
Confident or Seek Stronger: Exploring Uncertainty-Based Small LM Routing

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Small language models (SLMs) are increasingly deployed on edge devices for
2 personalized applications, offering low-latency inference and reduced energy con-
3 sumption. However, they often struggle with complex queries, leading to unre-
4 liable responses. Uncertainty-based SLM routing addresses this by offloading
5 low-confidence queries to stronger large language models (LLMs), following the
6 principle “if uncertain, seek stronger support” to improve reliability. While leverag-
7 ing LLMs enhances accuracy, it also incurs high invocation costs, making it crucial
8 to balance efficiency and efficacy. In this paper, we conduct a comprehensive inves-
9 tigation into *benchmarking of uncertainty-driven routing strategies from SLMs to*
10 *LLMs over 5000+ settings*. Our findings highlight: *First*, uncertainty-correctness
11 alignment in different uncertainty quantification (UQ) methods significantly im-
12 pacts routing performance. The extracted uncertainty distribution is primarily
13 influenced by the selected SLM and uncertainty quantification (UQ) method, show-
14 ing minimal dependence on the downstream dataset. The source code is available
15 at <https://anonymous.4open.science/r/quodlibeta>.

16 1 Introduction

17 Large language models (LLMs) have gained increasing attention for deployment on edge devices due
18 to their potential for low-latency and privacy-preserving inference. However, given the computational
19 and memory constraints of edge hardware, small language models (SLMs) (e.g., Phi2-mini [33],
20 Llama3.2-3B [64]) are designed for resource-efficient deployment on smartphones, wearables, and
21 similar devices. Their overarching goal is to democratize LM deployment, making it accessible and
22 affordable across diverse settings [46, 78, 76]. Despite this, SLMs often lack the robustness and
23 scalability of LLMs [8] (e.g., GPT-4o [2], Llama-3.1-405B), especially under diverse and complex
24 queries in edge deployments, leading to degraded overall performance and motivating the need for
25 improved response reliability.

26 To address this limitation, recent work proposes partial offloading of challenging queries from SLMs
27 to stronger LLMs [10, 53, 30, 60], forming hybrid systems that intelligently route queries for more
28 reliable and deterministic responses. However, these methods face significant challenges when
29 encountering new downstream tasks, as the data fall outside the distribution of the existing training
30 data, making them less practical for real-world scenarios, such as in personal edge device deployment,
31 where adaptability to unseen conditions is crucial. While LLMs achieve superior performance, their
32 deployment incurs high inference and infrastructure costs, e.g., a single NVIDIA A100 GPU can
33 cost \sim \$2,000 per month. Moreover, inaccurate routing by SLMs can unnecessarily increase LLM
34 query traffic, requiring greater bandwidth and further raising operational expenses, particularly under
35 continuous deployment scenarios. Therefore, developing an effective routing strategy is crucial for
36 fully leveraging SLMs [53, 60, 10], as it enhances response reliability while reducing service and
37 data transmission costs.

Leveraging SLMs’ self-uncertainty estimation emerges as a robust strategy for enhancing routing effectiveness [10, 15]. 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:

What is the best practice of uncertainty estimation for query routing from SLMs to LLMs? 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 [31, 14, 67], 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 strategies¹. Our contributions include 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 uncertainty-based SLM routing.
- **Uncertainty distribution characterization:** The extracted uncertainty distribution is determined by the choice of SLM and UQ method, exhibiting minimal dependence on the downstream dataset.

2 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 2) 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.

2.1 Benchmark Coverage and Setup

The 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 [49], Llama-3.2-3B-Instruct [49], Phi-3.5-mini-instruct [1], Mistral-7B-Instruct-v0.3 [34], Qwen2.5-7B-Instruct [71], Llama-3.1-8B-Instruct [18], and Granite-3.1-8B-Instruct [24]. The reasoning models are Qwen3-0.6B, Qwen3-1.7B, Qwen3-4B, and Phi-4-mini-reasoning [70]. The RNN model is RWKV-7-2.9B [55]. We also include four LLMs: three open-source models—Llama-3.1-70B-Instruct [18], Qwen3-32B, and DeepSeek-R1 [27]—and one proprietary API model, GPT-4.1 mini [32]. Qwen3-32B and DeepSeek-R1 are reasoning LLMs, whereas Llama-3.1-70B and GPT-4.1 mini are non-reasoning.

Datasets and Experiment Settings. Experiments span 15 datasets from four domains: (1) *Mathematical Reasoning* (AQuA [44], GSM8K [12], MultiArith [57], SVAMP [54], MATH-500 [40]), (2) *Commonsense Reasoning* (CommonsenseQA [61], HellaSwag [73], OpenBookQA [51], PIQA [6], TruthfulQA [41], WinoGrande [58], BoolQ [11], Social IQa [59]), (3) *Conversational and Contextual Understanding* (CoQA [56]), and (4) *Problem Solving* (MMLU [28]). These cover free-form, multiple-choice, and True/False question answering and are available via Hugging Face. Appendix A provides further details about the settings of datasets and selected UQ methods.

2.2 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 C.1.

¹For the convenience of writing, we interchangeably use uncertainty and confidence, where low uncertainty refers to high confidence.

