
Confident or Seek Stronger: Exploring Uncertainty-Based Small LM Routing From Benchmarking to Generalization

Yu-Neng Chuang^{1*} Leisheng Yu^{1*} Guanchu Wang⁴ Lizhe Zhang¹ Ling Chang¹
Hongyi Liu¹ Zirui Liu³ Xuanning Cai⁵ Yang Sui¹ Vladimir Braverman²³ Xia Hu¹
Rice University¹, John Hopkins University², Google Research²,
University of Minnesota³, University of North Carolina at Charlotte⁴, Meta AI⁵
{yc146, ly50}@rice.edu

Abstract

Small language models (SLMs) are increasingly deployed on edge devices for personalized applications, offering efficient decoding latency and reduced energy consumption. However, these SLMs often generate inaccurate responses when handling complex queries. One promising solution is uncertainty-based SLM routing, offloading high-stakes queries to stronger large language models (LLMs) when resulting in low-confidence responses on SLM. This follows the principle of *If you lack confidence, seek stronger support* to enhance reliability. Relying on more powerful LLMs is yet effective but increases invocation costs. Therefore, striking a routing balance between efficiency and efficacy remains a critical challenge. Additionally, efficiently generalizing the routing strategy to new datasets remains under-explored. In this paper, we conduct a comprehensive investigation into **benchmarking and generalization of uncertainty-driven routing strategies from SLMs to LLMs over 5000+ settings**. Our findings highlight: *First*, uncertainty-correctness alignment in different uncertainty quantification (UQ) methods significantly impacts routing performance. *Second*, uncertainty distributions depend more on both the specific SLM and the chosen UQ method, rather than on downstream data. Building on the insight, we propose a proxy routing data construction pipeline and open-source a hold-out set to enhance the generalization on predicting the routing curve for new downstream data. Experimental results indicate that proxy routing data effectively bootstraps routing performance without any new data. The source code is available at <https://github.com/ThunderbornSakana/quodlibeta>.

1 Introduction

Large language models (LLMs) deployment on edge devices has gained increasing attention in recent years, primarily due to their potential for low-latency, privacy-preserving inference. Given the computational and memory constraints of edge devices, small language models (SLMs) (e.g., Phi2-mini [35] or Llama3.2-3B [70]) are designed for resource-efficient deployment, particularly on devices such as smartphones and wearable devices. Their overarching goal is to democratize the deployment of LMs, making it accessible and affordable to users across diverse settings and at any time [52, 86, 83]. However, these SLMs often lack the robustness and scalability of LLMs [8] (e.g., GPT-4o [2] and Llama-3.1-405B), especially when faced with diverse and complex input queries

*Equal contribution

under the deployment on edge devices, which eventually degrade the overall performance. This limitation raises a critical need for exploring solutions to increase the response reliability of SLMs.

To mitigate this unreliability, a line of work proposes to partially offload challenging and complex queries from SLMs to LLMs [11, 59, 32, 66]. A hybrid system is then established to wisely route the queries from SLMs and seek more reliable and deterministic responses from stronger LLMs. Although LLMs can exhibit superior performance, they incur high maintenance and inference costs given the large scale of model size and their infrastructure (i.e., a single NVIDIA A100 GPU can cost approximately \$2,000 per month for deployment). Inaccurate routing by SLMs increases the volume of queries forwarded to LLMs, necessitating greater bandwidth allocation for maintaining the service of LLMs. As a result, operational costs and budgetary requirements rise accordingly, especially when continuous deployment is required. Hence, developing an effective routing strategy is crucial for fully deploying SLMs [59, 66, 11], as it both enhances response reliability and reduces the costs associated with services and data transmission.

Leveraging SLMs’ self-uncertainty estimation emerges as a robust strategy for enhancing routing effectiveness [11, 16]. By relying on the self-assessed uncertainty, the system can better decide whether to handle a query locally or delegate it to a larger model without the aid of extra routers, ensuring that only queries deemed unreliable by the SLMs are routed to LLMs. As a result, the uncertainty-based routing approach not only generalizes well to new datasets, as only self-assessed information from SLM is needed, but it also reduces the high operational costs associated with accurately running LLMs. To this end, we aim to explore two open and nontrivial research questions for uncertainty-based SLM routing:

1) What is the best practice of uncertainty estimation for query routing from SLMs to LLMs?

In this research question, we benchmark the uncertainty-correctness alignment of each uncertainty quantification (UQ) method under its impact on SLM routing. A good alignment is a key factor for successful routing decisions, as any misalignment can cause unnecessary offloading with extra cost. However, SLMs may struggle to provide reliable uncertainty estimates [33, 15, 73], making them less effective as indicators for query routing. Thus, we benchmark the alignment between uncertainty and correctness, paving the insights for establishing more effective routing strategies².

2) What is the best practice to initially establish an effective routing strategy when generalizing to new datasets?

In this research question, we explore how to generalize routing strategies to new datasets. Existing approaches [59, 32] rely on sufficient new downstream data to make routing decisions for optimal performance-cost trade-offs, but this process is time-consuming and labor-intensive. Broadly speaking, collecting and analyzing full downstream datasets under varying SLM configurations can be prohibitively costly, delaying implementation, which is not practical in real-world scenarios. This delay is particularly problematic in high-stakes scenarios, such as medical wearable devices, where reliability is critical, and inaccuracies are unacceptable even in early deployment stages. Based on our findings, we provide a data construction pipeline to predict the routing curves in new downstream scenarios without any new downstream data. A generated proxy routing dataset as a data-agnostic hold-out set enables the estimation of effective routing decisions via the predicted routing curves. We further benchmark the benefits of this proxy routing dataset, demonstrating its generalization ability in predicting the routing curve to new datasets.

This work offers an accessible and reproducible pipeline for uncertainty-based routing from benchmarking to generalization. Our main contributions are summarized as follows:

- **Comprehensive benchmarking and detailed analysis:** This benchmark evaluates 8 UQ methods across 15 datasets to examine the alignment between uncertainty and correctness in routing tasks. We incorporate 8 SLMs and 2 LLMs to emulate real-world deployment scenarios. We then delve into key observations from the extensive results and conclude the insights for developing uncertainty-based SLM routing.
- **Proxy routing data for generalizing routing to new data:** Building on our benchmarking pipeline, we introduce a proxy routing data construction pipeline designed to generalize the routing curve prediction in new downstream scenarios. Empirical results show that this proxy routing data generalizes effectively the routing prediction to new datasets *without relying on any new downstream data*.

²For the convenience of writing, we interchangeably use uncertainty and confidence, where low uncertainty refers to high confidence.

Table 1: Uncertainty quantification (UQ) methods evaluated in our benchmark. “Model Access” specifies whether a method views the LM’s weights/logits (white-box) or only its generated output (black-box). “Require Training?” indicates if additional training is needed. See Subsection 2.1 for taxonomy details and Subsection 3.1 for method descriptions.

Uncertainty Quantification (UQ) Methods	Taxonomy	Model Access	Require Training?
Average Token Prob [53]	Token/sequence probabilities	White-box	No
$p(\text{True})$ [39]	Token/sequence probabilities	White-box	No
Perplexity [21]	Token/sequence probabilities	White-box	No
Jaccard Degree [47]	Output consistency	Black-box	No
Verbalization-1s [76, 69]	Verbalized uncertainty	Black-box	No
Verbalization-2s [69]	Verbalized uncertainty	Black-box	No
Trained Probe [4, 39, 53]	Uncertainty probe	White-box	Yes
OOD Probe [39, 53]	Uncertainty probe	White-box	Yes

2 Reviewing Different Schools of Uncertainty Quantification and LLM Routing

2.1 Uncertainty Quantification for LMs

Uncertainty quantification methods estimate a model’s confidence in its predictions [31]. For traditional classification and regression, uncertainty estimation is well-established [23]. However, for LLMs generating free-form responses to complex queries, estimating uncertainty is more challenging because the output space can grow exponentially with vocabulary size, and each sequence spans multiple tokens [20]. Existing uncertainty quantification approaches for LLMs can be grouped into the following four categories.

Via verbalizing uncertainty. This line of work prompts language models to report linguistic confidence [53, 56]. To enable LMs to verbalize confidence, researchers have proposed fine-tuning them to express uncertainty [46] or teaching them to verbalize confidence through in-context learning [17]. Verbalized confidence can take the form of linguistic expressions of uncertainty or numerical scores [24]. Multiple studies find that LLMs tend to be overconfident when reporting confidence [76, 69]. To mitigate this overconfidence, prompting strategies such as multi-step elicitation, top- k , and Chain-of-Thought [72] have been explored [69]. Sampling multiple response-confidence pairs and designing more effective aggregation strategies can also help mitigate overconfidence [76]. Moreover, [69] reports that verbalized confidence is typically better calibrated than the model’s conditional probabilities.

Via analyzing token/sequence probabilities. This line of research derives confidence scores from model logits for output tokens [24, 33, 38]. The confidence of a generated sequence is computed by aggregating the log-probabilities of its tokens. Common aggregation strategies include arithmetic average, minimum, perplexity, and average entropy [20, 21, 71]. Because not all tokens in a sequence equally reflect semantic content, SAR reweights token likelihoods to emphasize more meaningful tokens [18]. However, different surface realizations of the same claim can yield different probabilities, implying that the calculated confidence reflects how a claim is articulated rather than the claim itself [53]. To combine LM self-assessment with token probabilities, $p(\text{True})$ is proposed: the model is asked whether its generated response is correct, and the probabilities of True/False tokens serve as the confidence score [39, 69].

Via gauging output consistency. This line of research (e.g., SelfCheckGPT [54]) assumes that high-confidence LLMs produce consistent outputs [53]. A typical approach samples m responses for a given input query, measures inter-response similarity, and calculates a confidence score from meaning diversity [20]. Common ways to measure pairwise similarity include Natural Language Inference (NLI) and Jaccard similarity [24]. Consistency is then assessed by analyzing the similarity matrix, for instance, by counting semantic sets, summing eigenvalues of the graph Laplacian or computing eccentricity [47]. Because different sentences can express the same meaning, semantic entropy [40] first clusters responses by semantic equivalence before measuring consistency.

Via training uncertainty probes. This approach trains classifiers to predict whether an LLM will arrive at the correct answer for a particular query, using predicted probabilities as confidence

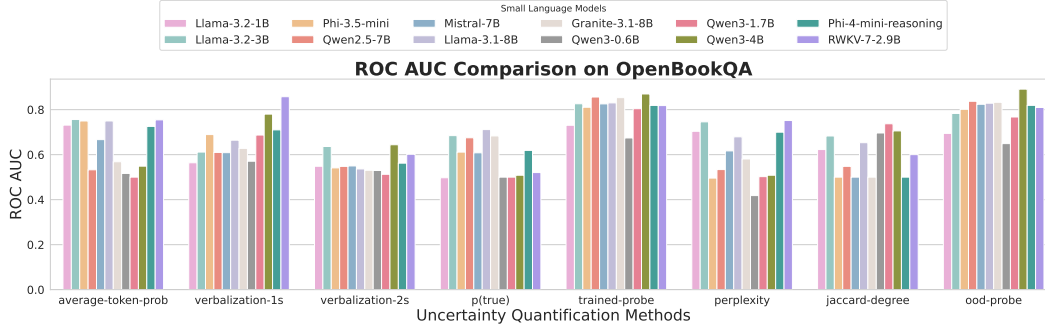


Figure 1: The ROC AUC scores measure the alignment between confidence and correctness across different SLMs and uncertainty quantification methods on OpenBookQA. A higher ROC AUC indicates a stronger alignment.

scores [24]. Training data is often obtained by sampling multiple answers per question at a fixed temperature and labeling each for correctness [39]. A probe (commonly a multi-layer perceptron) then takes hidden states as inputs to predict correctness [4, 42]. Because in-domain training data is not always available, Contrast-Consistent Search trains probes unsupervisedly by maximizing representation distances between contradictory answers on Yes/No questions [7]. Furthermore, whether probes trained on out-of-distribution data remain effective is still under debate [39, 53, 40].

2.2 LLM Routing

In query-routing scenarios, recent approaches train additional classifiers to direct queries to different SLMs or LLMs based on historical performance metrics and user feedback data [16, 59, 66, 37, 84]. For instance, RouterBench [32] collects inference outputs from selected LLMs to aid in the development of routing classifiers. However, these methods face significant challenges when encountering new downstream tasks, as such data falls outside the distribution of the existing training data. This limitation makes them less practical for real-world scenarios, such as on personal edge device deployment, where adaptability to unseen conditions is crucial. Our work focuses on how to establish routing systems between SLMs and LLMs and generalize to new downstream tasks. In this manner, uncertainty-based routing is an appropriate solution to overcome these challenges, as uncertainty is directly extracted from SLMs themselves. Furthermore, we propose a proxy routing data construction pipeline to initialize a routing system that generalizes to unseen datasets.

3 Benchmarking Uncertainty-based SLM Routing

In this section, we systematically evaluate 12 SLMs and 4 LLMs on 15 datasets using 8 UQ methods (see Table 1) for uncertainty-based SLM routing. This section details the datasets, models, and UQ methods, followed by several key findings and practical considerations. All experiments are conducted on four 80GB NVIDIA A100 GPUs.

