_s$ using only the strong model's greedy-decoding responses, subject to the following restrictions:

$$\forall x^u \in Q_u, \delta(S(x^u)) = 0,$$
$$\forall x^s \in Q_s, \delta(S(x^s)) = 1. \quad (6)$$

### 3.2.2  Soft Blocking

A closer examination of Eq.2 and the concept of the Pass Rate reveals that $r_i$ functions as a noisy indicator, capturing the behaviors of the two models when processing the same query. A key insight is that the performance of the strong model is independent of whether the weak model answers correctly. Instead of treating the two models' performances as a joint distribution, we can more effectively leverage the distributional information obtained from multiple samplings. By treating the two independent events separately, we can more accurately estimate $r_i$ through **Pass Rate**. To maximize the use of this information, we define the joint event for routing the query to the weak model by combining two conditions: *the weak model is correct* and *even if the weak model is incorrect, the strong model also fails*. This method allows us to offer a more refined and informative estimate of overall performance:

$$label_i = pr_w(x_i) + (1 - pr_w(x_i))(1 - pr_s(x_i))$$
$$= 1 - (1 - pr_w(x_i))pr_s(x_i), \quad (7)$$

where $label_i$ is the soft label used in router training to substitute $r_i$ in Eq.2, and $label_i$ is the observed frequency that the strong model fail to overperform the weak model.

## 4  Experiments

### 4.1  Settings

**Datasets**   In this study, we conduct experiments on generative tasks commonly used to assess the capabilities of large language models (LLMs). These

| Datasets | TriviaQA | | | | GSM8K | | | | HumanEval | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | APGR↑ | Pass Rate↑ | | | APGR↑ | Pass Rate↑ | | | APGR↑ | Pass Rate↑ | | |
| | | 20% | 50% | 80% | | 20% | 50% | 80% | | 20% | 50% | 80% |
| Linear Interpolation | 50.00 | 19.95 | 34.97 | 49.99 | 50.00 | 11.69 | 17.16 | 22.62 | 50.00 | 8.12 | 10.10 | 12.08 |
| Hybrid LLM | 49.17 | 18.98 | 34.38 | 49.99 | 62.08 | 14.38 | 20.75 | 24.79 | 51.94 | 8.10 | 10.50 | **12.35** |
| RouteLLM (MF) | 51.58 | 20.69 | 36.27 | 51.09 | 49.39 | 11.37 | 17.13 | 22.27 | 47.08 | 7.81 | 9.95 | 12.23 |
| Margin Sampling | 50.02 | 19.78 | 35.01 | 50.15 | 46.01 | 10.85 | 16.02 | 21.70 | 44.88 | 7.74 | 9.81 | 11.53 |
| Ours (Hard Block) | 53.16 | 22.09 | 37.85 | 50.96 | **67.37** | **16.34** | **22.46** | 24.67 | **54.36** | 8.17 | **10.77** | 12.27 |
| Ours (Soft Block) | **55.00** | **22.48** | **38.99** | **52.88** | 66.65 | 15.53 | 22.15 | **25.32** | 53.13 | **8.23** | 10.69 | 12.27 |

Table 1: Zero-shot performance of different methods across selected datasets. The weak model is Llama3.2-1B, and the strong model is Llama3.1-70B. Linear Interpolation represents the combined performance of the two LLMs to simulate random routing. **Bolded values** indicate the best-evaluated results.

tasks include TriviaQA (Joshi et al., 2017) for common sense question answering (QA), GSM8K (Cobbe et al., 2021) for mathematical reasoning, and HumanEval (Chen et al., 2021) for code synthesis. Our training set is derived from the TriviaQA and GSM8K training datasets, while HumanEval is utilized exclusively for evaluation, serving as an Out-Of-Domain task benchmark.

**Prompts** For all datasets, we employ a straightforward zero-shot prompt format without using any system prompts. Specifically, each query is formatted as: "Question: {question}\nAnswer: ".

**Models** In this study, we utilize two large language models (LLMs) from the Llama family (Grattafiori et al., 2024) for our experiments: Llama3.2-1B serves as the weak model, while Llama3.1-70B is employed as the strong model for training the router. Furthermore, to assess the generalizability of the trained router, we test it on an alternative model pair, substituting Llama3.2-3B as the weak model, to evaluate its adaptability to varying model configurations.

**Routers** Aligned with prior studies, we adopt DeBERTa-v3-large (He et al., 2023) as the router model, enhanced with an additional linear layer to predict the probability of routing queries to either the weak or strong model. The router is trained for 10 epochs using the specified loss function, with the best-performing checkpoints selected based on validation set performance for the final evaluation. All experiments are conducted on 8 NVIDIA A100 GPUs, each with 80GB of memory, to facilitate data parallelization; however, the same experiments can be reproduced on a single NVIDIA A100 GPU. The training code will be made publicly available.

**Baselines** We compare our approach with several state-of-the-art methods, including Hybrid LLM

(Ding et al., 2024), RouteLLM (Ong et al., 2024), and Margin Sampling (Ramírez et al., 2024). For Hybrid LLM, we reproduce the methodology and hyperparameter selection as outlined in the original paper. For RouteLLM, we employ the best practices and downloadable pre-trained weights, utilizing Matrix Factorization (MF) with OpenAI's text-embedding-3-small to embed the queries. For Margin Sampling, we treat it as a train-free baseline, where the routing decision is based on detecting the probability difference between the first and second most likely tokens.