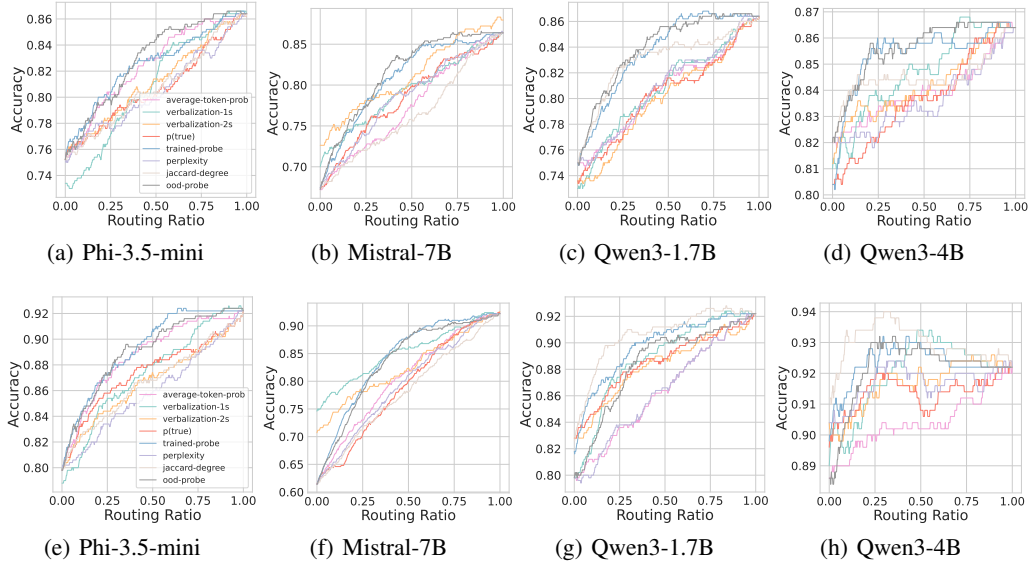


Figure 1: Overall accuracy vs. routing ratio with different UQ methods and SLMs. (a) and (b) show the results of routing to DeepSeek-R1 on the CommonsenseQA dataset; c and (d) show the results of routing to GPT4-mini on the Openbook QA dataset.

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 [31, 10]. 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 6). 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 6 and Figure 5, we can observe that the alignment may vary across datasets for the same SLM and UQ method. For instance, Perplexity [20] 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 [63, 47], and verbalization-2s [63], consistently withhold low alignment between uncertainty and prediction correctness. More experimental results can be found in Appendix C.1.

Observation ②: Verbalization-based UQ methods struggle to extract uncertainty in SLMs for query routing. We find that verbalization methods like verbalization-2s [63] 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 1. Recent advancements [68, 72] 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 [47], 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 [47], OOD Probe [36, 47], and Perplexity [19] 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.

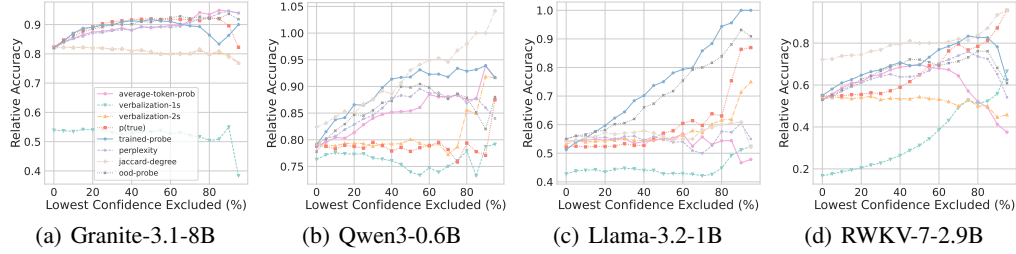


Figure 2: Relative accuracy of SLMs vs. LLMs on top- k % confident queries. “Relative accuracy” is the ratio of SLM accuracy to LLM accuracy. The x-axis “Lowest Conf. Excluded” shows the percentage of low-confidence queries removed; for example, 80 means 80% of queries with the lowest confidence are 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.

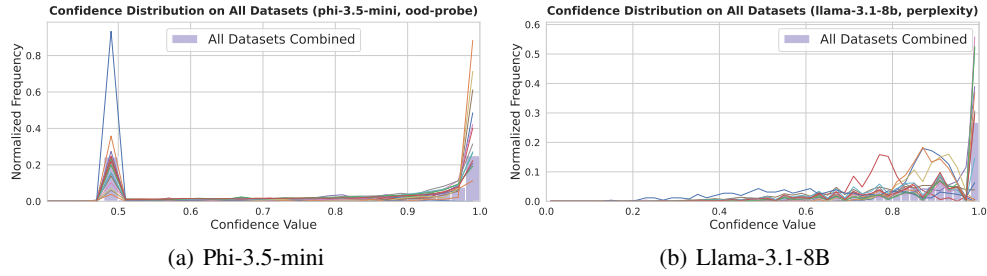


Figure 3: Uncertainty score distributions across 15 datasets. The histogram depicts the aggregated distribution from all datasets, while each curve represents a single dataset. (a) Confidence of Phi-3.5-mini by OOD Probe; (b) Confidence of Llama-3.1-8B by Perplexity.

Observation 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 2). 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 [47]) yield higher relative accuracy than weaker ones (e.g., verbalization-2s) across all query exclusion rates. Additional results appear in Appendix C.2.

Insights 5: The extracted confidence distribution is predominantly determined by the chosen SLM and UQ method, with minimal dependence on the downstream dataset. As shown in Figure 3, the confidence scores aggregated from 15 different tasks maintain a consistent overall shape across diverse datasets, indicating that dataset-specific factors have limited influence. In contrast, the distribution varies substantially across different SLMs and UQ methods, highlighting their dominant role in shaping the resulting confidence profiles.

3 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, we also observe that extracted uncertainty is primarily influenced by the selected SLM and uncertainty quantification (UQ) method, showing minimal dependence on the downstream dataset. Future work will focus on leveraging these insights to develop high-quality routing datasets that enable efficient SLM deployment on edge devices without re-calculating the uncertainty scores for routing.

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 Altenschmidt, 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] 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.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- [14] Gianluca Detommaso, Martin Bertran, Riccardo Fogliato, and Aaron Roth. Multicalibration for confidence scoring in llms. *arXiv preprint arXiv:2404.04689*, 2024.
- [15] 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.
- [16] 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.

- [17] 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.
- [18] 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.
- [19] 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.
- [20] 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.
- [21] Yarin Gal et al. Uncertainty in deep learning. 2016.
- [22] 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.
- [23] Xinyang Geng and Hao Liu. Openllama: An open reproduction of llama, May 2023.
- [24] IBM Granite Team. Granite 3.0 language models, 2024.
- [25] RLHF Griffin and Gemma Teams. Recurrentgemma: Moving past transformers for efficient open language models, 2024.
- [26] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafford, 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.
- [27] 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.
- [28] 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.
- [29] 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.
- [30] 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.
- [31] 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.

- [32] 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.
- [33] Mojan Javaheripi and Sébastien Bubeck. Phi-2: The surprising power of small language models, December 2023.
- [34] 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.
- [35] 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.
- [36] 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.
- [37] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*, 2023.
- [38] 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.
- [39] 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.
- [40] 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.
- [41] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- [42] Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words. *arXiv preprint arXiv:2205.14334*, 2022.
- [43] Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantification for black-box large language models. *arXiv preprint arXiv:2305.19187*, 2023.
- [44] Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. *ACL*, 2017.
- [45] 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.
- [46] 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.
- [47] 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.
- [48] 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.