3.1 Benchmark Coverage and Setup

Language Models. We evaluate 12 open-source SLMs, organized into three categories: non-reasoning LMs, reasoning LMs, and a recurrent neural network (RNN) model. The non-reasoning models are Llama-3.2-1B-Instruct [55], Llama-3.2-3B-Instruct [55], Phi-3.5-mini-instruct [1], Mistral-7B-Instruct-v0.3 [36], Qwen2.5-7B-Instruct [78], Llama-3.1-8B-Instruct [19], and Granite-3.1-8B-Instruct [26]. The reasoning models are Qwen3-0.6B, Qwen3-1.7B, Qwen3-4B, and Phi-4-mini-reasoning [77]. The RNN model is RWKV-7-2.9B [61]. These SLMs come from Alibaba (four models), Meta (three), Microsoft, Mistral AI, IBM, and LF AI & Data. Except for RWKV-7-2.9B, all adopt decoder-only Transformer architectures and are available on Hugging Face. We also include four LLMs: three open-source models—Llama-3.1-70B-Instruct [19], Qwen3-32B, and DeepSeek-R1 [29]—and one proprietary API model, GPT-4.1 mini [34]. Qwen3-32B and DeepSeek-R1 are reasoning LLMs, whereas Llama-3.1-70B and GPT-4.1 mini are non-reasoning.

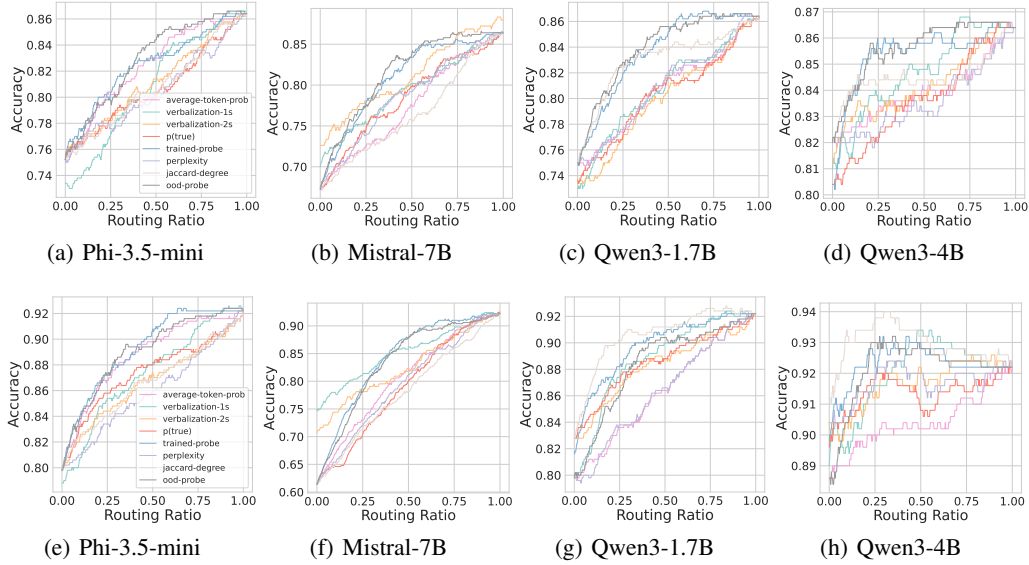


Figure 2: Overall accuracy vs. routing ratio with different UQ methods and SLMs. (a)-(d) show the results of routing to DeepSeek-R1 on the CommonsenseQA dataset; and (e)-(h) demonstrate the results of routing to GPT-4.1 mini on the OpenBookQA dataset.

Datasets. Experiments span 15 datasets from four domains: (1) *Mathematical Reasoning* (AQuA [48], GSM8K [13], MultiArith [63], SVAMP [60], MATH-500 [43]), (2) *Commonsense Reasoning* (CommonsenseQA [67], HellaSwag [80], OpenBookQA [57], PIQA [6], TruthfulQA [45], Wino-Grande [64], BoolQ [12], Social IQa [65]), (3) *Conversational and Contextual Understanding* (CoQA [62]), and (4) *Problem Solving* (MMLU [30]). These cover free-form, multiple-choice, and True/False question answering and are available via Hugging Face. Table 2 in Appendix B provides further details.

UQ Methods and Hyperparameters. We evaluate 8 approaches from the four categories in Section 2.1. (1) *Average token probability* uses the probability of the chosen option token (e.g., “A”) for multiple-choice tasks or the mean probability of all generated tokens for free-form tasks. (2) *Perplexity* is computed for a sequence of N output tokens $\{y_i\}_{i=1}^N$ with probabilities $\{p(y_i)\}_{i=1}^N$ as $\exp(\frac{1}{N} \sum_{i=1}^N \ln p(y_i))$, and its reciprocal serves as the confidence score. (3) $p(\text{True})$ is a method where the LM first outputs an answer, then evaluates the generated response using only “True” or “False.” The probabilities for these two tokens are normalized to sum to 1, and the probability of “True” is used as confidence. (4) *Verbalized confidence in a single response* (denoted as verbalization-1s) prompts the model to output both the answer and numeric confidence in one step. (5) *Verbalized confidence in the second round* (denoted as verbalization-2s) obtains the confidence in a separate, follow-up query after the model has provided an answer. (6) *The degree matrix* (denoted as jaccard-degree) generates $m = 5$ samples (temperature 1.0) for one query, computes pairwise Jaccard similarities, and sets confidence to $\text{trace}(mI - D)/m^2$, where D is the degree matrix. (7) *Trained probe* is a four-layer MLP with LeakyReLU activations, trained on a fixed subsample of the in-domain training set for each dataset, taking as input the hidden states from the eighth-to-last transformer layer. We train for 20 epochs (learning rate 5×10^{-4}). (8) *Trained probe on out-of-distribution data* (denoted as ood-probe) is identical in architecture but trained on all other datasets. e.g., if AQuA is evaluated, the ood-probe is trained on the remaining 15 datasets (20 epochs, learning rate 1×10^{-4}).

For verbalization-based methods, we discard queries when the model does not follow instructions to produce a confidence score. For free-form question answering, we use GPT-4.1 mini to evaluate whether a response is essentially equivalent to the ground truth answer [85].

3.2 Report Observations

In this section, we present our benchmarking results analyzing the impact of uncertainty-correctness alignment on routing tasks. More observations and experimental results on proxy routing and routing can be found in Appendix D.1.

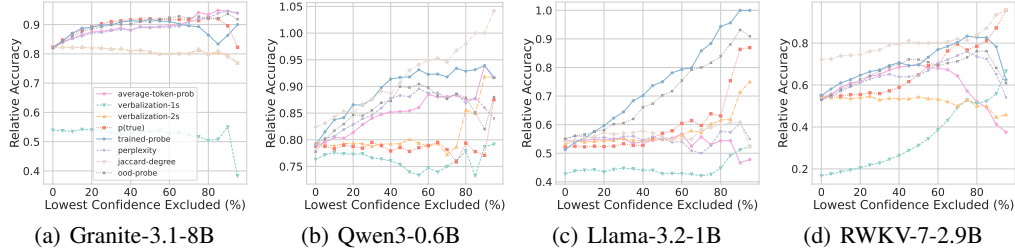


Figure 3: Relative accuracy of SLMs vs. LLMs on top- k % confident queries. “Relative accuracy” is the ratio of SLM accuracy to LLM accuracy. The x-axis “Lowest Conf. Excluded” shows the percentage of low-confidence queries removed; for example, 80 means 80% of queries with the lowest confidence are excluded, leaving the top 20%. (a) and (b) compare SLMs to Llama-3.1-70B on GSM8K, while (c) and (d) compare SLMs to Qwen3-32B on BoolQ.

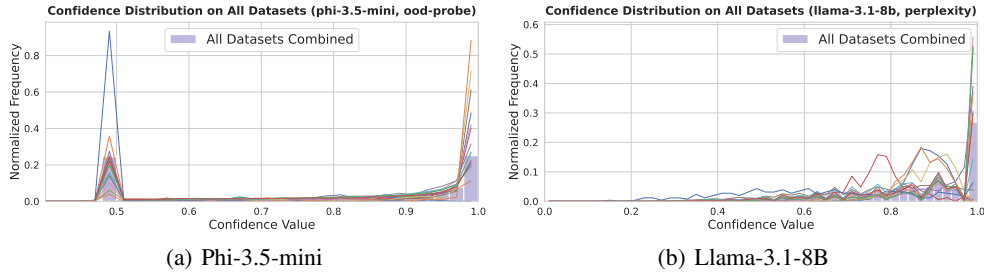


Figure 4: Confidence distributions across 15 datasets. The histogram depicts the aggregated distribution from all datasets, while each curve represents a single dataset. (a) Confidence of Phi-3.5-mini by OOD Probe; (b) Confidence of Llama-3.1-8B by Perplexity.

Observation ①: Uncertainty estimation in SLMs may exhibit misalignment with prediction correctness. From the theoretical perspective, well-calibrated uncertainty scores do not necessarily imply a strong correlation with the correctness of the predictions [33, 11]. The predictions of models might be perfectly calibrated yet still display relatively low accuracy (i.e., confidently provide wrong answers). This phenomenon is also evident in our benchmark results (illustrated in Figure 1). We compute AUC scores to quantify the correlation between extracted uncertainty and prediction correctness, treating correctness as a binary ground truth and using confidence values as the ranking metric. The results show that not all UQ methods effectively exhibit a strong alignment between confidence and prediction correctness. Moreover, from Figure 1 and Figure 9, we can observe that the alignment may vary across datasets for the same SLM and UQ method. For instance, Perplexity [21] demonstrates strong alignment for Phi-3.5-mini on the MultiArith dataset but fails on the OpenBookQA dataset. On the other hand, OOD Probe, Trained Probe, and Perplexity obtain consistently decent alignment compared to other UQ methods across different SLMs and domains of datasets. Conversely, we notice that verbalization-based methods, namely verbalization-1s [69, 53], and verbalization-2s [69], consistently withhold low alignment between uncertainty and prediction correctness. More experimental results can be found in Appendix D.1.

Observation ②: Verbalization-based UQ methods struggle to extract uncertainty in SLMs for query routing. We find that verbalization methods like verbalization-2s [69] obtain poor alignment between confidence and prediction correctness, and this misalignment can lead to inferior routing performance in SLMs, where the conclusion can be found in Figure 2. Recent advancements [75, 79] also show that uncertainty scores derived from verbalization may exhibit good reflection on models’ intrinsic uncertainty of prediction across multiple models and datasets. This discrepancy poses a significant challenge for establishing effective routing performance since queries that are actually correct may be unnecessarily routed from SLMs to LLMs, thereby increasing the overall cost of deploying routing systems.

Observation ③: A good routing standard highly depends on UQ methods with good uncertainty-correctness alignment. A notable phenomenon occurs when UQ methods, such as Trained Probe [53], exhibit strong alignment, leading to significant improvements in routing performance. This is because the extracted uncertainty scores from these UQ methods more effectively indicate whether SLMs produce correct predictions. Among all UQ methods evaluated for routing tasks, we find that Trained Probe [53], OOD Probe [39, 53], and Perplexity [20] consistently rank as the top three methods for SLM routing. Therefore, a comprehensive analysis of UQ methods before deploying a routing system in SLMs is highly recommended to ensure efficient query routing.

Observation ④: SLMs can match LLM performance on high-confidence queries. Although SLMs generally underperform LLMs, we find that for queries where SLMs exhibit high confidence, their accuracy approaches that of LLMs. To illustrate, we progressively remove queries starting from those with the lowest SLM confidence and compute the ratio of SLM to LLM accuracy on the remaining top- $k\%$ queries (Figure 3). As more low-confidence queries are excluded, SLMs achieve comparable performance to LLMs. For instance, on GSM8K, Qwen3-0.6B achieves performance nearly equal to Llama-3.1-70B on the top 20% highest-confidence queries. Moreover, the effectiveness of this selection depends on the uncertainty quantification (UQ) method: approaches with stronger alignment (e.g., Trained Probe [53]) yield higher relative accuracy than weaker ones (e.g., verbalization-2s) across all query exclusion rates. Additional results appear in Appendix D.2.

4 Generalize SLM Routing for New Downstream Scenarios

In this section, we first describe the training-free pipeline for constructing proxy routing data for SLM routing with experimental details. We then investigate how well the proxy routing data can predict the routing curve for new downstream scenarios without accessing the new datasets. Finally, we discuss our results and offer several insights into the proxy routing data for establishing routing in early-stage deployments.

4.1 Proxy Routing Data Construction Pipeline

We aim to evaluate the effectiveness of proxy routing data in generalizing routing curve predictions to new downstream scenarios without relying on any downstream data. The proxy routing data functions as a data-agnostic hold-out set tailored to a particular SLM, enabling the transfer of routing standards across diverse datasets. By leveraging this proxy data, we establish a generalizable routing framework that simplifies deployment and removes the need for dataset-specific routing analysis, while demonstrating that proxy routing data can effectively transfer routing behavior across domains.

Overview of Construction. Let $\mathbb{D} = \{\mathcal{D}_i\}_{i=1}^N$ denote a diverse collection of datasets spanning multiple domains (e.g., commonsense, mathematics), following the setup in [50]. Each instance in \mathbb{D} is processed with a selected UQ method to obtain an uncertainty representation $\mathcal{F}_{\mathcal{D}_i}$ that captures the SLM’s confidence characteristics in that domain. These aggregated representations provide the foundation for constructing proxy routing data.

Proxy Data Formation. We construct the proxy routing dataset by pooling all domains in \mathbb{D} into a unified source distribution $\mathbb{D} \sim \mathbb{P}_{\text{pool}}$. Due to on-device storage and memory constraints, we uniformly sample m ($k\%$) instances from pooled distribution:

$$X_1, X_2, \dots, X_m \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{\text{pool}}.$$

This uniform sampling preserves the natural diversity of uncertainty patterns across domains while maintaining computational and storage efficiency.