**Metrics** We evaluate the performance of the hybrid inference system using the **Pass Rate**, defined as pass@1 (Chen et al., 2021), based on $n = 32$ sampling iterations. The system's performance is assessed at different proportions (20%, 50%, 80%) of queries routed to the strong model. Furthermore, we incorporate the Average Performance Gap Recovered (APGR) metric from RouteLLM (Ong et al., 2024), which measures the system's ability to recover performance gaps between two LLMs. The APGR is computed across various proportions (0%, 10%, . . . , 100%) of queries routed to the strong model, with values ranging from 0% to 100%, reflecting the extent to which performance discrepancies are mitigated through dynamic routing.

## 4.2 Main Results

### 4.2.1 Overall Performance

Table 1 summarizes the overall performance of various routing methods within a hybrid inference system utilizing Llama3.2-1B and Llama3.1-70B. Methods achieving higher APGR also exhibit improved performance across different proportions of queries routed to the strong model. Our proposed methods outperform existing approaches, with a notable improvement of **3.72% on TriviaQA, 5.29%**

| Datasets | TriviaQA | | | | GSM8K | | | | HumanEval | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | APGR↑ | Pass Rate↑ | | | APGR↑ | Pass Rate↑ | | | APGR↑ | Pass Rate↑ | | |
| | | 20% | 50% | 80% | | 20% | 50% | 80% | | 20% | 50% | 80% |
| Linear Interpolation | 50.00 | 25.75 | 38.59 | 51.44 | 50.00 | 12.60 | 17.72 | 22.85 | 50.00 | 10.27 | 11.44 | 12.61 |
| Hybrid LLM | 49.15 | 24.86 | 38.04 | 51.50 | 61.09 | 14.38 | 20.75 | 24.79 | 51.94 | **10.42** | 11.47 | 12.63 |
| RouteLLM (MF) | 51.22 | 26.14 | 39.45 | 52.28 | 49.11 | 15.11 | 20.72 | 24.66 | 50.33 | 10.10 | 11.26 | 12.65 |
| Margin Sampling | 51.21 | 26.07 | 39.15 | 52.49 | 43.82 | 11.71 | 16.11 | 21.42 | 44.88 | 9.97 | 11.41 | 12.42 |
| Ours (Hard Block) | 53.29 | 27.61 | 41.15 | 52.28 | **65.97** | **16.68** | **22.30** | 24.54 | 50.86 | 10.04 | 11.62 | 12.63 |
| Ours (Soft Block) | **55.38** | **27.91** | **42.28** | **54.28** | 65.48 | 16.01 | 22.04 | **25.24** | **52.37** | 10.33 | **11.72** | **12.80** |

Table 2: Zero-shot performance of various methods across selected datasets, generalizing to different model pairs. Trained on the hybrid inference system of Llama3.2-1B and Llama3.1-70B, and evaluated on the hybrid inference system of Llama3.2-3B and Llama3.1-70B. Linear Interpolation simulates random routing by combining the performance of the two LLMs. **Bolded values** indicate the best-evaluated results.

**on GSM8K, and 2.42% on HumanEval**, demonstrating robustness across diverse query scenarios. Additional visualizations of these results are provided in Appendix B.

On TriviaQA, the Soft Blocking method delivers the best performance, while the Hard Blocking method also outperforms all existing approaches. The poor performance of Hybrid LLM in this context is not surprising, as it relies on win-rate based on BartScore, which has proven unreliable across datasets in Appendix A. Specifically, responses with higher BartScore do not consistently outperform those with lower scores. In contrast, the other methods exceed random routing, demonstrating their effectiveness.

Across both GSM8K and HumanEval, routing methods exhibit consistent patterns, typically either performing well on both datasets or underperforming on both. In contrast, our methods achieve state-of-the-art performance. Although based on the same training data, the difference in training objectives sets our approach apart from Hybrid LLM, underscoring the effectiveness of our method. On the other hand, RouteLLM and Margin Sampling show weaker performance on these datasets, indicating potential generalization challenges. For RouteLLM, these limitations may stem from domain shifts in evaluation tasks and Out-of-Domain challenges associated with LLM selection. As for Margin Sampling, the results indicate that reasoning tasks—such as math, where multiple valid solutions exist—pose difficulties, as they conflict with the fundamental assumption of Margin Sampling, particularly when smaller LLMs are used. Besides, uncertainty in the responses is not the only factor that needs to be considered when deciding whether to cascade to the strong model. This oversight contributes to the failure of Margin Sampling on HumanEval, where multiple valid paths may lead to the correct answer, making the assumption underlying Margin Sampling less reliable in this context.

### 4.2.2 Generalizing to Different Model Pairs

In Table 2, we evaluate the performance of the hybrid inference system configured with Llama3.2-3B and Llama3.1-70B, utilizing routers trained in prior experiments without any additional retraining. Our methods, particularly Soft Blocking, consistently demonstrate superior performance in this configuration, achieving an APGR improvement of **4.16% on TriviaQA, 4.88% on GSM8K, and 0.43% on HumanEval**, which highlights the generalization capability of our method, where routers trained on one model pair exhibit consistent performance when applied to another, confirming its adaptability. Additional visualizations of these results are provided in Appendix B.

On TriviaQA, both of our methods continue to outperform the other approaches, with Hybrid LLM performing worse than random routing. RouteLLM and Margin Sampling show improved performance in this setting. For Margin Sampling, as the weak model scales up, the probability difference becomes a more reliable indicator of response uncertainty in common-sense QA tasks.

On GSM8K and HumanEval, most methods maintain their performance as observed in Table 1. For RouteLLM, since the weakest model in its training setting is Llama-13B, replacing the weak model with Llama3.2-3B likely reduces the gap between the training and evaluation conditions.