- [49] AI Meta. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. *Meta AI Blog*. Retrieved December, 20:2024, 2024.
- [50] 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.
- [51] 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.
- [52] 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.
- [53] 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.
- [54] 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.
- [55] 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.
- [56] 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.
- [57] Subhro Roy and Dan Roth. Solving general arithmetic word problems. *arXiv preprint arXiv:1608.01413*, 2016.
- [58] 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.
- [59] 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.
- [60] 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.
- [61] 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.
- [62] 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.
- [63] 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.

- [64] 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.
- [65] 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.
- [66] 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.
- [67] 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.
- [68] 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.
- [69] 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.
- [70] 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.
- [71] 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.
- [72] Gal Yona, Roei Aharoni, and Mor Geva. Can large language models faithfully express their intrinsic uncertainty in words? *arXiv preprint arXiv:2405.16908*, 2024.
- [73] 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.
- [74] Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*, 2024.
- [75] 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.
- [76] 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.
- [77] 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.
- [78] Zhengping Zhou, Lezhi Li, Xinxin Chen, and Andy Li. Mini-giants: "small" language models and open source win-win. *arXiv preprint arXiv:2307.08189*, 2023.

A Benchmark Details

A.1 Details about Datasets

The details of the 15 datasets are further listed in Tabel 1. 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.

Table 1: 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

A.2 Details about Uncertainty Quantification Methods

We evaluate 8 approaches from the four categories in Section B.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(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 [77].

B Related Work

B.1 Reviewing Different Schools of Uncertainty Quantification and LLM Routing

Uncertainty quantification methods estimate a model’s confidence in its predictions [29]. For traditional classification and regression, uncertainty estimation is well-established [21]. 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 [19]. 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 [47, 50]. To enable LMs to verbalize confidence, researchers have proposed fine-tuning them to express uncertainty [42] or teaching them to verbalize confidence through in-context learning [16]. Verbalized confidence can take the form of linguistic expressions of uncertainty or numerical scores [22]. Multiple studies find that LLMs tend to be overconfident when reporting confidence [69, 63]. To mitigate this overconfidence, prompting strategies such as multi-step elicitation, top- k , and Chain-of-Thought [66] have been explored [63]. Sampling multiple response-confidence pairs and designing more effective aggregation strategies can also help mitigate overconfidence [69]. Moreover, [63] 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 [22, 31, 35]. 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 [19, 20, 65]. Because not all tokens in a sequence equally reflect semantic content, SAR reweights token likelihoods to emphasize more meaningful tokens [17]. 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 [47]. 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 [36, 63].

Via gauging output consistency. This line of research (e.g., SelfCheckGPT [48]) assumes that high-confidence LLMs produce consistent outputs [47]. A typical approach samples m responses for a given input query, measures inter-response similarity, and calculates a confidence score from meaning diversity [19]. Common ways to measure pairwise similarity include Natural Language Inference (NLI) and Jaccard similarity [22]. Consistency is then assessed by analyzing the similarity matrix, for instance, by counting semantic sets, summing eigenvalues of the graph Laplacian or computing eccentricity [43]. Because different sentences can express the same meaning, semantic entropy [37] 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 scores [22]. Training data is often obtained by sampling multiple answers per question at a fixed

Table 2: 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 B.1 for taxonomy details and Subsection 2.1 for method descriptions.

Uncertainty Quantification (UQ) Methods	Taxonomy	Model Access	Require Training?
Average Token Prob [47]	Token/sequence probabilities	White-box	No
$p(\text{True})$ [36]	Token/sequence probabilities	White-box	No
Perplexity [20]	Token/sequence probabilities	White-box	No
Jaccard Degree [43]	Output consistency	Black-box	No
Verbalization-1s [69, 63]	Verbalized uncertainty	Black-box	No
Verbalization-2s [63]	Verbalized uncertainty	Black-box	No
Trained Probe [4, 36, 47]	Uncertainty probe	White-box	Yes
OOD Probe [36, 47]	Uncertainty probe	White-box	Yes

453 temperature and labeling each for correctness [36]. A probe (commonly a multi-layer perceptron)
454 then takes hidden states as inputs to predict correctness [4, 39]. Because in-domain training data
455 is not always available, Contrast-Consistent Search trains probes unsupervisedly by maximizing
456 representation distances between contradictory answers on Yes/No questions [7]. Furthermore,
457 whether probes trained on out-of-distribution data remain effective is still under debate [36, 47, 37].

458 **B.2 Small Language Models**

459 Small Language Models (SLMs) are designed for deployment on resource-constrained devices like
460 desktops, smartphones, and wearables. Specifically, we consider the Transformer-based SLMs in this
461 work due to their state-of-the-art performance, like Phi-3-mini [1], TinyLlama [74], MobileLLM [45],
462 and Qwen-1.5B [5], LiteLLaMa-460M, OPT-125M [75], BLOOMZ (560M, 1.1B, 1.7B, 3B) [38],
463 SmolLM (135M, 360M, 1.7B) [3], OLMo (1B) [26], OLMoE (1B) [52], MobiLlama (0.5B, 1B) [62],
464 MobileLLaMA (1.4B, 2.7B) [9], OpenLLaMA (3B) [23]. These models are designed with lightweight
465 architectures to operate effectively within the computational and storage limitations of mobile devices
466 and edge hardware.

467 Recurrent Neural Networks (RNNs), like RWKV (1B, 3B, 7B) [55], Mamba (1.4B, 6.9B) [13], and
468 RecurrentGemma-2B [25], can provide promising solutions for on-device inference in resource-
469 constrained environments. These models leverage the recurrent nature of RNNs to process sequential
470 data efficiently without requiring a KV cache, which is suitable for resource-constrained on edge
471 devices. Specifically, RWKV introduces a hybrid RNN-Transformer backbone to capture long-term
472 dependencies while maintaining computational efficiency. Similarly, Mamba and RecurrentGemma
473 design recurrent layers for low-power consumption and high throughput inference, which can
474 significantly reduce memory and computational requirements, fostering low-latency applications
475 directly on devices.