4.2 Proxy Routing Data Setups

Benchmark Settings. We evaluate the constructed proxy routing data on 15 SLMs and 4 LLMs across 15 datasets. Based on the observations and results from the previous benchmark section, we select 2 UQ methods that demonstrate the strongest alignment between predicted uncertainty and actual correctness: "OOD Probe" [39, 53] and "Perplexity" [20] method. We consider the routing performance evaluated on the entire new dataset as the ground truth. To simulate new dataset

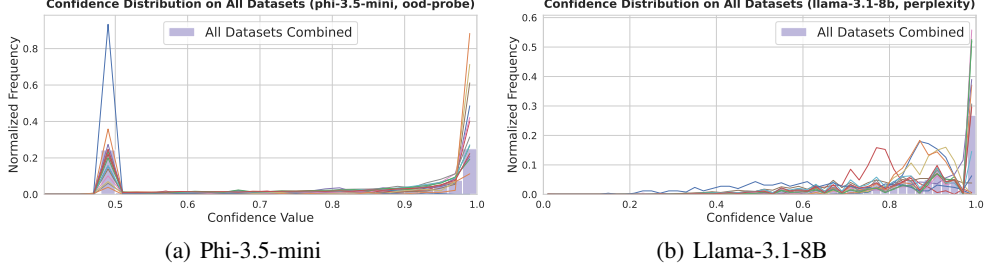


Figure 5: Confidence distributions across 15 datasets. The histogram depicts the aggregated distribution from all datasets, while each curve represents a single dataset. (a) Confidence of Phi-3.5-mini by OOD Probe; (b) Confidence of Llama-3.1-8B by Perplexity.

scenarios, we introduce two evaluation settings: (1) fully out-of-domain and (2) partially in-domain. First, for the out-of-domain setting, we evaluate a target dataset using proxy routing data derived from source datasets with no domain overlap. Second, in the partially in-domain setting, we designate one dataset as the target and construct its proxy routing data using the remaining 14 datasets, where the domain of the dataset may partially overlap. The target dataset’s generalization performance is then evaluated using this proxy routing set, which does not contain any information from the target dataset. All reported results represent the average across three individual experimental runs.

Data Construction Settings. The proxy routing data is weighted-sampled from each bin of the proxy routing data distributions, with the number of bins set to 30. We sample $k\%$ of the instances from each bin to form the final proxy routing data. The temperature is fixed at 0 with a fixed random seed of 50 to ensure reproducibility. In this work, we use $k=10$ to meet the limited resources of on-device machine routing. We also provide the sensitivity check and results comparison with baselines of $k\%$ sample rate in Appendix G and E.

4.3 Theoretical Analysis and Intuition

We provide a theoretical justification for the proposed proxy routing framework through a sequence of three theorems. The first two theorems formalize the regularity conditions of the UQ function, while the final theorem derives an end-to-end generalization bound that quantifies how proxy routing performance transfers to unseen downstream datasets. The proof are provided in Appendix ??.

Theorem 1 (Lipschitz continuity of practical UQ mappings). *Let $g: \mathcal{X} \rightarrow \mathbb{R}^d$ denote the SLM representation and $u: \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ denote a UQ mapping. Assume that (i) logits $\ell = g(x)$ are bounded, $\|\ell\| \leq M$; and (ii) u is a finite composition of continuously differentiable functions with bounded Jacobians on compact domains (e.g., exp, log, softmax, max, linear maps, or polynomials). Then u is Lipschitz continuous on its bounded domain; that is, there exists $L_u < \infty$ such that $\forall x, x' \in \mathcal{X}$,*

$$|u(g(x)) - u(g(x'))| \leq L_u L_g \|x - x'\|.$$

Intuition of Theorem 1. In Theorem 1, we showcase that under mild and realistic conditions such as bounded logits and smooth functional composition, the uncertainty mapping $u(g(x))$ is globally Lipschitz on compact domains. Intuitively, this means that similar inputs yield similar uncertainty scores, so small shifts in the representation or domain lead to small and predictable changes in uncertainty. When the input differs substantially from the proxy domain, the change in uncertainty may increase proportionally, but the Lipschitz property still guarantees that this growth remains bounded rather than arbitrary. This controlled sensitivity ensures that routing decisions remain stable under moderate distribution shifts.

Theorem 2 (Bounded density of uncertainty distribution near routing threshold). *Let $U = u(g(X))$ be the scalar uncertainty score of $X \sim \mathbb{P}_X$ with continuous density f_X bounded by B_X . If u is Lipschitz with constant $L_u < \infty$, and U is smoothed by additive uniform noise $U_\sigma = U + \zeta$ with $\zeta \sim \text{Unif}[-\sigma, \sigma]$, $\sigma > 0$, then U_σ admits a continuous density f_{U_σ} satisfying*

$$\|f_{U_\sigma}\|_\infty \leq \|f_U\|_\infty \leq \frac{B_X}{L_u}.$$

Hence, the uncertainty distribution remains smooth and bounded around any routing threshold.

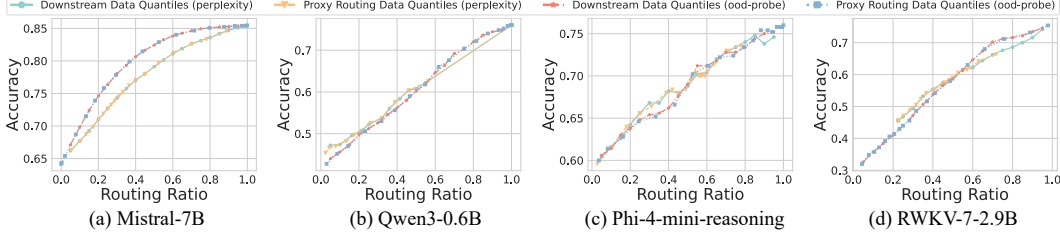


Figure 6: Routing results from four SLMs to Llama-3.1-70B on HellaSwag, with the remaining 14 other datasets constituting the proxy routing data.

Intuition of Theorem 2. Theorem 2 provides a theoretical justification for using uniform sampling when constructing the proxy routing dataset. Since the uncertainty scores $U = u(g(X))$ admit a bounded and continuous density after mild uniform smoothing, the samples in the pooled distribution \mathbb{P}_{pool} are already well distributed across the uncertainty space. Uniform sampling therefore preserves this smooth uncertainty landscape without introducing selection bias, ensuring that the proxy dataset accurately reflects the diversity of uncertainty patterns across domains and supports stable threshold estimation for cross-domain routing.

Theorem 3 (End-to-End Proxy Routing). Let \mathbb{P}_{pool} be a pooled multi-domain input distribution and $X_1, \dots, X_m \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{\text{pool}}$. Define $U_j = u(g(X_j))$, the empirical routing objective

$$\hat{R}_m(\tau) = \frac{1}{m} \sum_{j=1}^m r(\mathbf{1}\{U_j \geq \tau\}, X_j),$$

and population objective

$$R_{\text{pool}}(\tau) = \mathbb{E}_{\mathbb{P}_{\text{pool}}}[r(\mathbf{1}\{U \geq \tau\}, X)].$$

Let $\hat{\tau} \in \arg \max_{\tau} \hat{R}_m(\tau)$ and $R_{\mathbb{Q}}(\tau)$ denote the routing objective under a new distribution \mathbb{Q} . Assume: (i) $r(1, x) = r(0, x) + \Delta(x)$ with $|\Delta(x)| \leq \Delta_{\max}$; (ii) the uncertainty densities under \mathbb{P}_{pool} and \mathbb{Q} are bounded by B ; (iii) g, u are Lipschitz with constants L_g, L_u . Then, for universal constants $C, c > 0$ and any $\alpha \in (0, 1)$, with probability at least $1 - \alpha$,

$$\begin{aligned} \max_{\tau} R_{\mathbb{Q}}(\tau) - R_{\mathbb{Q}}(\hat{\tau}) &\leq 2C \sqrt{\frac{\log(c/\alpha)}{m}} \\ &\quad + \Delta_{\max} 2\sqrt{2B L_u L_g W_1(\mathbb{P}_{\text{pool}}, \mathbb{Q})}, \end{aligned}$$

where $W_1(\cdot, \cdot)$ denotes the 1-Wasserstein distance.

Intuition of Theorem 3. Given these regularity conditions, Theorem 3 provides a complete generalization bound for proxy routing across domains. It decomposes the cross-domain routing regret into two interpretable terms: (i) a *finite-sample term* of order $O(1/\sqrt{m})$, which arises from estimating the routing threshold on a limited proxy dataset, and (ii) a *domain-shift term* of order $O(\sqrt{W_1(\mathbb{P}_{\text{pool}}, \mathbb{Q})})$, which quantifies how the expected routing performance changes with the Wasserstein distance W_1 between the proxy and target input distributions. The second term explicitly depends on the Lipschitz constants (L_g, L_u) and the density bound B , linking the generalization behavior to the smoothness of the model and the UQ mapping. Together, these results ensure that a uniformly sampled proxy dataset can approximate the routing behavior of unseen downstream tasks without requiring additional labeled data.

Practical UQ Methods. As the discussion and details shown in Section 2.1, these standard assumptions are well and trivially satisfied by all practical UQ methods we consider, such as token probability, perplexity, $p(\text{True})$, verbalized confidence, diversity-based, and probe-based measures. Some other UQ methods, such as Jaccard Degree, may apply lightweight preprocessing steps like logit clipping or temperature scaling to meet the assumption. These standard treatments ensure that the theoretical conditions are not only mathematically justified but empirically realizable in practice.

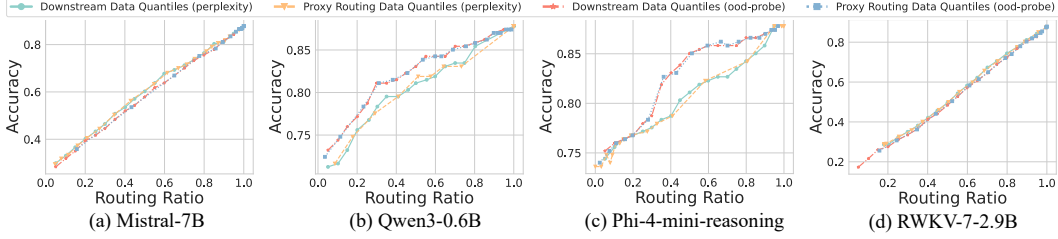


Figure 7: Routing results from four SLMs to DeepSeek-R1 on AQuA (mathematical reasoning), with eight commonsense-reasoning datasets and one conversational & contextual understanding dataset constituting the proxy routing data, demonstrate the strong generalization capability of the proxy routing data.

4.4 Routing Curve Prediction by Proxy Data

We provide several key insights into the generalization ability of proxy routing data as follows.

Insights ①: The extracted confidence distribution is predominantly determined by the chosen SLM and UQ method, with minimal dependence on the downstream dataset. As illustrated in Figure 5, confidence scores aggregated from 15 different tasks exhibit a nearly identical shape regardless of the specific dataset. Instead, they vary notably with different SLMs and UQ methods. This finding suggests that the confidence distribution is largely data-agnostic, enabling the construction of proxy routing data that generalizes to new tasks without any new datasets.

Insights ②: Proxy routing data helps SLM routing to predict an accurate routing curve without any new data, allowing routing strategies to be initialized on SLMs without accessing new datasets. Building on our findings about uncertainty distributions, we sampled a data subset to create a final proxy routing dataset using the pipeline described in Section 4.1. We then utilized this proxy routing dataset to predict all thresholds for different routing ratios in new downstream scenarios. The experimental results (see Figure 6 and Figure 7) show that the routing curves from the proxy routing data closely match those from the entire new downstream dataset in both evaluation settings, indicating that the proxy routing data provides strong capability for establishing routing strategies on unseen downstream datasets. An identical phenomenon is observed across multiple UQ methods and different SLMs, highlighting the potential of proxy routing data to initiate for any new dataset, independent of the UQ method or SLM used. More results are in Appendix D.3.

5 Conclusion

This paper investigates the routing accuracy of SLMs in estimating their uncertainty and establishing best practices for initiating effective routing strategies. Through comprehensive benchmarking of 15 SLMs, 4 LLMs, 8 UQ methods, and 15 datasets across 5000+ settings, we found that the alignment between uncertainty and correctness significantly impacts routing performance. Additionally, our experiments show that uncertainty distributions depend primarily on the specific SLM and UQ method rather than the downstream data. Building on the insights, we introduced a proxy routing data construction pipeline and a hold-out dataset to generalize routing strategies without prior knowledge of new downstream data. The results confirm that the proxy routing data effectively bootstraps routing, indicating its strong potential for benefiting in resource-efficient SLM deployment.

References

- [1] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

- [3] Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Leandro von Werra, and Thomas Wolf. Smollm - blazingly fast and remarkably powerful, 2024.
- [4] Amos Azaria and Tom Mitchell. The internal state of an llm knows when it’s lying. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 967–976, 2023.
- [5] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [6] 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, pages 7432–7439, 2020.
- [7] Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. *arXiv preprint arXiv:2212.03827*, 2022.
- [8] Lihu Chen and Gaël Varoquaux. What is the role of small models in the llm era: A survey. *arXiv preprint arXiv:2409.06857*, 2024.
- [9] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198, 2024.
- [10] Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yiming Hu, Xinyang Lin, Bo Zhang, et al. Mobilevlm v2: Faster and stronger baseline for vision language model. *arXiv preprint arXiv:2402.03766*, 2024.
- [11] Yu-Neng Chuang, Helen Zhou, Prathusha Kameswara Sarma, Parikshit Gopalan, John Boccio, Sara Bolouki, and Xia Hu. Learning to route with confidence tokens. *arXiv preprint arXiv:2410.13284*, 2024.
- [12] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936, 2019.
- [13] 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.
- [14] Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- [15] Gianluca Detommaso, Martin Bertran, Riccardo Fogliato, and Aaron Roth. Multicalibration for confidence scoring in llms. *arXiv preprint arXiv:2404.04689*, 2024.
- [16] Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Rühle, Laks VS Lakshmanan, and Ahmed Hassan Awadallah. Hybrid llm: Cost-efficient and quality-aware query routing. In *The Twelfth International Conference on Learning Representations*, 2024.
- [17] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, et al. A survey on in-context learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1107–1128, 2024.
- [18] Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5050–5063, 2024.