### 4.3 Ablation Study

#### 4.3.1 Router Models

An alternative choice for the router model backbone is causal LLMs (Ong et al., 2024). However, we argue that using a router model larger than the weak model incurs unnecessary computa-

| Datasets | TriviaQA | | | | GSM8K | | | | HumanEval | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | APGR↑ | Pass Rate↑ | | | APGR↑ | Pass Rate↑ | | | APGR↑ | Pass Rate↑ | | |
| | | 20% | 50% | 80% | | 20% | 50% | 80% | | 20% | 50% | 80% |
| Weak Model Pass Rate | 50.96 | 20.00 | 35.88 | 50.94 | 51.17 | 12.18 | 17.48 | 22.63 | 53.42 | 8.19 | 10.56 | 12.31 |
| Strong Model Pass Rate | 54.31 | 21.44 | 38.53 | **53.28** | 65.98 | 15.24 | 21.82 | **25.35** | 49.16 | 7.87 | 9.95 | 12.27 |
| Hard Label | 52.05 | 20.80 | 36.51 | 51.51 | 63.26 | 14.68 | 21.02 | 24.91 | 50.63 | 8.61 | 10.12 | 11.97 |
| Hard Blocking w/o SMS | 54.48 | 22.23 | 38.62 | 52.61 | 63.43 | 14.95 | 21.09 | 25.05 | **54.44** | **8.86** | 10.48 | **12.60** |
| Hard Block | 53.16 | 22.09 | 37.85 | 50.96 | **67.37** | **16.34** | **22.46** | 24.67 | 54.36 | 8.17 | **10.77** | 12.27 |
| Soft Block | **55.00** | **22.48** | **38.99** | 52.88 | 66.65 | 15.53 | 22.15 | 25.32 | 53.13 | 8.23 | 10.69 | 12.27 |

Table 3: Zero-shot performance of various label designs across selected datasets. The models were trained and evaluated with the weak model being Llama3.2-1B and the strong model being Llama3.1-70B. **Bolded values** indicate the best-evaluated results.

| Datasets | TriviaQA | GSM8K | HumanEval |
|---|---|---|---|
| Metrics | | APGR↑ | |
| Hard Blocking (Causal) | 51.78 | 57.16 | 54.44 |
| Hard Blocking (Deberta) | 53.16 | **67.37** | 54.36 |
| Soft Blocking (Causal) | 52.44 | 58.55 | **55.31** |
| Soft Blocking (Deberta) | **55.00** | 66.65 | 53.13 |

Table 4: Zero-shot performance of different backbone models (DeBERTa-v3-large, Llama3.2-1B) across selected datasets. Trained and evaluated within the hybrid inference system of Llama3.2-1B and Llama3.1-70B. **Bolded values** indicate the best-evaluated results.

tional costs and impacts response time. As a result, we train the weak model as the router for comparison. As shown in Table 4, DeBERTa-v3-large (with 300M parameters) outperforms Llama3.2-1B, despite its smaller size, demonstrating better performance. Notably, Llama3.2-1B performs better on HumanEval, indicating its superior generalization ability.

### 4.3.2 Label Designs

We also conduct an ablation study on various label designs, as presented in Table 3. More visualized results can be found in Appendix B.

**Pass Rates of the Weak Model**   Training the router using only the weak model's Pass Rates as a soft label results in performance that is only marginally better than random routing. This outcome suggests that effective routing does not primarily rely on the weak model's capacity to provide correct answers, highlighting the need for additional factors to guide the routing process.

**Pass Rates of the Strong Model**   When trained using only the strong model's Pass Rates, the router achieves solid performance but remains outperformed by the two proposed methods. Interestingly, this labeling strategy demonstrates superior results when only a small fraction of queries are routed to the weak model, as it focuses on the strong

model's accuracy and effectively identifies long-tail queries. This suggests that the Pass Rate of the strong model is intrinsically tied to query complexity, as queries that are challenging for the strong model are equally difficult for the weak model. However, in broader scenarios, its performance is surpassed by the more robust Hard Blocking and Soft Blocking techniques.

**Hard Label**   Using hard labels, as defined in Eq 1, derived from greedy sampling leads to improved performance compared to relying solely on the weak model's Pass Rates. This improvement suggests that the router is not simply forwarding unresolved queries from the weak model to the strong model, but instead is capable of learning more sophisticated and nuanced routing strategies.

**Hard Blocking without Strong Model Sampling**   This variant of Hard Blocking, described in Eq 6, is a more economical alternative. By employing greedy decoding on the strong model while performing multiple samplings on the weak model, the router achieves comparable performance. This demonstrates the efficiency of the method, as it reduces computational overhead while maintaining robust routing performance.

## 5 Conclusions

In this work, we propose **Firewall Routing**, a dual-model hybrid inference framework that leverages multiple sampling and innovative blocking techniques to optimize query routing. Through extensive experiments across various benchmarks, our approach demonstrates state-of-the-art performance, significantly reducing computational costs while maintaining high response quality. These results highlight the effectiveness and robustness of the proposed framework in handling routing tasks.

## Limitations

The generalization of the proposed hybrid inference system across different model pairs and datasets remains an area for further exploration. While our approach demonstrates promising results on selected model combinations (e.g., Llama3.2-1B and Llama3.1-70B) and datasets (e.g., TriviaQA, GSM8K, HumanEval), the system's performance may vary when applied to other model pairs or domains. Specifically, the router's ability to generalize across models with varying sizes, architectures, and performance characteristics requires further investigation. Additionally, the limited range of datasets tested raises questions about the system's robustness in more specialized or domain-specific tasks. Future work should include a broader evaluation across diverse models and datasets to assess the scalability and applicability of the proposed approach in real-world, heterogeneous settings.