C Additional Benchmarking Results

In this section, we present additional experimental results on evaluating the impact of uncertainty-correctness alignment on small language model (SLM) routing. **Since our studies yield over 5,600 results, we here present a representative subset in the following section. The full set of results will be provided in the Github repository.** For the experiment (Section C.1 and Section C.2), we provide the complete set of results, including the AUC measurements for uncertainty-correctness alignment and the performance of uncertainty-based routing. Each dataset referenced in the experiments is treated as a novel dataset for evaluation.

C.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.

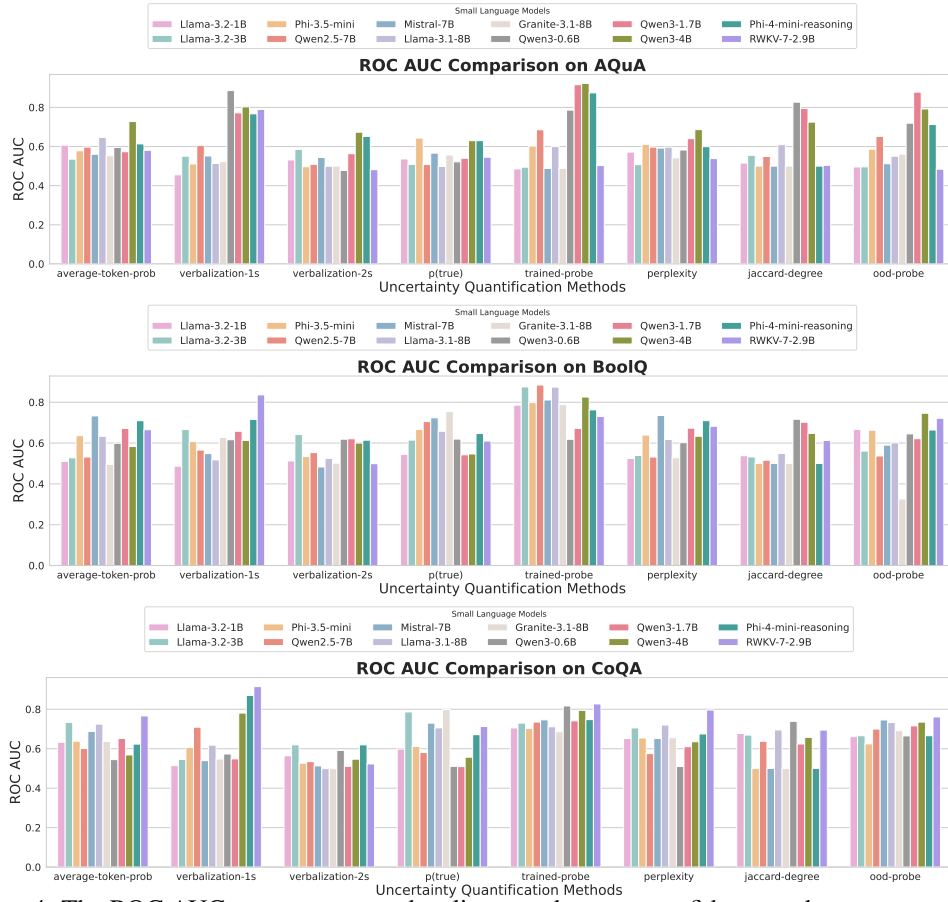


Figure 4: The ROC AUC scores measure the alignment between confidence and correctness across different SLMs and uncertainty quantification methods on AQUA, BoolQ, and CoQA. A higher ROC AUC indicates a stronger alignment.