- [19] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [20] Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, et al. Lm-polygraph: Uncertainty estimation for language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 446–461, 2023.
- [21] Marina Fomicheva, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:539–555, 2020.
- [22] Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pages 10323–10337. PMLR, 2023.
- [23] Yarin Gal et al. Uncertainty in deep learning. 2016.
- [24] Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koepl, Preslav Nakov, and Iryna Gurevych. A survey of confidence estimation and calibration in large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6577–6595, 2024.
- [25] Xinyang Geng and Hao Liu. Openllama: An open reproduction of llama, May 2023.
- [26] IBM Granite Team. Granite 3.0 language models, 2024.
- [27] RLHF Griffin and Gemma Teams. Recurrentgemma: Moving past transformers for efficient open language models, 2024.
- [28] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. Olmo: Accelerating the science of language models. *Preprint*, 2024.
- [29] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [30] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021.
- [31] Bairu Hou, Yujian Liu, Kaizhi Qian, Jacob Andreas, Shiyu Chang, and Yang Zhang. Decomposing uncertainty for large language models through input clarification ensembling. In *Forty-first International Conference on Machine Learning*, 2024.
- [32] Qitian Jason Hu, Jacob Bieker, Xiuyu Li, Nan Jiang, Benjamin Keigwin, Gaurav Ranganath, Kurt Keutzer, and Shriyash Kaustubh Upadhyay. Routerbench: A benchmark for multi-llm routing system. *arXiv preprint arXiv:2403.12031*, 2024.
- [33] Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma. Look before you leap: An exploratory study of uncertainty measurement for large language models. *arXiv preprint arXiv:2307.10236*, 2023.

- [34] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [35] Mojan Javaheripi and Sébastien Bubeck. Phi-2: The surprising power of small language models, December 2023.
- [36] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. *arXiv*, 2023.
- [37] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- [38] Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. How can we know when language models know? on the calibration of language models for question answering. *Transactions of the Association for Computational Linguistics*, 9:962–977, 2021.
- [39] Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- [40] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*, 2023.
- [41] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. 2023.
- [42] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. In *Advances in Neural Information Processing Systems*, 2024.
- [43] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
- [44] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:87–100, 2024.
- [45] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- [46] Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words. *arXiv preprint arXiv:2205.14334*, 2022.
- [47] Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantification for black-box large language models. *arXiv preprint arXiv:2305.19187*, 2023.
- [48] Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. *ACL*, 2017.
- [49] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [50] Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv preprint arXiv:2402.09353*, 2024.

- [51] Zechun Liu, Changsheng Zhao, Forrest Iandola, Chen Lai, Yuandong Tian, Igor Fedorov, Yunyang Xiong, Ernie Chang, Yangyang Shi, Raghuraman Krishnamoorthi, et al. Mobilellm: Optimizing sub-billion parameter language models for on-device use cases. *arXiv preprint arXiv:2402.14905*, 2024.
- [52] Zhenyan Lu, Xiang Li, Dongqi Cai, Rongjie Yi, Fangming Liu, Xiwen Zhang, Nicholas D Lane, and Mengwei Xu. Small language models: Survey, measurements, and insights. *arXiv preprint arXiv:2409.15790*, 2024.
- [53] Matéo Mahaut, Laura Aina, Paula Czarnowska, Momchil Hardalov, Thomas Müller, and Lluís Màrquez. Factual confidence of llms: on reliability and robustness of current estimators. *arXiv preprint arXiv:2406.13415*, 2024.
- [54] Potsawee Manakul, Adian Liusie, and Mark Gales. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017, 2023.
- [55] AI Meta. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. *Meta AI Blog*. Retrieved December, 20:2024, 2024.
- [56] Sabrina J Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. Reducing conversational agents’ overconfidence through linguistic calibration. *Transactions of the Association for Computational Linguistics*, 10:857–872, 2022.
- [57] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018.
- [58] Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, Yuling Gu, Shane Arora, Akshita Bhagia, Dustin Schwenk, David Wadden, Alexander Wettig, Binyuan Hui, Tim Dettmers, Douwe Kiela, Ali Farhadi, Noah A. Smith, Pang Wei Koh, Amanpreet Singh, and Hannaneh Hajishirzi. Olmoe: Open mixture-of-experts language models, 2024.
- [59] Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez, M Waleed Kadous, and Ion Stoica. Routellm: Learning to route llms with preference data. *arXiv preprint arXiv:2406.18665*, 2024.
- [60] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, 2021.
- [61] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, et al. Rwkv: Reinventing rnns for the transformer era. *arXiv preprint arXiv:2305.13048*, 2023.
- [62] Siva Reddy, Danqi Chen, and Christopher D. Manning. CoQA: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266, 2019.
- [63] Subhro Roy and Dan Roth. Solving general arithmetic word problems. *arXiv preprint arXiv:1608.01413*, 2016.
- [64] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- [65] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social iqa: Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, 2019.

- [66] Dimitris Stripelis, Zijian Hu, Jipeng Zhang, Zhaozhuo Xu, Alay Shah, Han Jin, Yuhang Yao, Salman Avestimehr, and Chaoyang He. Polyrouter: A multi-llm querying system. *arXiv preprint arXiv:2408.12320*, 2024.
- [67] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158. Association for Computational Linguistics, 2019.
- [68] Omkar Thawakar, Ashmal Vayani, Salman Khan, Hisham Cholakkal, Rao Muhammad Anwer, Michael Felsberg, Timothy Baldwin, Eric P. Xing, and Fahad Shahbaz Khan. Mobillama: Towards accurate and lightweight fully transparent gpt, 2024.
- [69] Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. *arXiv preprint arXiv:2305.14975*, 2023.
- [70] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [71] Artem Vazhentsev, Akim Tsvigun, Roman Vashurin, Sergey Petrakov, Daniil Vasilev, Maxim Panov, Alexander Panchenko, and Artem Shelmanov. Efficient out-of-domain detection for sequence to sequence models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1430–1454, 2023.
- [72] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in neural information processing systems*, pages 24824–24837, 2022.
- [73] Gwenth Portillo Wightman, Alexandra Delucia, and Mark Dredze. Strength in numbers: Estimating confidence of large language models by prompt agreement. In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, pages 326–362, 2023.
- [74] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR, 2023.
- [75] Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. *arXiv preprint arXiv:2306.13063*, 2023.
- [76] Miao Xiong, Zhiyuan Hu, Xinyang Lu, YIFEI LI, Jie Fu, Junxian He, and Bryan Hooi. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. In *The Twelfth International Conference on Learning Representations*, 2024.
- [77] Haoran Xu, Baolin Peng, Hany Awadalla, Dongdong Chen, Yen-Chun Chen, Mei Gao, Young Jin Kim, Yunsheng Li, Liliang Ren, Yelong Shen, et al. Phi-4-mini-reasoning: Exploring the limits of small reasoning language models in math. *arXiv preprint arXiv:2504.21233*, 2025.
- [78] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.

- [79] Gal Yona, Roei Aharoni, and Mor Geva. Can large language models faithfully express their intrinsic uncertainty in words? *arXiv preprint arXiv:2405.16908*, 2024.
- [80] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.
- [81] Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*, 2024.
- [82] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [83] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [84] Zesen Zhao, Shuowei Jin, and Z Morley Mao. Eagle: Efficient training-free router for multi-llm inference. *arXiv preprint arXiv:2409.15518*, 2024.
- [85] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems*, pages 46595–46623, 2023.
- [86] Zhengping Zhou, Lezhi Li, Xixi Chen, and Andy Li. Mini-giants: "small" language models and open source win-win. *arXiv preprint arXiv:2307.08189*, 2023.

A Challenges and Opportunities

❶ **How to cash-in routing efficiency on new edge devices?** Based on the benchmark results, proxy routing data provides a robust foundation for establishing routing policies on new edge devices without accessing prior knowledge at the early stage of deployment. This enables the routing policies with strong generalization to new dataset scenarios and enhances the efficiency across diverse deployments for personal edge devices. While proxy routing data holds a good performance in the early deployment stage, an important direction to explore is how to effectively leverage additional private on-device data to strengthen the quality of proxy routing data, aiming to continuously enhance the deployment of personalized routing strategies. With the aid of proxy routing data, less private data is required, but striking a balance between privacy and performance remains an open challenge.

❷ **How to effectively strike a balance between LLM routing efficiency and utility?** We empirically observe that by leveraging UQ methods with strong uncertainty-utility alignment (e.g., Perplexity and OOD Probe methods), routing thresholds can effectively be determined with the sweet points of efficiency and utility. However, achieving such sweet spots can be challenging due to the variability in downstream datasets and the sensitivity of UQ methods to LLM-specific characteristics. Additionally, discrepancies across different device types, such as variations between iOS and Android systems, further complicating the process, requiring tailored strategies and analytics to account for platform-specific constraints and capabilities. Based on these factors, providing a fair apple-to-apple comparison regarding routing performance is inherently challenging. Researchers should be mindful of these complexities and focus on developing methods that are not only efficient but also capable of handling long-context scenarios effectively.

❸ **How is the performance when conducting compression (e.g., pruning, quantization) on the on-device model?** As with the on-device models discussed in the above sections, we directly adopt a pre-trained small model without any modifications. Alternatively, on-device models can also be generated by compressing larger models. Specifically, numerous works have explored methods for compressing LLMs into smaller sizes using techniques such as pruning [22] and quantization [74, 44]. The advantage of employing compression methods is that the smaller models compressed from larger ones tend to retain similar distributions of the output, thereby mitigating the issue of distribution shift.

❹ **Uncertainty-aware routing in on-device multimodal language models.** While LLMs typically operate with a single modality for both input and output, a promising research direction involves exploring uncertainty-aware routing in multimodal language models (MLLMs). For instance, in vision-language models (VLMs) such as LLaVa [49] and InternVL [9], the inputs include both images/videos and text. By incorporating visual modalities, the properties of vision tokens significantly influence the output. As a result, the uncertainty in the generated text differs from that of language-only models. Benchmarking and generalizing uncertainty-aware routing for on-device MLLMs is a valuable direction for the research community.

B Details about Datasets

The details of the 15 datasets are further listed in Tabel 2. We applied the original dataset directly from the Huggingface dataset repositories without any further processing. A thorough examination of each dataset’s attributes, size, and notable characteristics is provided below.

C Related Work

C.1 Small Language Models

Small Language Models (SLMs) are designed for deployment on resource-constrained devices like desktops, smartphones, and wearables. Specifically, we consider the Transformer-based SLMs in this work due to their state-of-the-art performance, like Phi-3-mini [1], TinyLlama [81], MobileLLM [51], and Qwen-1.5B [5], LiteLLaMa-460M, OPT-125M [82], BLOOMZ (560M, 1.1B, 1.7B, 3B) [41], SmoLLM (135M, 360M, 1.7B) [3], OLMo (1B) [28], OLMoE (1B) [58], MobileLlama (0.5B, 1B) [68], MobileLLaMA (1.4B, 2.7B) [10], OpenLLaMA (3B) [25]. These models are designed with lightweight architectures to operate effectively within the computational and storage limitations of mobile devices and edge hardware.

Table 2: Details of the 15 datasets used in our benchmark. FF: Free-form question answering (including numerical answers for math tasks); MCQ: Multiple-choice question answering; TF: True/False question answering.

Dataset	Type	Domain	# Train	# Test	Description
GSM8K	FF	Mathematical Reasoning	7473	1319	Grade school math word problems
AQuA	MCQ	Mathematical Reasoning	97467	254	Algebraic word problems
MultiArith	FF	Mathematical Reasoning	420	180	Algebraic word problems
SVAMP	FF	Mathematical Reasoning	700	300	Algebraic word problems
MATH-500	FF	Mathematical Reasoning	—	500	Algebraic word problems
BoolQ	TF	Commonsense Reasoning	9427	3270	Commonsense and factual reasoning questions
CommonsenseQA	MCQ	Commonsense Reasoning	9741	1221	Questions assessing various types of commonsense knowledge
HellaSwag	MCQ	Commonsense Reasoning	39905	10042	Sentence completion based on narrative understanding
OpenBookQA	MCQ	Commonsense Reasoning	4957	500	Open-book science and commonsense questions
PIQA	MCQ	Commonsense Reasoning	16113	1838	Physical commonsense reasoning questions
Social IQa	MCQ	Commonsense Reasoning	33410	1954	Social commonsense intelligence questions
TruthfulQA	FF	Commonsense Reasoning	653	164	Assessing models' ability to prevent false information
WinoGrande	MCQ	Commonsense Reasoning	2558	1267	Pronoun ambiguity resolution with commonsense reasoning
CoQA	FF	Conversational & Contextual Understanding	7199	500	Conversational questions on text passages from diverse domains
MMLU	MCQ	Problem Solving	99842	14042	Problem solving across various subjects

Recurrent Neural Networks (RNNs), like RWKV (1B, 3B, 7B) [61], Mamba (1.4B, 6.9B) [14], and RecurrentGemma-2B [27], can provide promising solutions for on-device inference in resource-constrained environments. These models leverage the recurrent nature of RNNs to process sequential data efficiently without requiring a KV cache, which is suitable for resource-constrained on edge devices. Specifically, RWKV introduces a hybrid RNN-Transformer backbone to capture long-term dependencies while maintaining computational efficiency. Similarly, Mamba and RecurrentGemma design recurrent layers for low-power consumption and high throughput inference, which can significantly reduce memory and computational requirements, fostering low-latency applications directly on devices.