## References

Ralph Allan Bradley and Milton E Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345.

Lingjiao Chen, Matei Zaharia, and James Zou. 2023. Frugalgpt: How to use large language models while reducing cost and improving performance. *Preprint*, arXiv:2305.05176.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code. *Preprint*, arXiv:2107.03374.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *Preprint*, arXiv:2110.14168.

Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Ruhle, Laks V. S. Lakshmanan, and Ahmed Hassan Awadallah. 2024. Hybrid llm: Cost-efficient and quality-aware query routing. *Preprint*, arXiv:2404.14618.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-

9

hana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Neha Gupta, Harikrishna Narasimhan, Wittawat Jitkrittum, Ankit Singh Rawat, Aditya Krishna Menon, and Sanjiv Kumar. 2024. Language model cascades: Token-level uncertainty and beyond. *Preprint*, arXiv:2404.10136.

Surya Narayanan Hari and Matt Thomson. 2023. Tryage: Real-time, intelligent routing of user

10

prompts to large language models. *Preprint*, arXiv:2308.11601.

Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. Debertav3: Improving deberta using electra-style pretraining with gradient-disentangled embedding sharing. *Preprint*, arXiv:2111.09543.

Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *Preprint*, arXiv:1705.03551.

Anil Kag, Igor Fedorov, Aditya Gangrade, Paul Whatmough, and Venkatesh Saligrama. 2023. Efficient edge inference by selective query. In *The Eleventh International Conference on Learning Representations*.

Sehoon Kim, Karttikeya Mangalam, Suhong Moon, Jitendra Malik, Michael W Mahoney, Amir Gholami, and Kurt Keutzer. 2024. Speculative decoding with big little decoder. *Advances in Neural Information Processing Systems*, 36.

Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pages 19274–19286. PMLR.

Keming Lu, Hongyi Yuan, Runji Lin, Junyang Lin, Zheng Yuan, Chang Zhou, and Jingren Zhou. 2023. Routing to the expert: Efficient reward-guided ensemble of large language models. *Preprint*, arXiv:2311.08692.

Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E. Gonzalez, M Waleed Kadous, and Ion Stoica. 2024. Routellm: Learning to route llms with preference data. *Preprint*, arXiv:2406.18665.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

11

Guillem Ramírez, Alexandra Birch, and Ivan Titov. 2024. Optimising calls to large language models with uncertainty-based two-tier selection. *Preprint*, arXiv:2405.02134.

Marija Šakota, Maxime Peyrard, and Robert West. 2024. Fly-swat or cannon? cost-effective language model choice via meta-modeling. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pages 606–615.

Tal Shnitzer, Anthony Ou, Mírian Silva, Kate Soule, Yuekai Sun, Justin Solomon, Neil Thompson, and Mikhail Yurochkin. 2023. Large language model routing with benchmark datasets. *arXiv preprint arXiv:2309.15789*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. 2022. Orca: A distributed serving system for transformer-based generative models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, pages 521–538.

# A  Is BartScore a Reliable Metric of Response?

We calculate the BartScore for the responses of different LLMs on TriviaQA and GSM8K. The responses are sorted based on their BartScore, and the sorted responses are grouped into bins. Average accuracy is then calculated within each bin to assess the performance of the models at different levels of response correctness.

As shown in Figure 3, 4, a correlation between BartScore and accuracy is only observed on TriviaQA with Llama3.1-70B. In other cases, no consistent or discernible pattern is evident.



Figure 3: BartScore analysis of LLM responses on TriviaQA, GSM8K, and HumanEval. The responses are sorted by BartScore and grouped into bins, with accuracy calculated within each bin to evaluate performance at varying levels of response quality.



Figure 4: BartScore analysis of LLM responses on training set of TriviaQA and GSM8K. The responses are sorted by BartScore and grouped into bins, with accuracy calculated within each bin to evaluate performance at varying levels of response quality.

## B   Visualization of Route Method Performance

Similarly, we rank all queries based on the values predicted by the models, and patch them into distinct bins. For each bin, we compute the average pass rate of the strong model and the weak model. Additionally, we evaluate the improvement in pass rate achieved by routing the queries in each bin to the strong model, rather than to the weak model.