488 **Results of Alignment between uncertainty and correctness.**

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

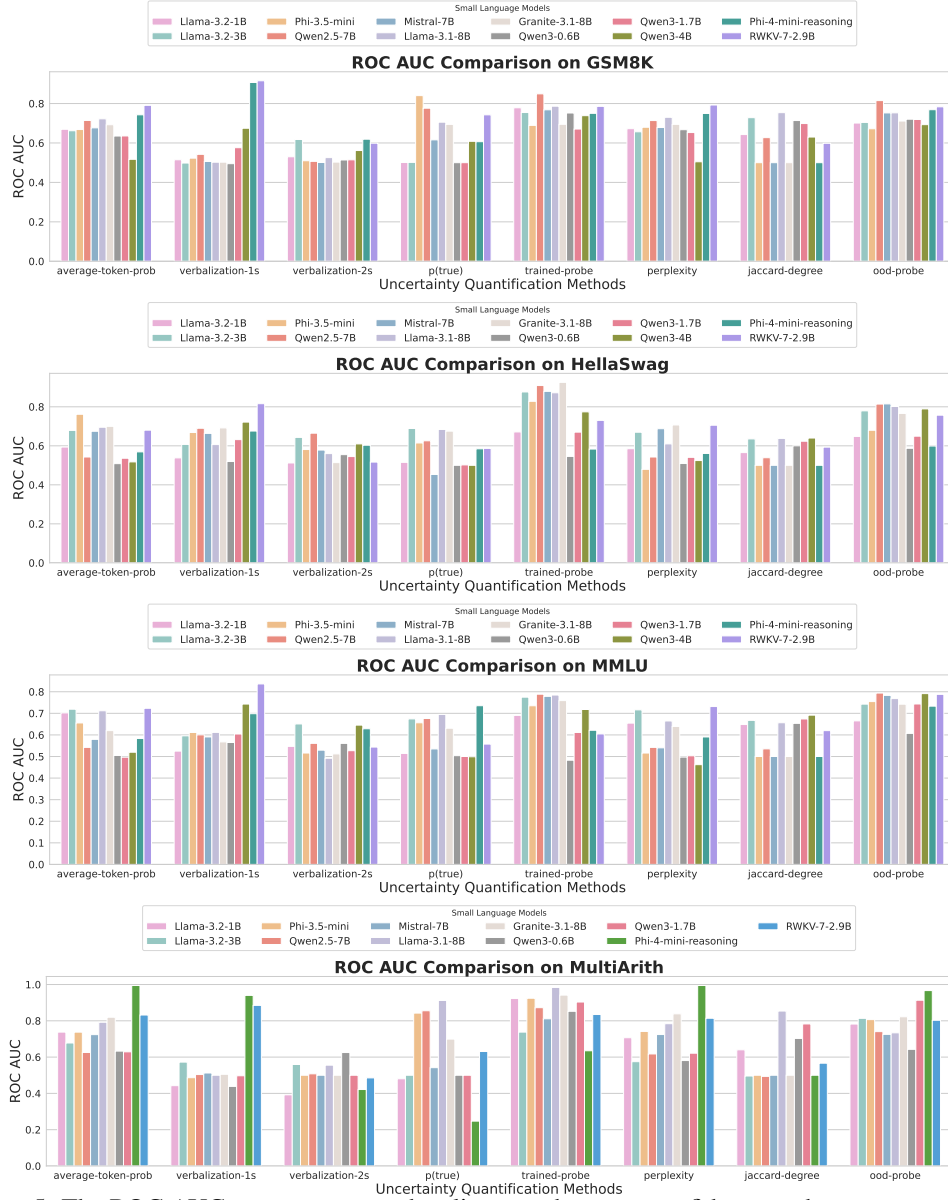


Figure 5: The ROC AUC scores measure the alignment between confidence and correctness across different SLMs and uncertainty quantification methods on GSM8K, HellaSwag, MMLU, and Multi-Arith. A higher ROC AUC indicates a stronger alignment.

491 **Results of Alignment between uncertainty and correctness.**

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

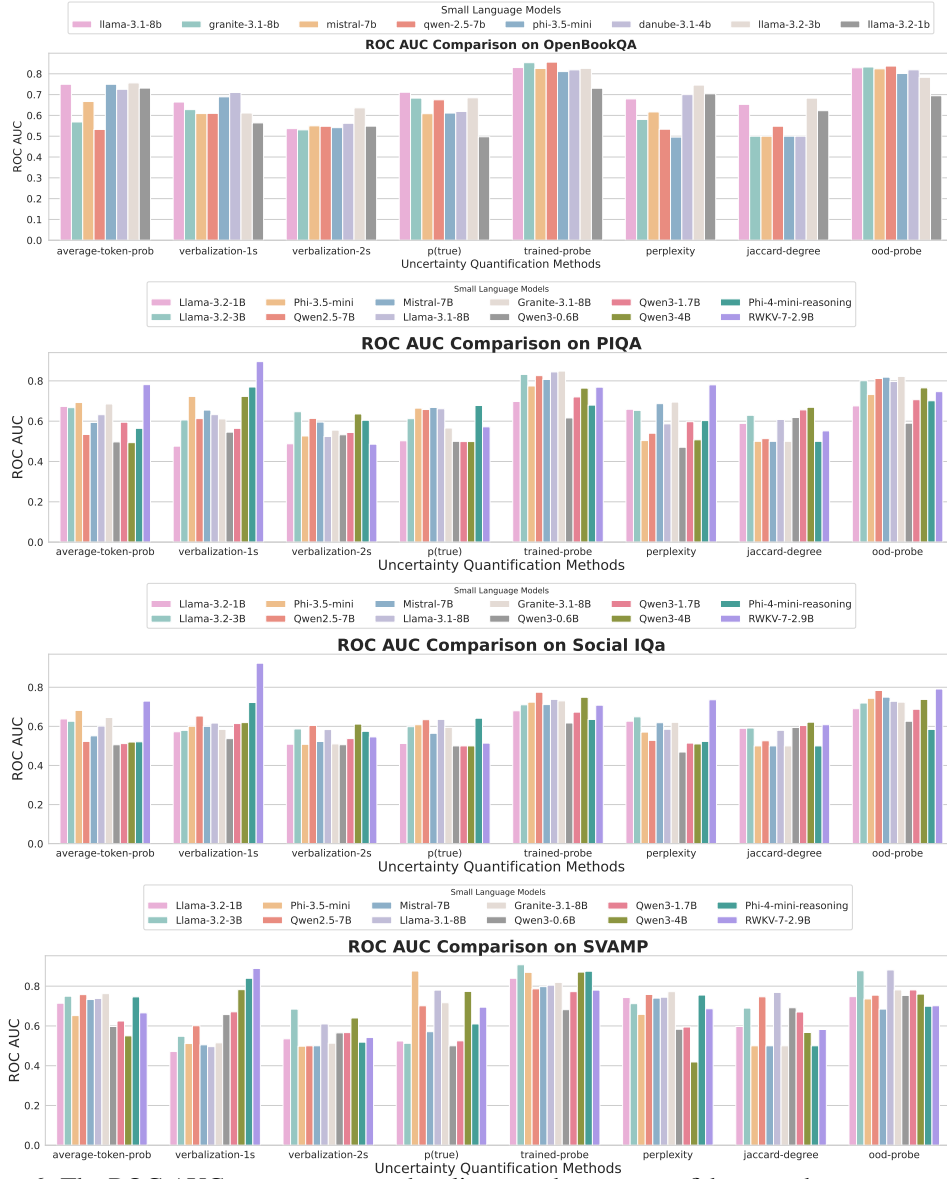


Figure 6: The ROC AUC scores measure the alignment between confidence and correctness across different SLMs and uncertainty quantification methods on OpenBookQA, PIQA, SocialIQA, and SVAMP datasets. A higher ROC AUC indicates a stronger alignment.