D Additional Experimental Results from Benchmarking to Generalization

In this section, we present additional experimental results on (1) evaluating the impact of uncertainty-correctness alignment on small language model (SLM) routing and (2) investigating the generalization capability of proxy routing data on novel datasets. **Since our studies yield over 5,600 results, we here present a representative subset in the following section. The full set of results is provided in the supplementary materials.**

For the first experiment (Section D.1 and Section D.2), we provide the complete set of results, including the AUC measurements for uncertainty-correctness alignment and the performance of uncertainty-based routing. For the second experiment (Section D.3), we present a comprehensive experimental results of proxy routing prediction under partially in-domain setting. Each dataset referenced in the experiments is treated as a novel dataset for evaluation.

D.1 Evaluation on Uncertainty-correctness Alignment

Results of Alignment between uncertainty and correctness.

All the experiments shown on this page are conducted under AQUA, BoolQ, and CoQA datasets with all 8 UQ methods.

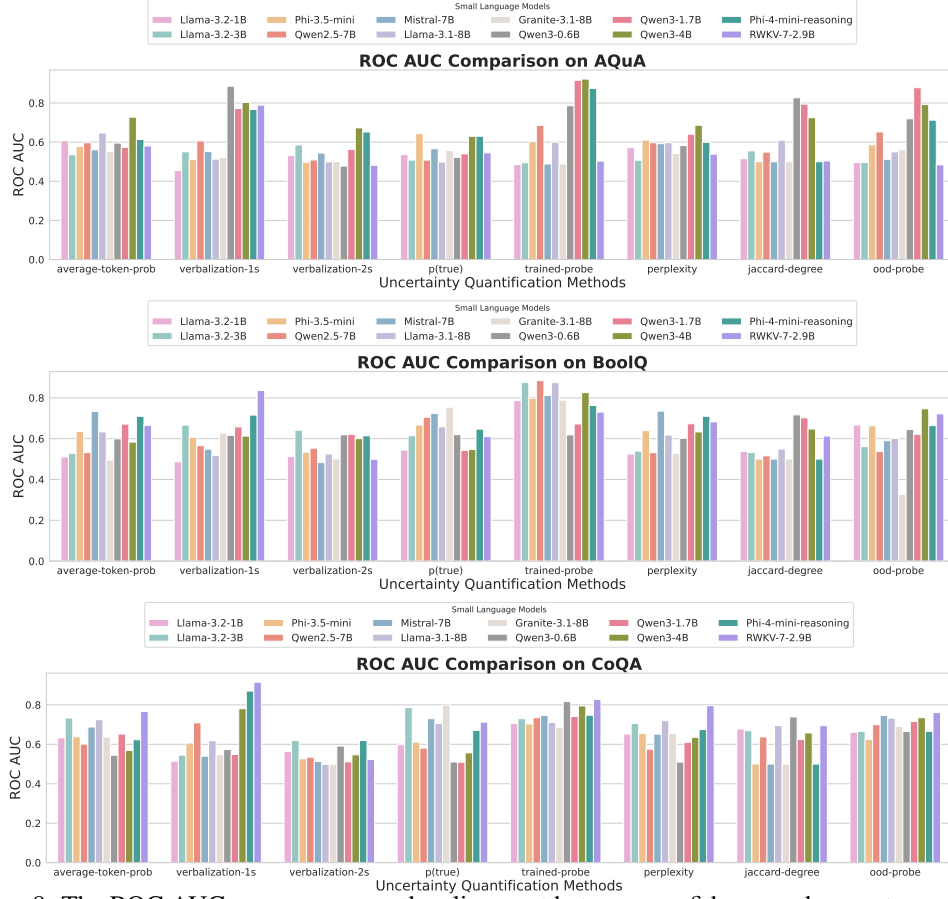


Figure 8: The ROC AUC scores measure the alignment between confidence and correctness across different SLMs and uncertainty quantification methods on AQUA, BoolQ, and CoQA. A higher ROC AUC indicates a stronger alignment.

Results of Alignment between uncertainty and correctness.

All the experiments shown on this page are conducted under GSM8K, HellaSwag, MMLU, and MultiArith datasets with all 8 UQ methods.

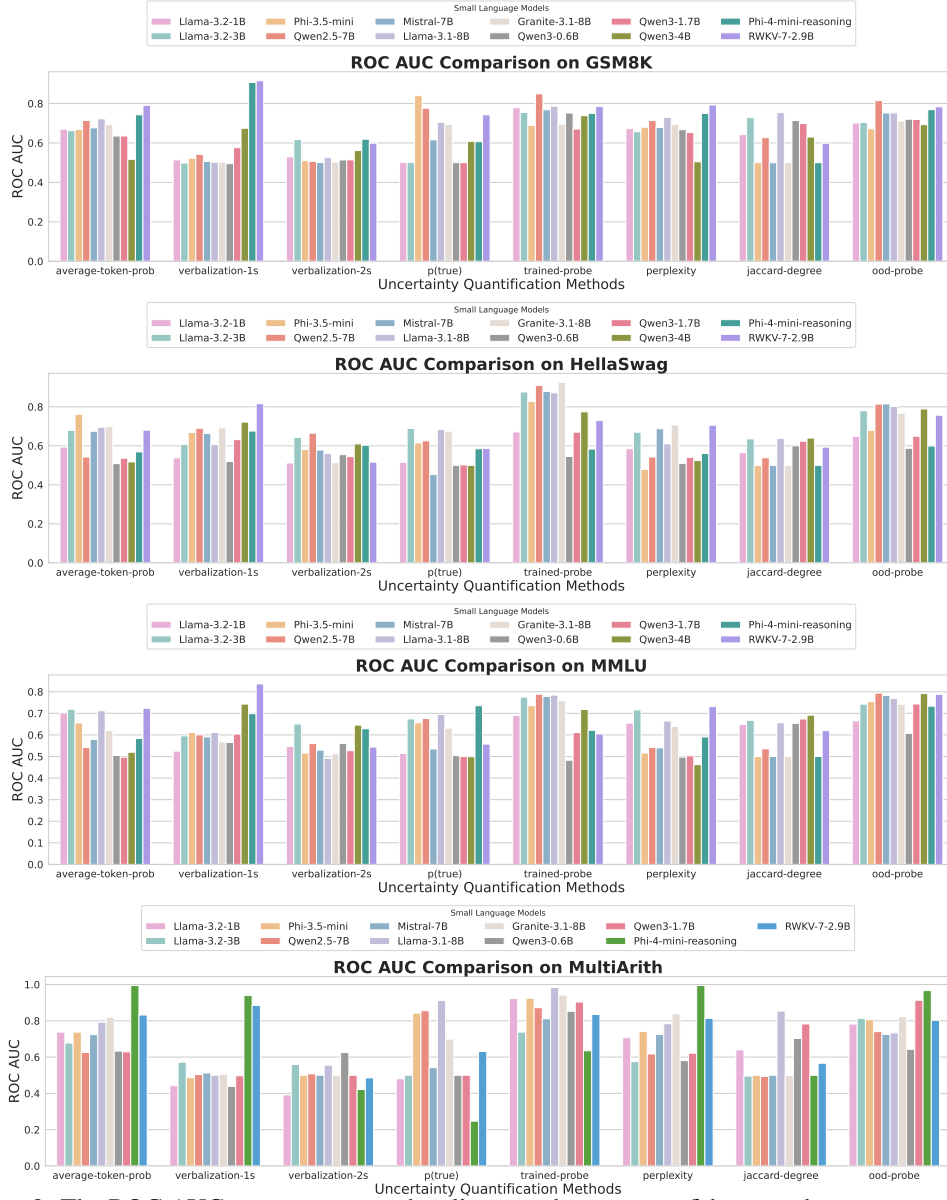


Figure 9: The ROC AUC scores measure the alignment between confidence and correctness across different SLMs and uncertainty quantification methods on GSM8K, HellaSwag, MMLU, and Multi-Arith. A higher ROC AUC indicates a stronger alignment.

Results of Alignment between uncertainty and correctness.

All the experiments shown on this page are conducted under OpenBookQA, PIQA, SocialIQA, and SVAMP datasets with all 8 UQ methods.

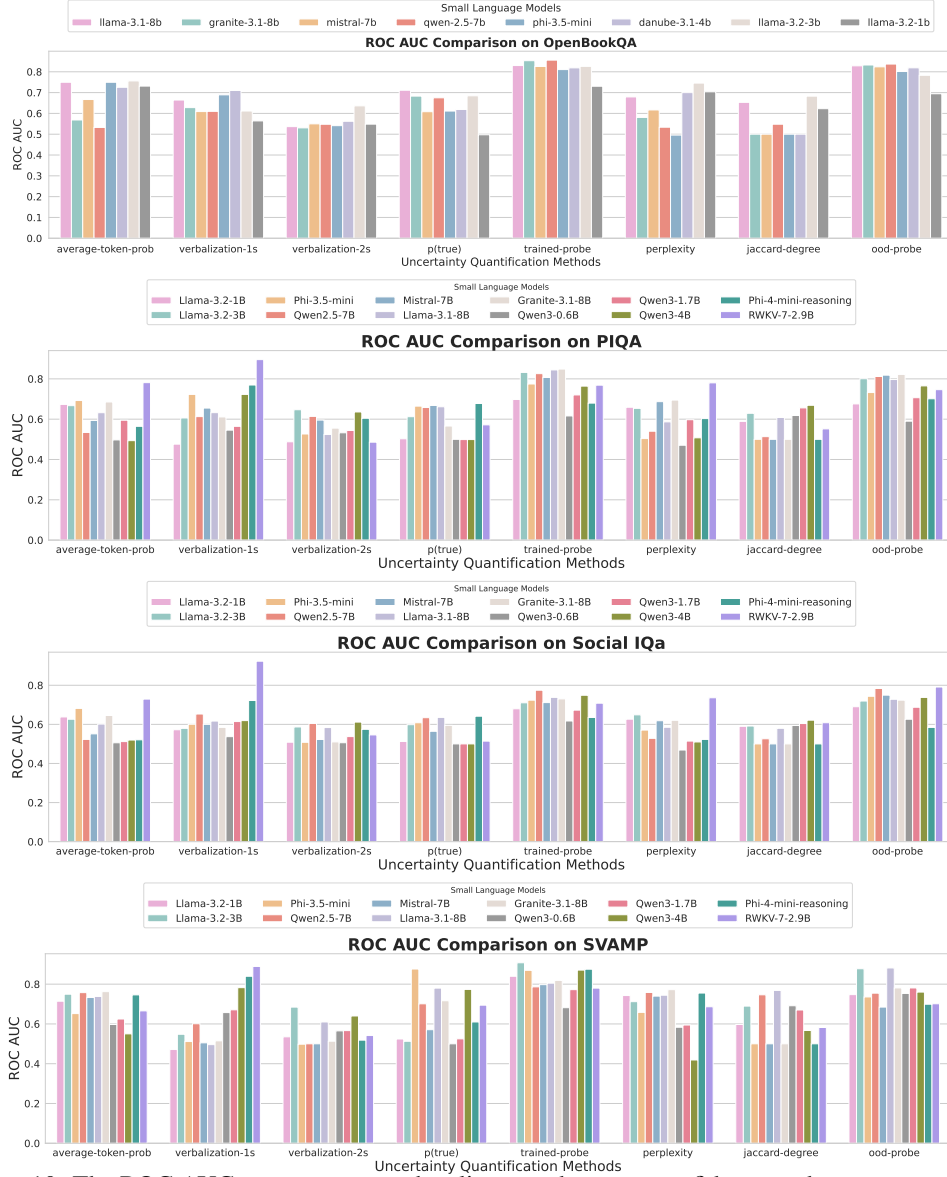


Figure 10: The ROC AUC scores measure the alignment between confidence and correctness across different SLMs and uncertainty quantification methods on OpenBookQA, PIQA, SocialIQA, and SVAMP datasets. A higher ROC AUC indicates a stronger alignment.

Results of Alignment between uncertainty and correctness.

All the experiments shown on this page are conducted under CommonsenseQA, SVAMP, TruthfulQA, and Math500 dataset with all 8 UQ methods.