Figure 5: Performance evaluation of reproduced Hybrid LLM on selected datasets. The system utilizes Llama3.2-1B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.
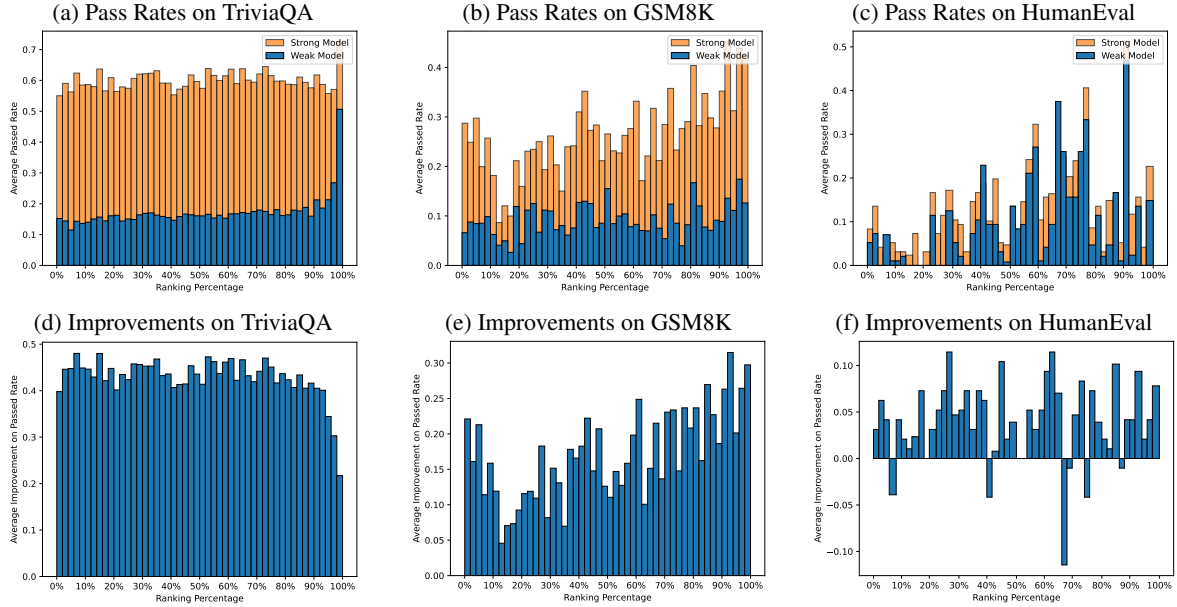


Figure 6: Performance evaluation on generalization of reproduced Hybrid LLM on selected datasets. Evaluated on a system with Llama3.2-3B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.
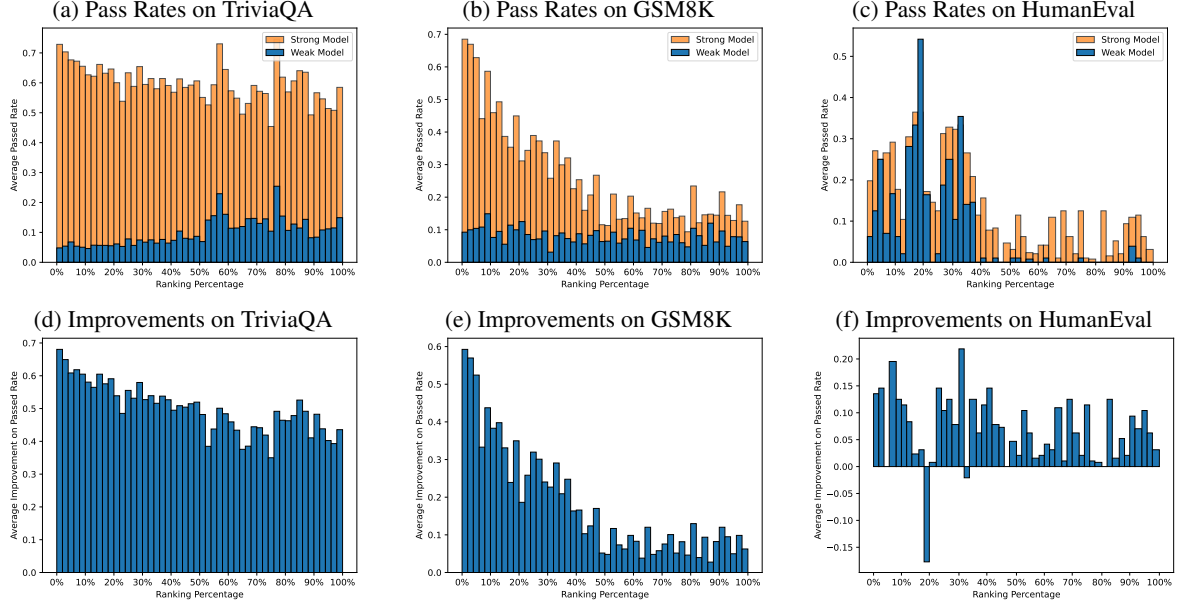
Figure 7: Performance evaluation of Matrix Factorization from RouteLLM on selected datasets. The system utilizes Llama3.2-1B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.
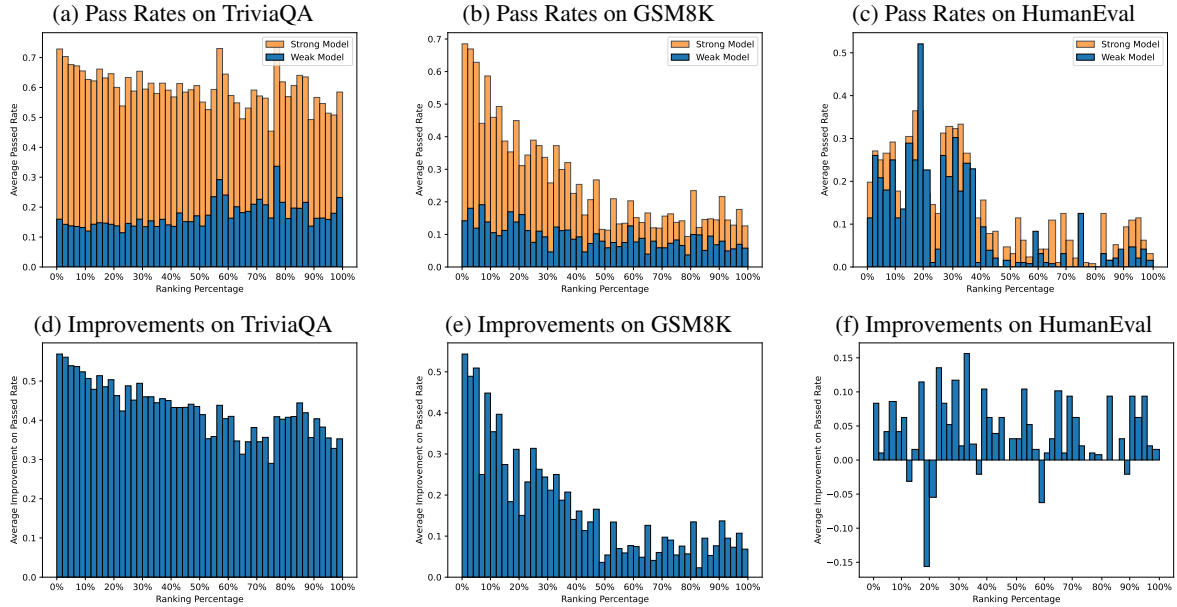