494 **Results of Alignment between uncertainty and correctness.**

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

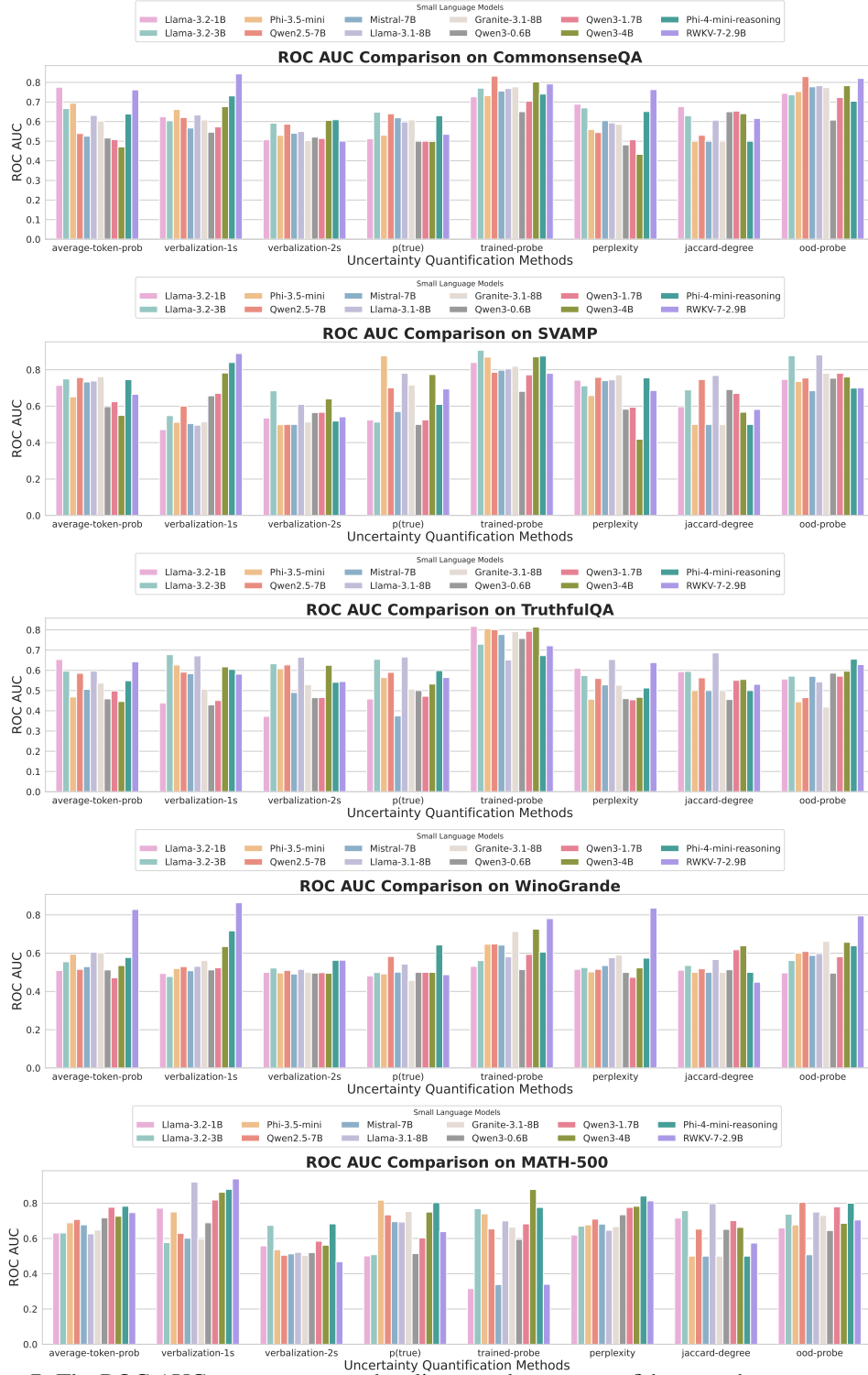


Figure 7: The ROC AUC scores measure the alignment between confidence and correctness across different SLMs and uncertainty quantification methods on CommonsenseQA. A higher ROC AUC indicates a stronger alignment.

497 C.2 Evaluation on Uncertainty-based Routing Approaches

498 Results of routing to GPT-4.1-Mini

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

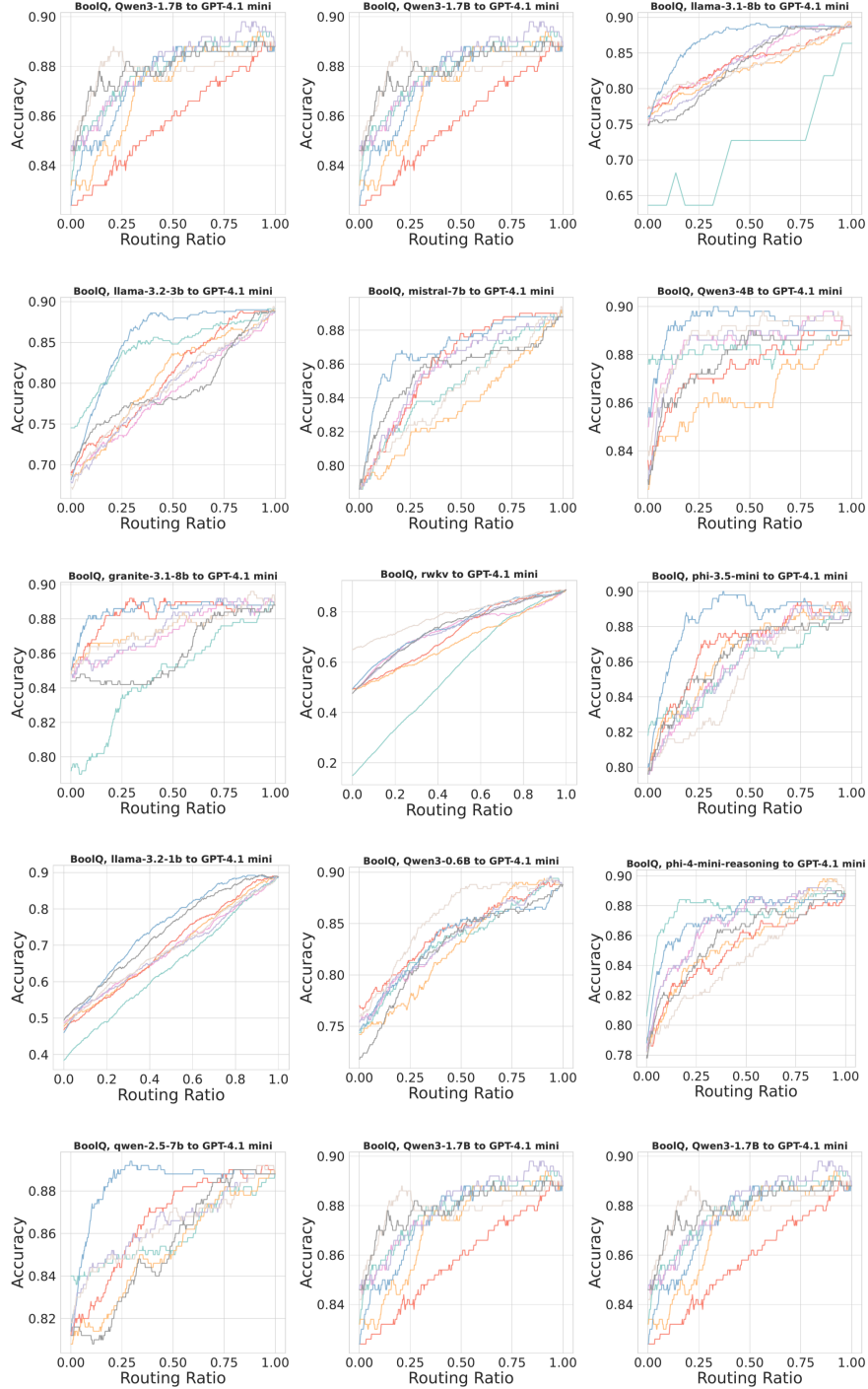


Figure 8: Overall accuracy vs. routing ratio with different UQ methods and SLMs.