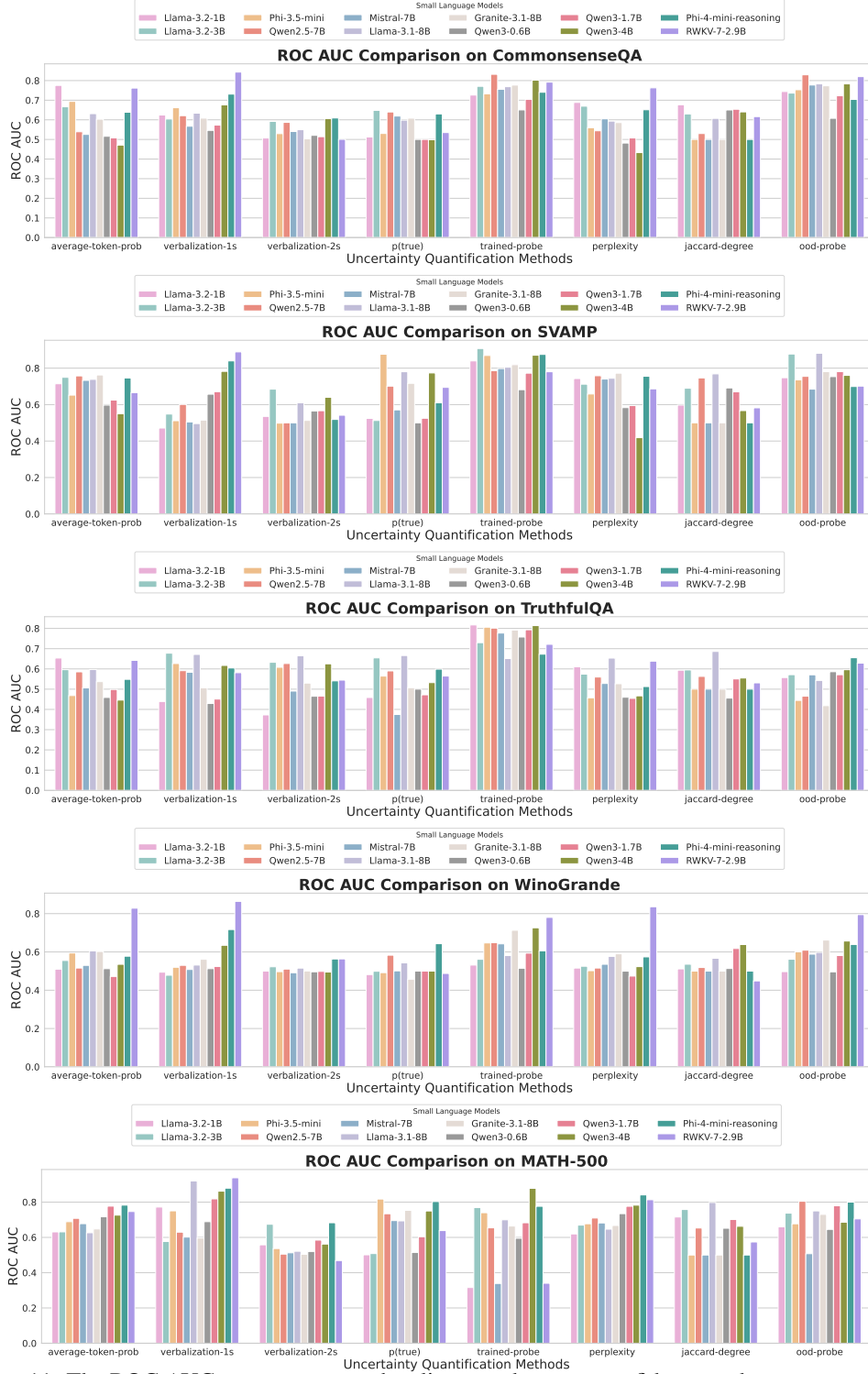


Figure 11: The ROC AUC scores measure the alignment between confidence and correctness across different SLMs and uncertainty quantification methods on CommonsenseQA. A higher ROC AUC indicates a stronger alignment.

D.2 Evaluation on Uncertainty-based Routing Approaches

Results of routing to GPT-4.1-Mini

All the experiments shown on this page are conducted under all benchmark datasets with selected SLMs. We only showcase partial of the experimental results.

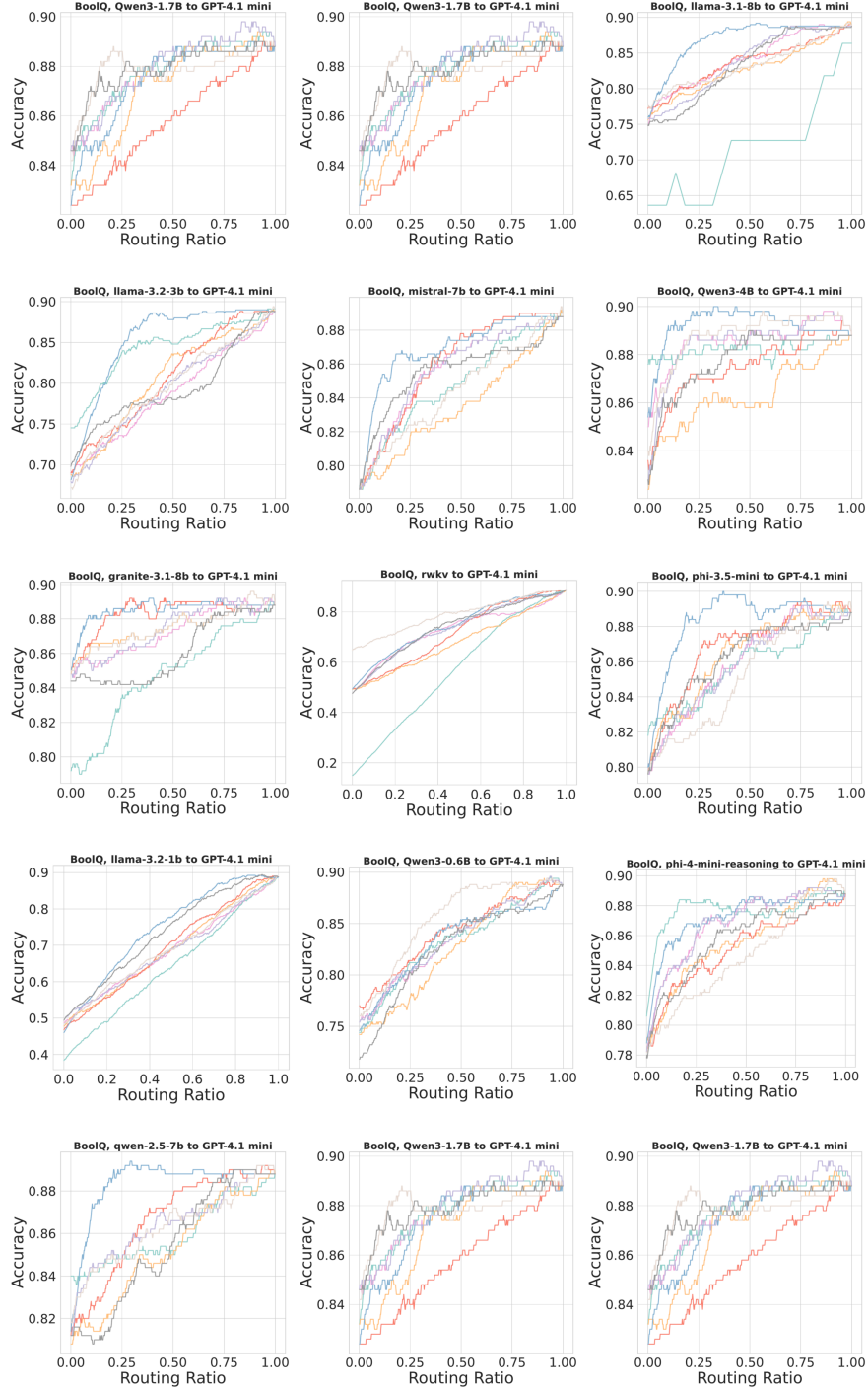


Figure 12: Overall accuracy vs. routing ratio with different UQ methods and SLMs.

Results of routing to DeepSeek-R1

All the experiments shown on this page are conducted under all benchmark datasets with selected SLMs. We only showcase partial of the experimental results.

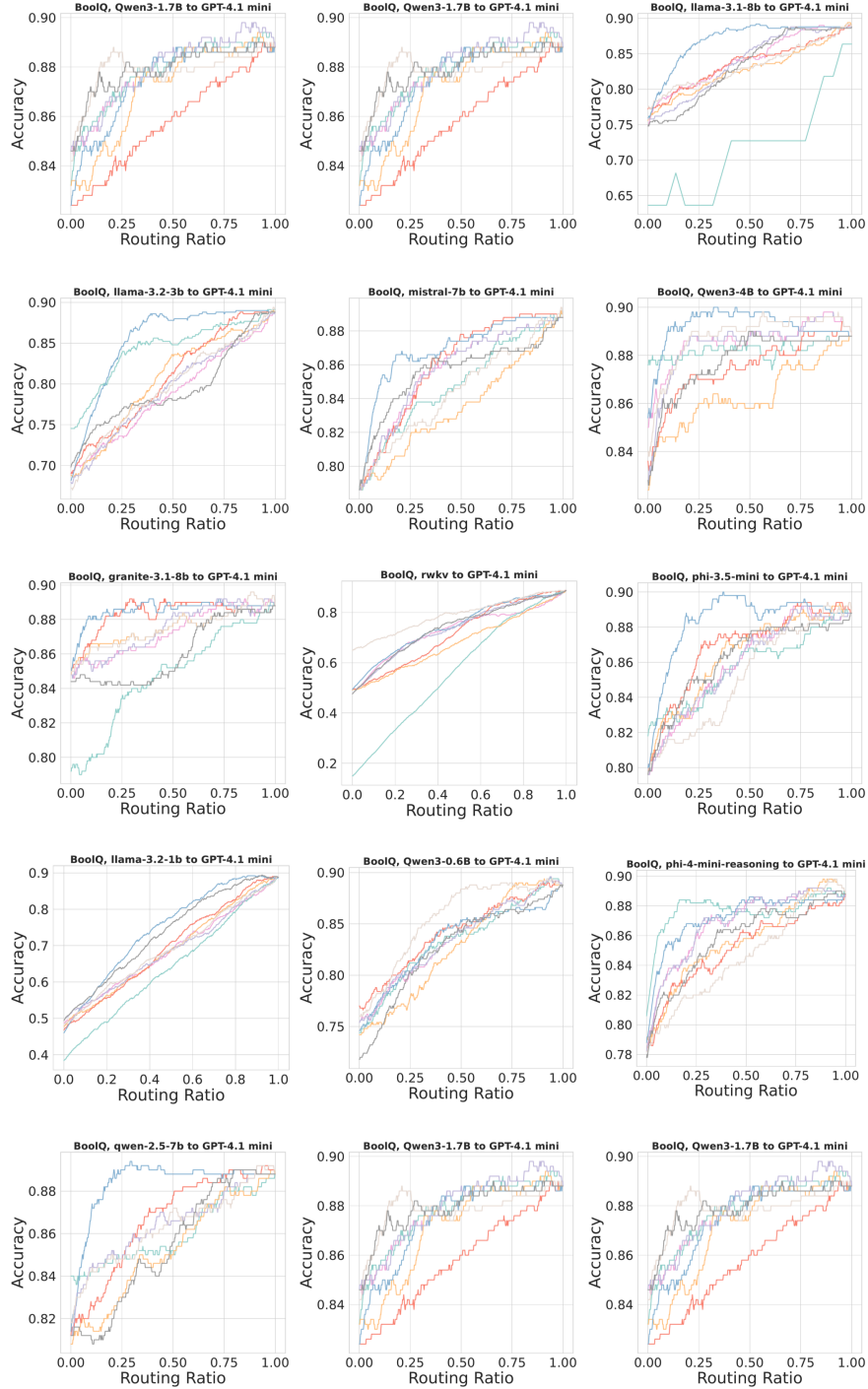


Figure 13: Overall accuracy vs. routing ratio with different UQ methods and SLMs.

Results of routing to Llama-3.1-70B-Instruct

All the experiments shown on this page are conducted under all benchmark datasets with selected SLMs. We only showcase partial of the experimental results.

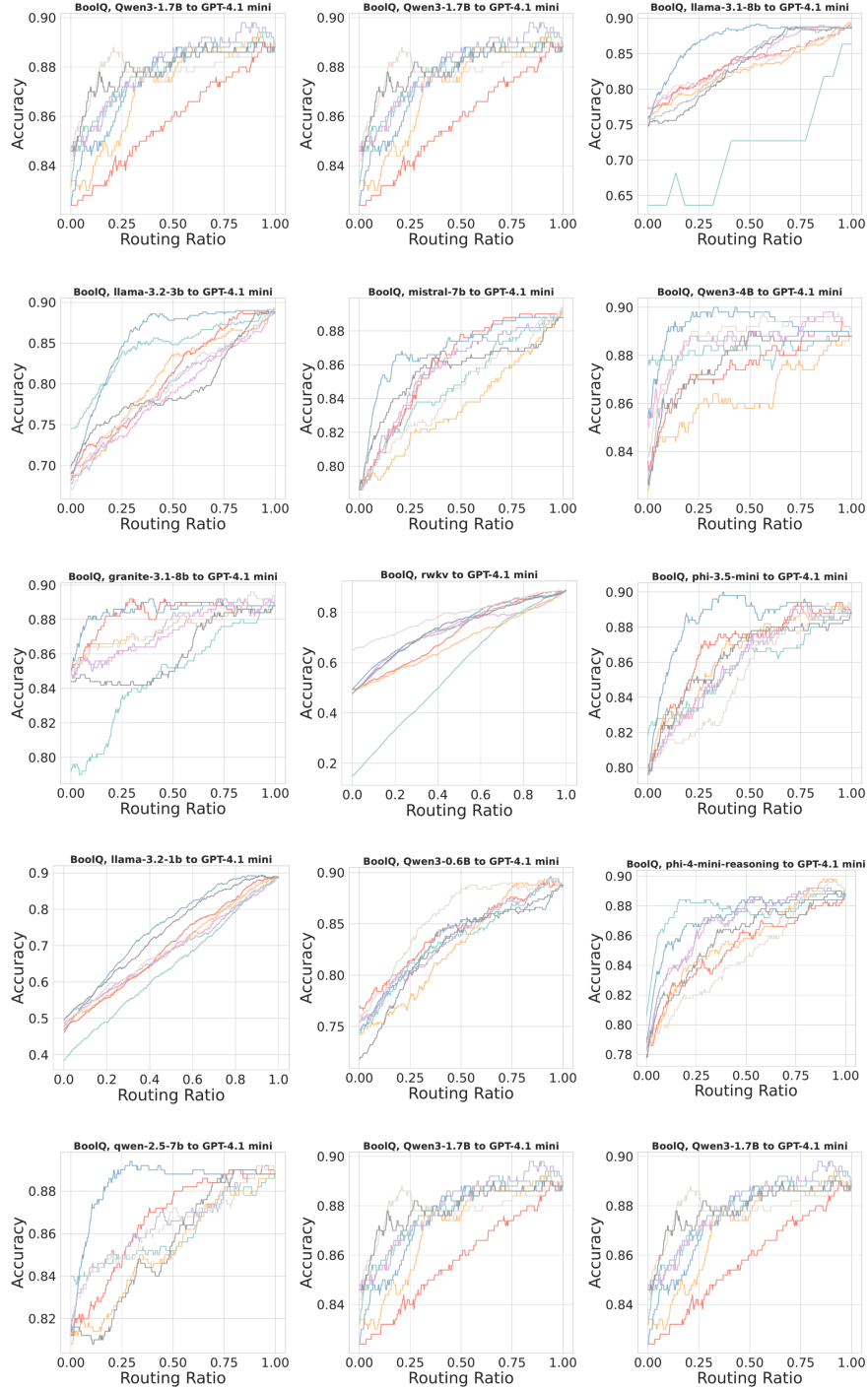


Figure 14: Overall accuracy vs. routing ratio with different UQ methods and SLMs.