Figure 8: Performance evaluation on generalization of Matrix Factorization from RouteLLM on selected datasets. Evaluated on a system with Llama3.2-3B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.
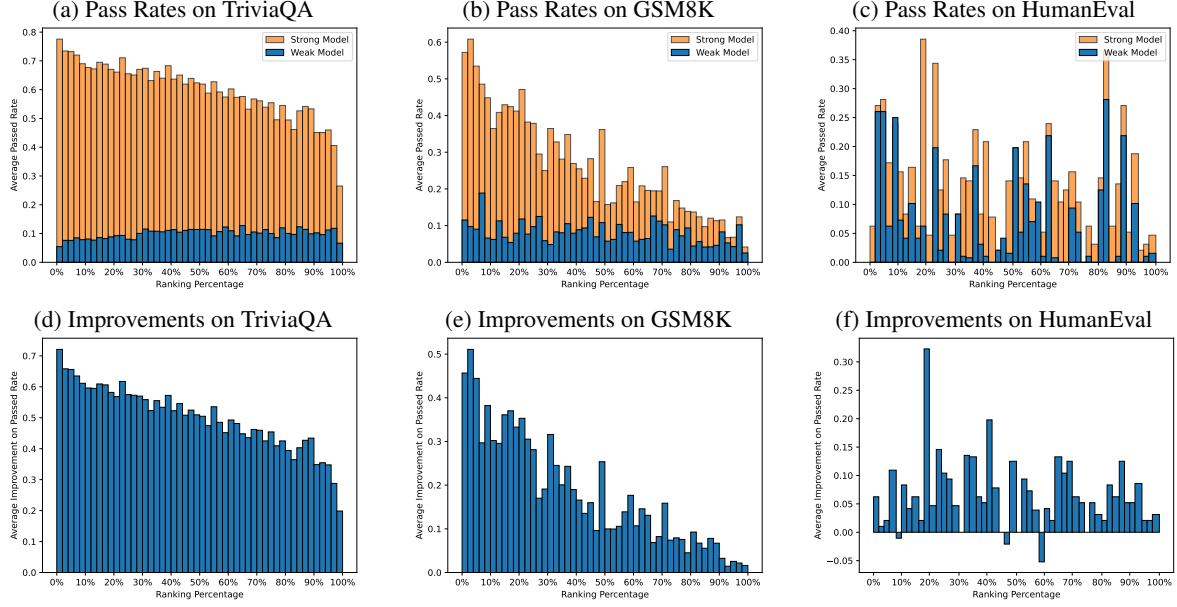
Figure 9: Performance evaluation of Margin Sampling on selected datasets. The system utilizes Llama3.2-1B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.



Figure 10: Performance evaluation on generalization of Margin Sampling on selected datasets. Evaluated on a system with Llama3.2-3B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.

Figure 11: Performance evaluation of Hard Blocking on selected datasets. The system utilizes Llama3.2-1B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.
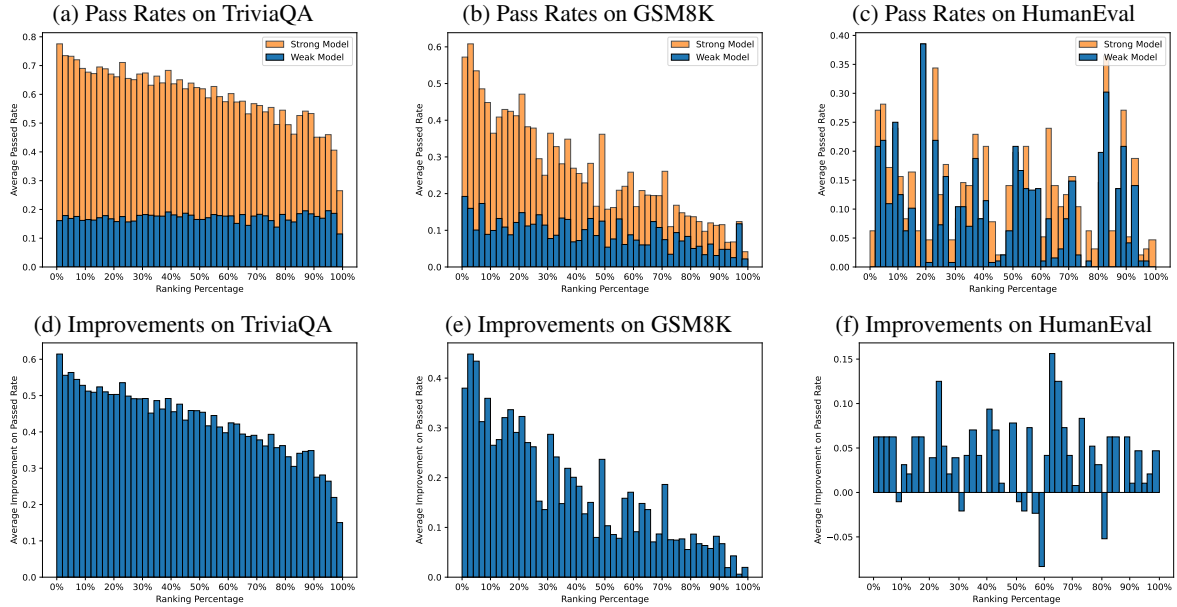