501 **Results of routing to DeepSeek-R1**

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

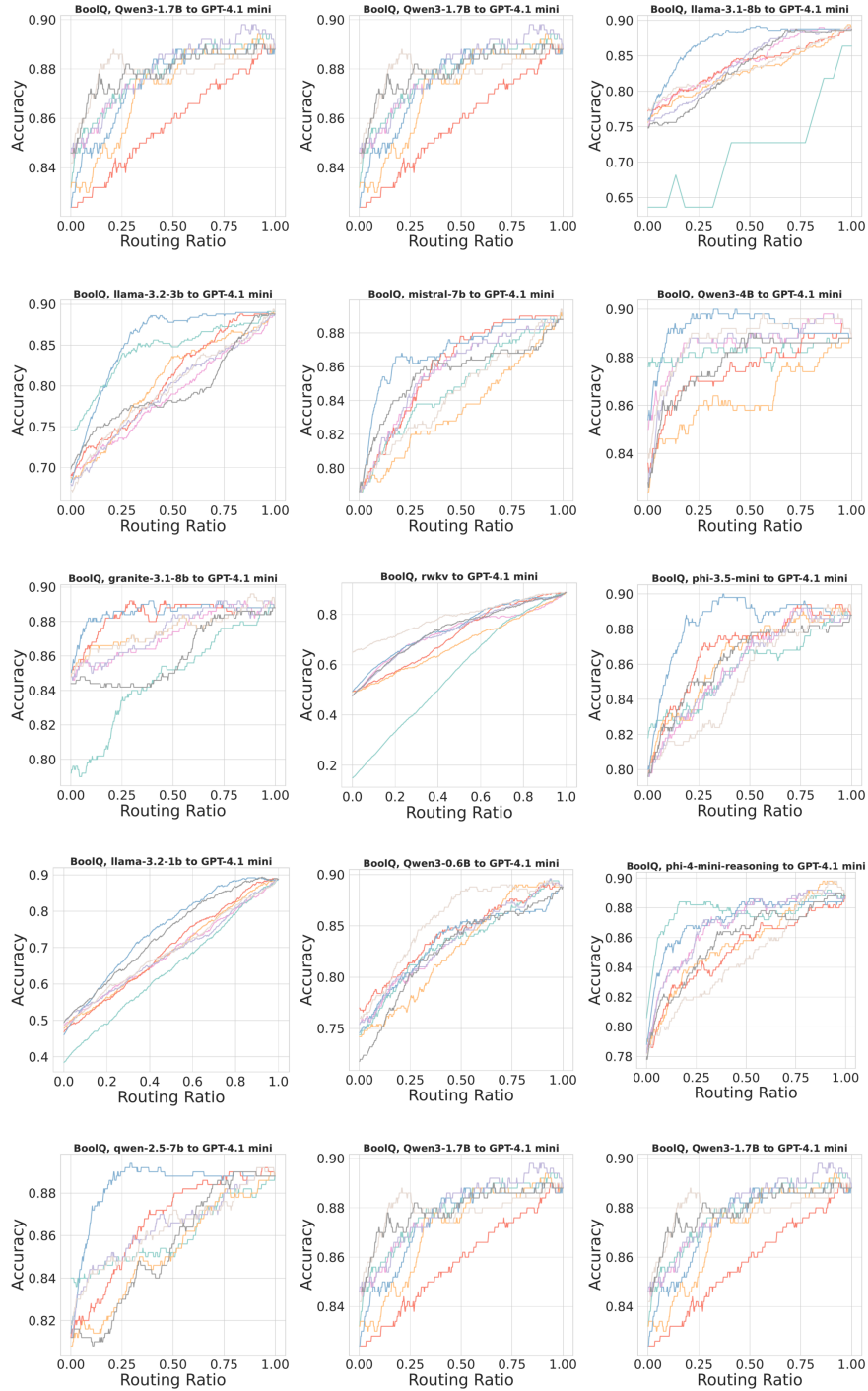


Figure 9: Overall accuracy vs. routing ratio with different UQ methods and SLMs.