Results of routing to Qwen3-32B

All the experiments shown on this page are conducted under all benchmark datasets with selected SLMs. We only showcase partial of the experimental results.

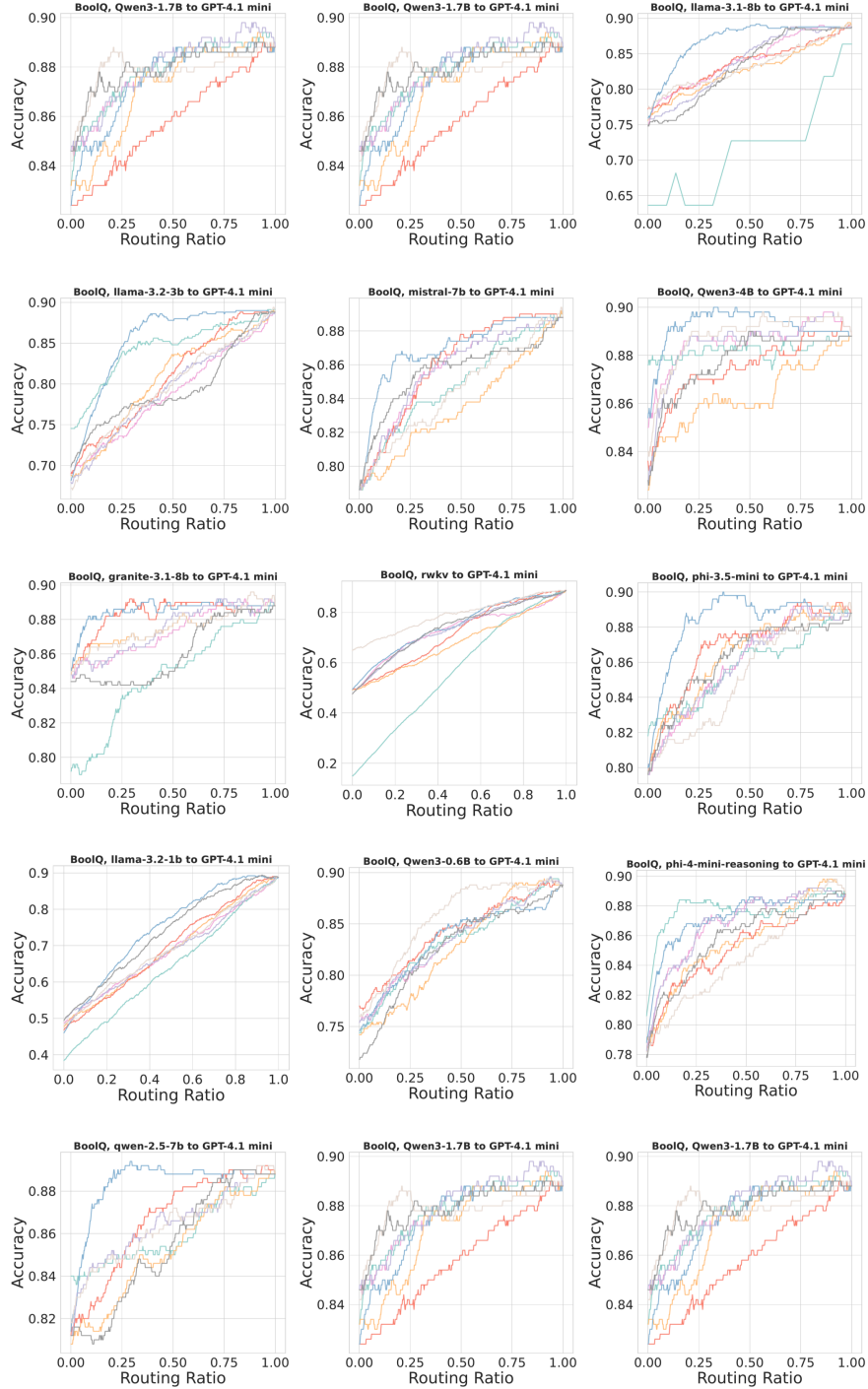


Figure 15: Overall accuracy vs. routing ratio with different UQ methods and SLMs.

D.3 Evaluation of Proxy Routing Data on New Downstream Scenario

Evaluation results on routing to GPT-4.1-Mini

The experiments shown on this page are conducted under all 15 datasets with different SLMs.

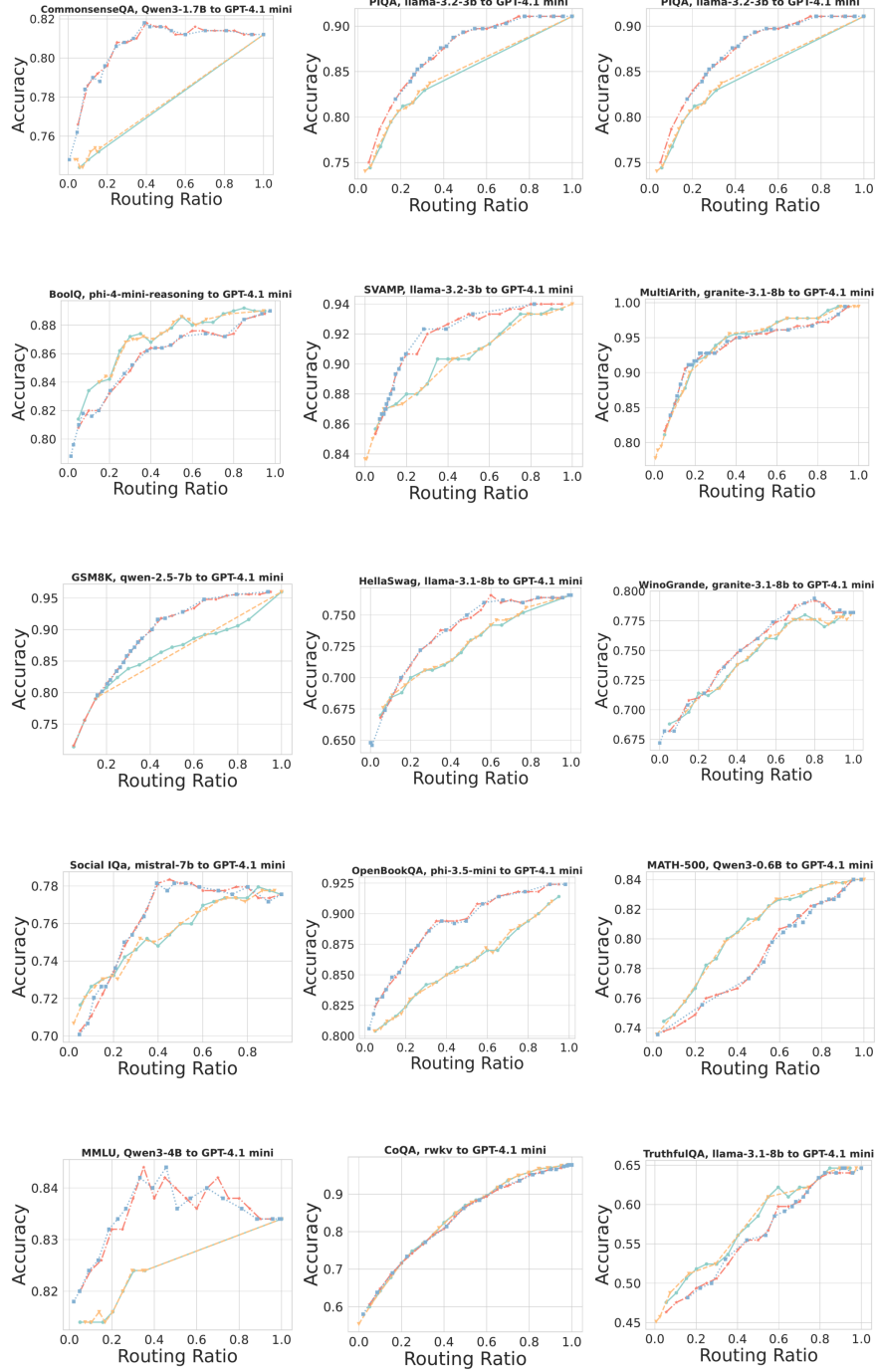


Figure 16: Assessing the generalization of proxy routing data to new downstream data for routing 12 SLMs to GPT-4.1-Mini on 15 datasets using two UQ methods (OOD Probe & Perplexity). The legend in Figure 6 is also used here.

Evaluation results on routing to DeepSeek-R1

The experiments shown are conducted under all math reasoning datasets with different SLMs.

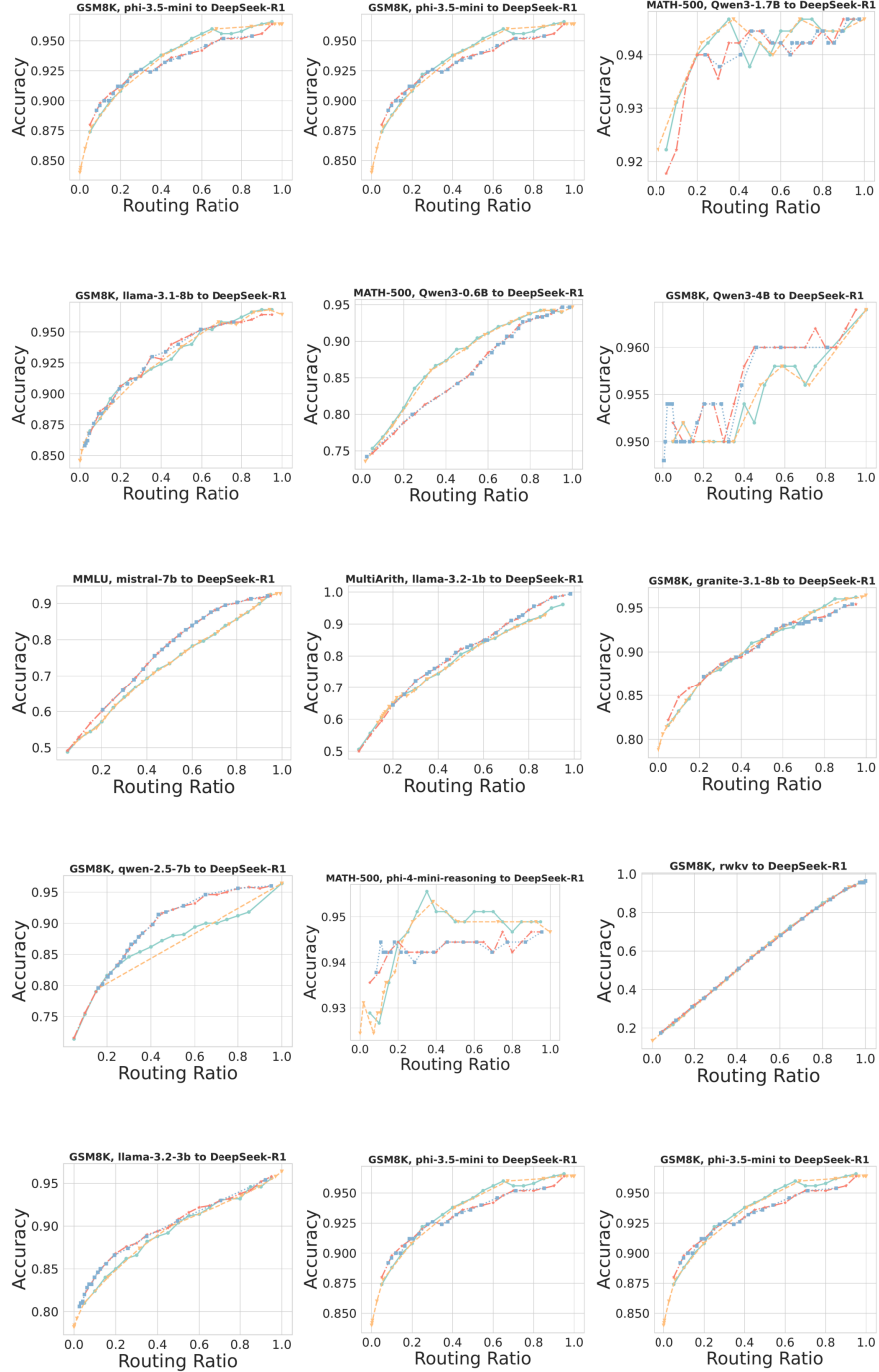


Figure 17: Assessing the generalization of proxy routing data to new downstream data for routing 12 SLMs to DeepSeek-R1 on 15 datasets using two UQ methods (OOD Probe & Perplexity). The legend in Figure 6 is also used here.