Figure 12: Performance evaluation on generalization of the Hard Blocking on selected datasets. Evaluated on a system with Llama3.2-3B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.
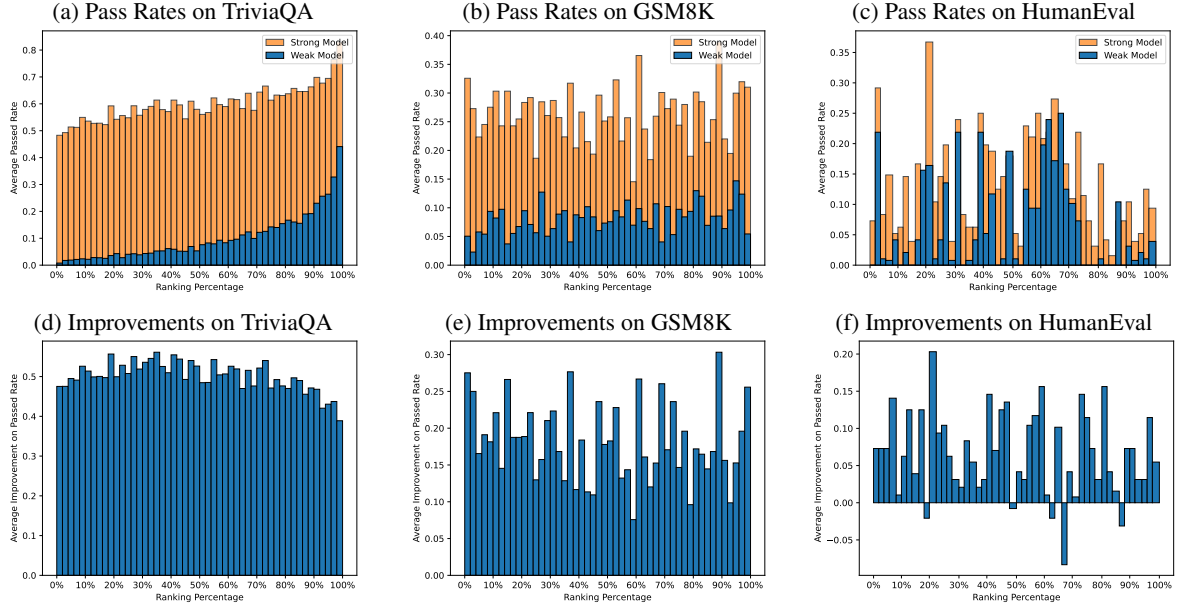
Figure 13: Performance evaluation of Soft Blocking on selected datasets. The system utilizes Llama3.2-1B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.



Figure 14: Performance evaluation on generalization of the Soft Blocking on selected datasets. Evaluated on a system with Llama3.2-3B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.

Figure 15: Performance evaluation of the router trained on Weak Model's Pass Rates across selected datasets. The system utilizes Llama3.2-1B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.
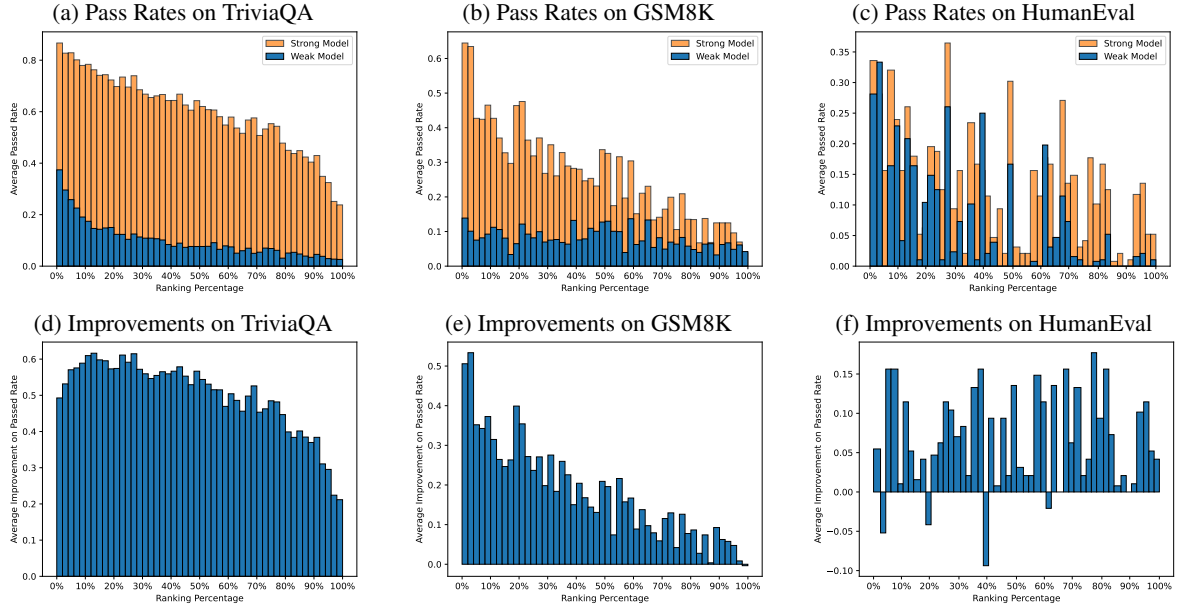