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

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

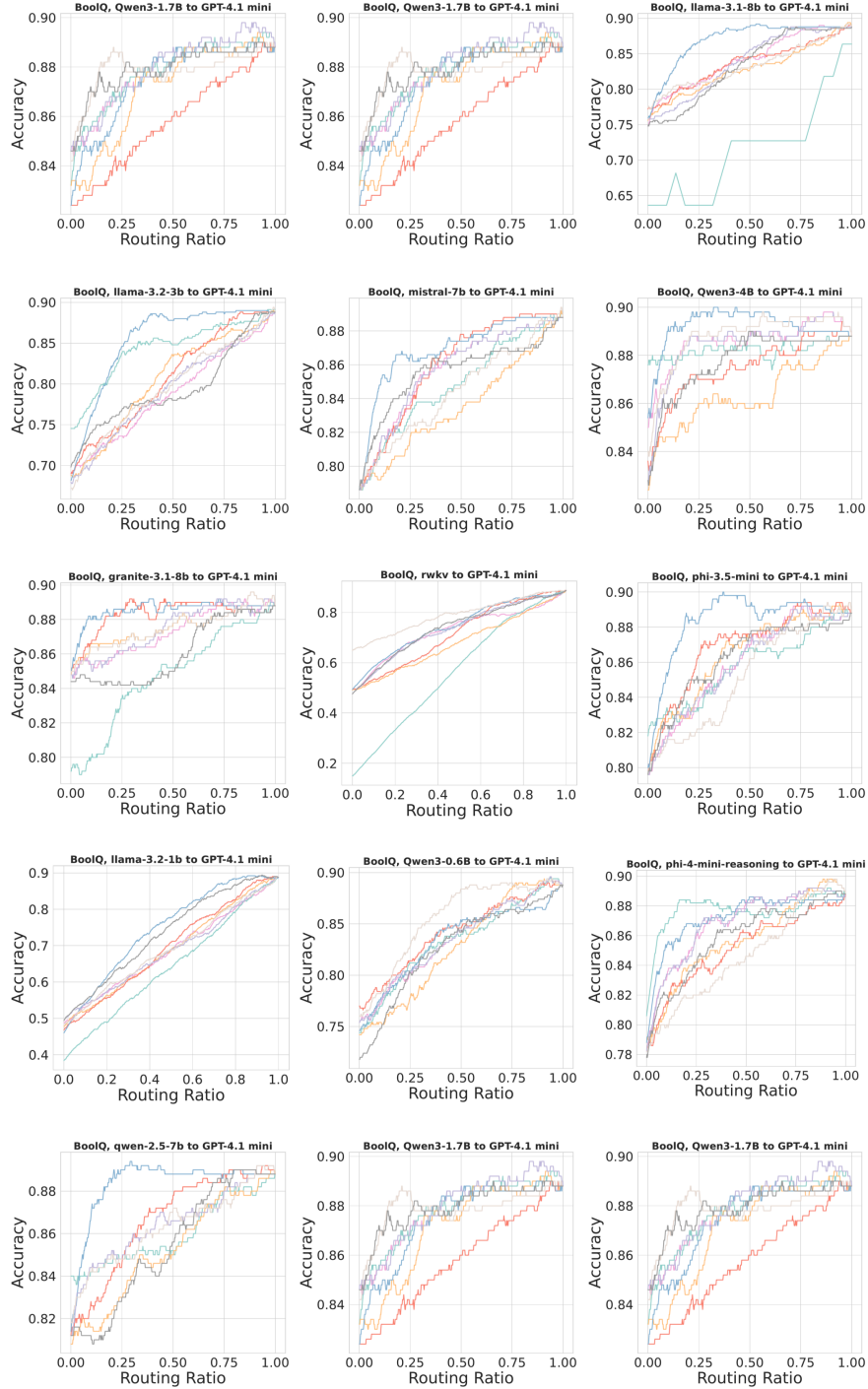


Figure 10: Overall accuracy vs. routing ratio with different UQ methods and SLMs.

507 **Results of routing to Qwen3-32B**

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

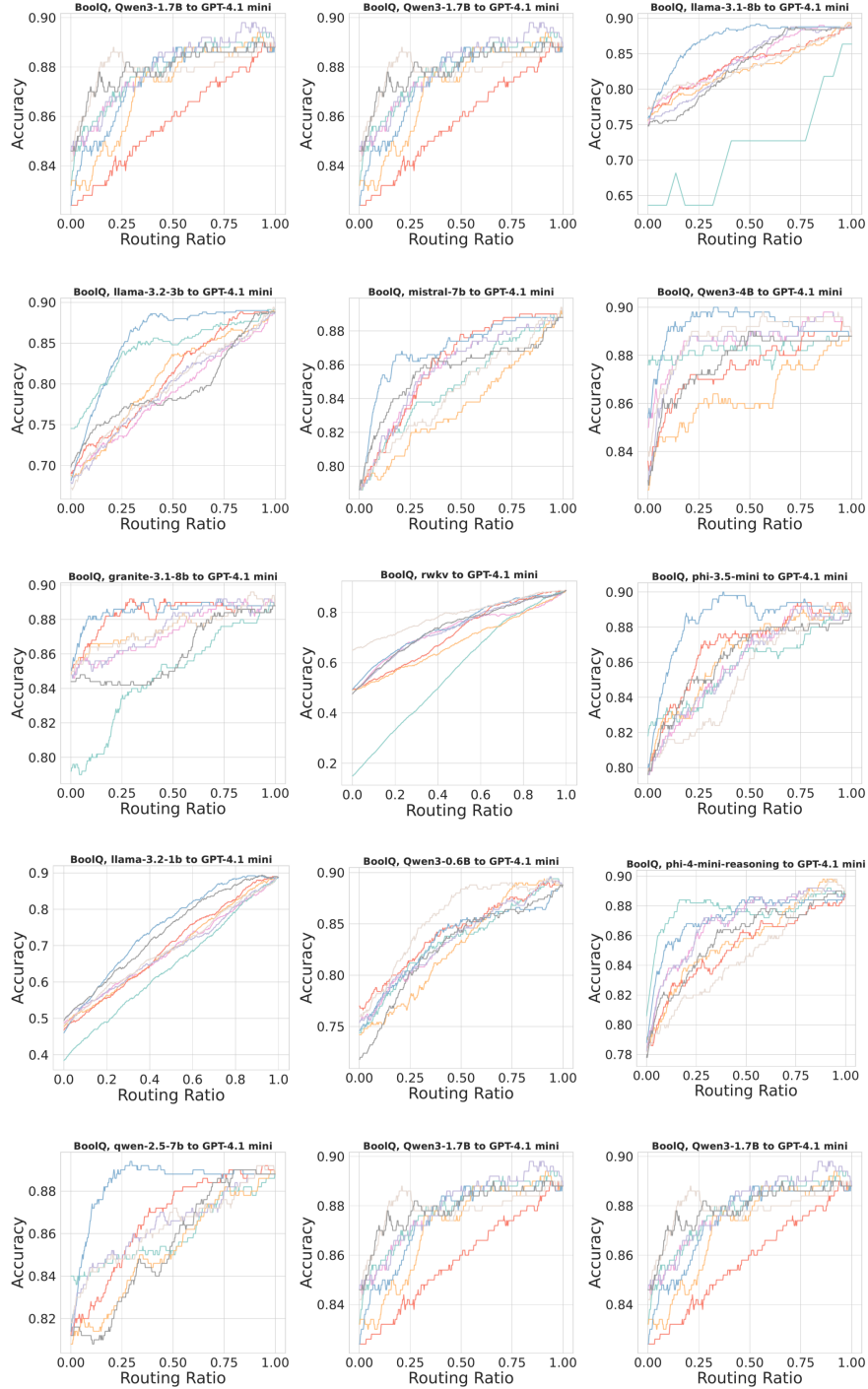


Figure 11: Overall accuracy vs. routing ratio with different UQ methods and SLMs.