Evaluation results on routing to Llama-3.1-70B-Instruct

The experiments shown on this page are conducted under all 15 datasets with different SLMs.

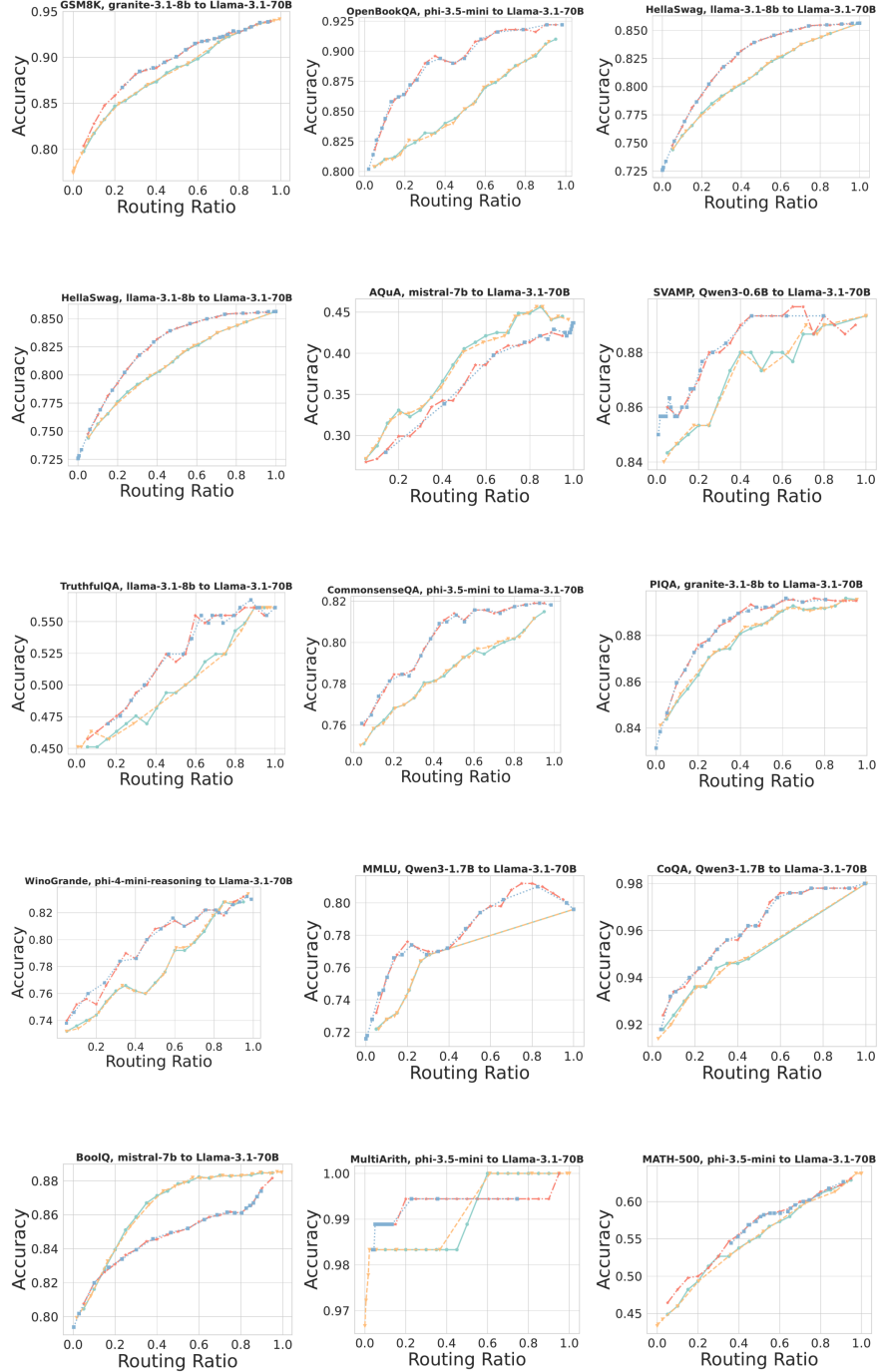


Figure 18: Assessing the generalization of proxy routing data to new downstream data for routing 12 SLMs to Llama-3.1-70B-Instruct on 15 datasets using two UQ methods (OOD Probe & Perplexity). The legend in Figure 6 is also used here.

Evaluation results on routing to Qwen3-32B

The experiments shown on this page are conducted under all 15 datasets with different SLMs.

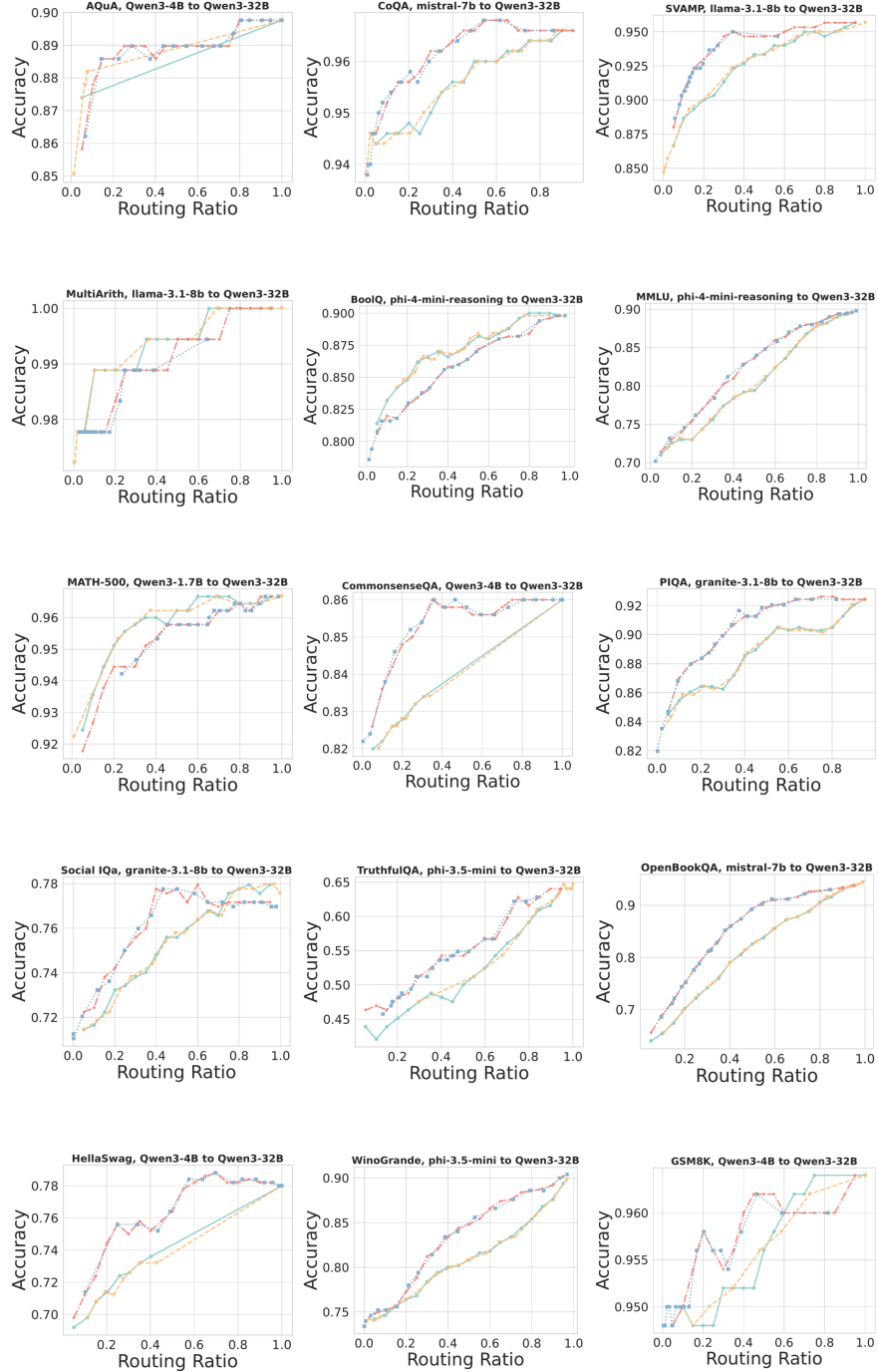


Figure 19: Assessing the generalization of proxy routing data to new downstream data for routing 12 SLMs to Qwen3-32B on 15 datasets using two UQ methods (OOD Probe & Perplexity). The legend in Figure 6 is also used here.

E Sensitivity to Sampling Ratio

We conducted a sensitivity analysis of the sampling ratio to assess our method’s generalizability and robustness. As shown in Figure 20, we varied the sampling ratio from 0.01 to 0.3 and computed the RMS distance between the oracle routing curve derived from the full downstream dataset and the curves obtained at different sampling ratios. We observe that our method is highly stable and robust to the choice of sampling ratio. Similar results hold across other datasets and routing scenarios.

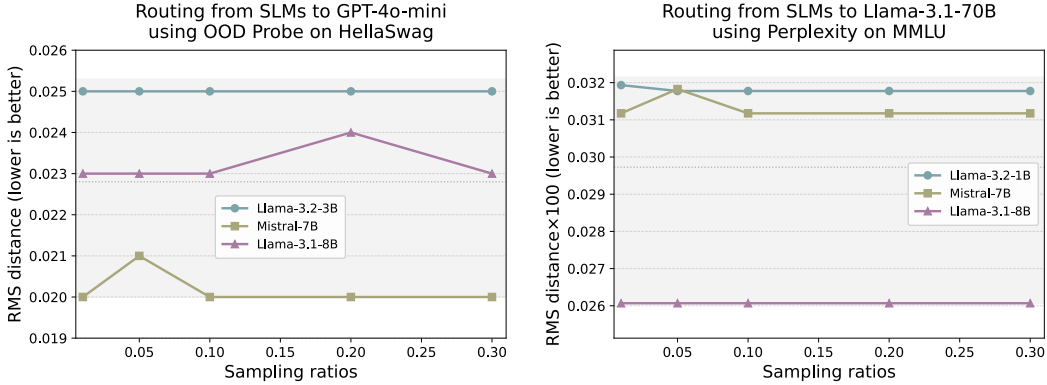


Figure 20: RMS distance versus sampling ratio across two routing scenarios.

F Computational Infrastructure

The computational infrastructure information is given in Table 3.

Table 3: Experiment configuration and computing infrastructure.

Name	Value
Data type	torch.bfloat16
Flash-Attention	True
Computing Infrastructure	GPU
GPU Model	NVIDIA-A100
GPU Memory	80GB
GPU Number	4
CUDA Version	12.1
CPU Memory	512GB

G Routing with Proxy Routing Data vs. Random Routing

Our approach sets a confidence threshold for routing on a new dataset without prior access, and no existing work has addressed this scenario to the best of knowledge. To quantify its effectiveness we compare it with a *random routing* baseline. On HellaSwag, we route three SLMs to GPT-4.1 mini using Perplexity. We compute the average root mean squared (RMS) distance between the oracle routing curve derived from the full downstream dataset and those obtained using either our proposed method or random routing. A lower RMS distance indicates closer alignment with the oracle and therefore better routing quality. As shown in Table 4, routing based on proxy routing data dramatically outperforms random routing. Similar gains are observed with other uncertainty quantification methods (e.g., a 34.14% improvement with OOD Probe).

Robustness under strong OOD shifts. We further test our approach based on proxy routing data in two challenging out-of-distribution (OOD) scenarios:

Table 4: Routing with proxy routing data vs. random routing on HellaSwag. Lower RMS is better.

Method	LLAMA-3.2-3B	MISTRAL-7B	LLAMA-3.1-8B
Ours	0.001	0.001	0.001
Random routing	0.029	0.031	0.019

- **Math** \rightarrow **Commonsense**: Proxy routing data drawn solely from math datasets (GSM8K, AQuA, MultiArith, SVAMP); evaluation on commonsense reasoning (TruthfulQA), routing various SLMs to Llama-3.1-70B.
- **Commonsense** \rightarrow **Math**: Proxy data drawn solely from commonsense datasets; evaluation on the math dataset AQuA with the same routing setup.

Table 5: Routing with proxy routing data vs. random routing under strong OOD shifts.

OOD Setting	Method	PHI-3.5-MINI	MISTRAL-7B	LLAMA-3.1-8B
Math \rightarrow Commonsense	Ours	0.0148	0.0048	0.0132
	Random routing	0.0187	0.0090	0.0176
Commonsense \rightarrow Math	Ours	0.0057	0.0060	0.0022
	Random routing	0.0085	0.0082	0.0037

Across both OOD shifts (Table 5), routing with proxy routing data consistently yields smaller RMS distances than the random baseline, underscoring its strong generalization capability.