Figure 16: Performance evaluation of the router trained on Strong Model's Pass Rates across selected datasets. The system utilizes Llama3.2-1B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.
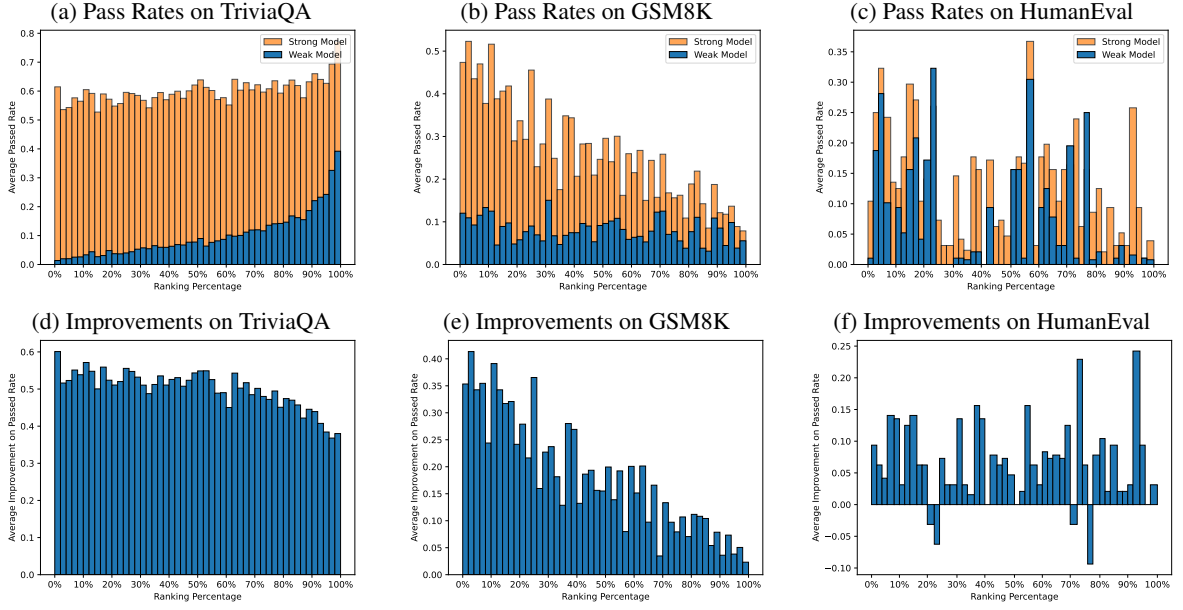
Figure 17: Performance evaluation of the router trained on Hard Labels attained with greedy decoding across selected datasets. The system utilizes Llama3.2-1B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.
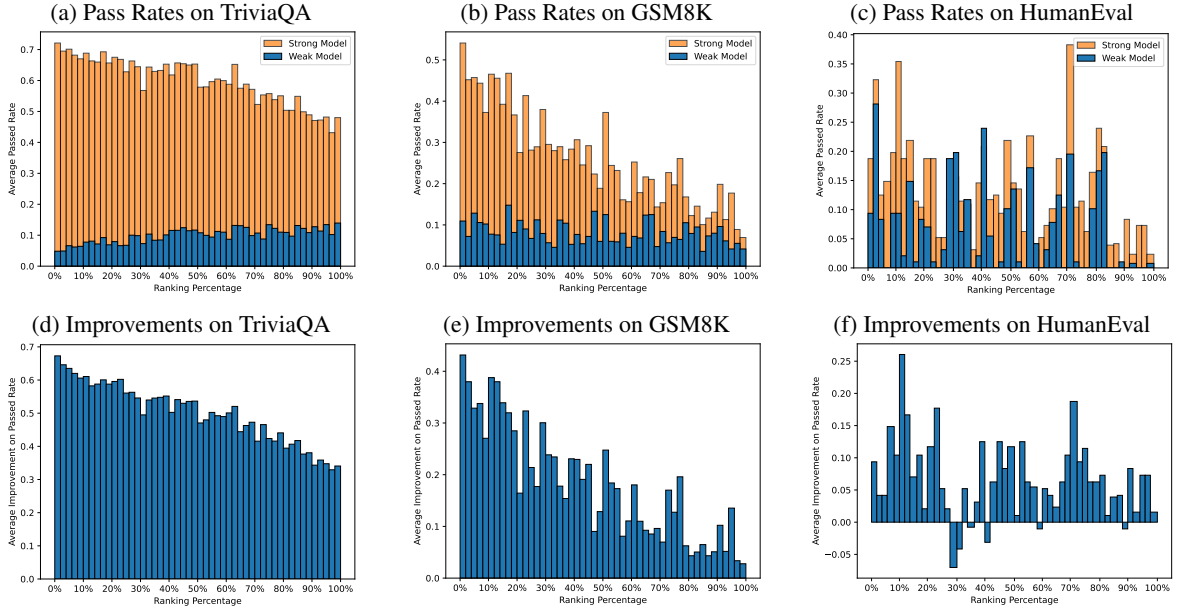


Figure 18: Performance evaluation of the router trained using Hard Blocking without conducting sampling on the strong model across selected datasets The system utilizes Llama3.2-1B as weak model and Llama3.1-70B as strong model. Results are presented in a zero-shot setting.