

ENHANCING MULTILINGUAL REASONING IN LLMs: INSIGHTS FROM CROSS-LINGUISTIC CORRELATIONS AND OPTIMAL DATA PROPORTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) typically rely on fine-tuning to enhance their reasoning capabilities across various languages. However, limited research has been conducted on the optimal balance of language proportions within multilingual reasoning datasets. To fill this gap, we performed a systematic study to examine how different proportions of language data in multilingual reasoning datasets influence fine-tuning performance. Our study revealed a clear relationship between language proportions in datasets and the fine-tuning performance of LLMs. By fine-tuning multiple LLMs using the appropriate language distributions and data volumes identified in our study, we achieved state-of-the-art performance in both multilingual mathematical reasoning and solving mathematical problems using Python code. Furthermore, our approach significantly reduced data volume requirements and translation costs compared to existing methods, providing a valuable reference for future research.

1 INTRODUCTION

Despite recent significant advancements in LLMs (OpenAI, 2023; Chowdhery et al., 2023; Touvron et al., 2023; Brown et al., 2020; Workshop et al., 2023), LLMs still encounter considerable challenges in reasoning tasks for non-English languages, particularly low-resource languages (Shi et al., 2022b; Huang et al., 2023b; Qin et al., 2023b). A common scenario is that during the post-training process of LLMs, it is necessary to simultaneously enhance their reasoning capabilities across multiple languages (Zhu et al., 2024a). A critical research question thus emerges: how can LLMs, predominantly trained on high-resource languages like English, effectively generalize their reasoning capabilities to low-resource languages that suffer from insufficient training data? Addressing this question is essential for developing more inclusive and globally applicable AI systems.

Fine-tuning has been the commonly adopted solution to this issue. Prior studies have primarily focused on fine-tuning models using translated mathematical reasoning datasets across multiple languages Yoon et al. (2024). For instance, one approach enhanced LLMs' multilingual reasoning capabilities using translated datasets (Chen et al., 2023), while another improved model performance through question alignment (Zhu et al., 2024b;a). While these methods have improved multilingual reasoning performance, they still exhibit significant limitations. First, previous work often involved translating English datasets into multiple languages in equal proportions for fine-tuning. This approach is costly, overlooks the impact of language proportions in the fine-tuning datasets Zhu et al. (2024b). Second, in scenarios requiring extensive multilingual data, it is impractical to translate English data into multiple languages in equal amounts Hendrickson et al. (2013). Thus, the key is to efficiently leverage a small amount of low-resource language data to broadly enhance the multilingual reasoning capabilities of LLMs.

Moreover, earlier studies mainly focused on a limited number of languages and rarely tested generalization capabilities across more than 20 languages Zhang et al. (2024b); Zhu et al. (2024b); Yoon et al. (2024). Additionally, most prior research was conducted on only a few LLMs, and it remains unclear whether these approaches generalize well to other, more advanced models Zhang et al. (2024a); Zhu et al. (2024b); Yoon et al. (2024).

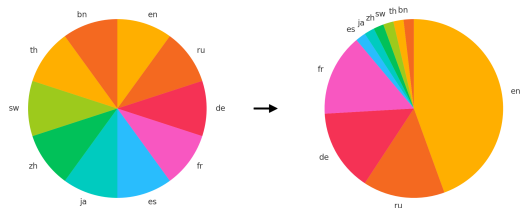
054
055
056
057
058
059
060
061
062

Figure 1: Two charts represent the proportion of data for different languages in the multilingual mathematical reasoning dataset. In previous methods, the amount of data for each language was the same (as shown in the left chart). Through our research, we determined the optimal proportion for each language (as shown in the right chart).

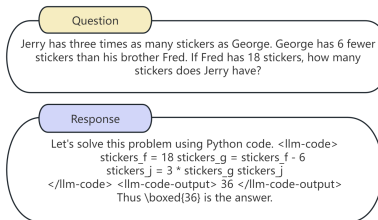
063
064
065
066
067
068
069
070
071

Figure 2: Following Toshniwal et al. (2024)’s study, we use this data format to teach model to solve mathematical reasoning task with executable code.

072
073
074
075
076
077
078
079
080
081
082
083

In this paper, we address these challenges through a systematic investigation of how varying language proportions in multilingual mathematical reasoning datasets affect fine-tuning outcomes, and determine the optimal amount of data needed for fine-tuning, as illustrated in Figure 1. Directly exploring the optimal ratio among tens of languages is extremely challenging, as it requires determining the optimal proportions among tens of variables, which results in an extremely large search space. So we adopted an innovative approach to effectively reduce the search space and determined the appropriate ratio for the 10 languages commonly used in previous works (Zhu et al., 2024b;a), which was further extended to 25 languages. Additionally, we explored the impact of data volume on the fine-tuning performance of LLMs. Ultimately, building on previous research, we created the largest multilingual mathematical reasoning dataset, HighMath-350k, alongside a multilingual code reasoning dataset (Figure 2), HighCode-350k. Fine-tuning multiple LLMs on these datasets resulted in state-of-the-art performance, clearly demonstrating the efficacy of our approach.

084
085
086
087
088
089
090
091
092
093

In summary, our key contributions are as follows:

- Through more than 600 groups of experiments, we analyzed and modeled how the proportions of different languages in multilingual reasoning datasets affect fine-tuning performance, and determined the appropriate volume of fine-tuning data.
- We constructed multilingual chain-of-thought mathematical reasoning dataset, HighMath-350k, as well as the multilingual code reasoning dataset, HighCode-350k.
- Based on multiple LLMs, we trained state-of-the-art multilingual chain-of-thought reasoning models and code reasoning models.

094
095
096
097

2 RELATED WORK

098
099
100
101
102
103
104
105
106
107

2.1 MULTILINGUAL MATHEMATICAL REASONING

Mathematical reasoning is a core task for assessing the intelligence of LLMs (Zhang et al., 2024b;a), requiring them to understand mathematical problems and generate answers through step-by-step reasoning (Ahn et al., 2024; Zhang et al., 2022; Liu et al., 2023). (Shi et al., 2022b) expanded this evaluation to a multilingual setting by translating English math questions from the GSM8K test set (Cobbe et al., 2021) into various non-English languages, introducing the multilingual benchmark MGSM. Efforts to improve LLMs’ multilingual reasoning performance have been continuing. For example, (Huang et al., 2023a) and (Qin et al., 2023a) explored prompting ChatGPT (OpenAI, 2023) to translate non-English questions into English and generate answers based on the translations. However, (Hu et al., 2024) found that this prompting strategy is not consistently effective for open-source LLMs. To enhance the multilingual capabilities of these models, researchers like (Nguyen et al., 2024) have explored continued pretraining on large-scale non-English corpora. However,

108 this approach is resource-intensive and data-inefficient, highlighting the need for more optimized
109 fine-tuning strategies—a gap our study aims to address.
110

111 2.2 LLM’S LANGUAGE PREFERENCE 112

113 LLMs, with their large-scale parameters (Aki, 1967; Huo & Kassab, 2009; Wang et al., 2006; Rosen-
114 feld, 1999), pre-trained on vast corpora and fine-tuned on comprehensive instruction datasets (Huang
115 et al., 2016; Zhao et al., 2020; Lindström & Abraham, 2022; Koncel-Kedziorski et al., 2016), have
116 demonstrated impressive intelligence (Touvron et al., 2023; Floridi & Chiriatti, 2020; Kalyan, 2023;
117 Lagler et al., 2013). However, empirical studies show that LLMs still struggle with multilingual
118 scenarios, particularly for low-resource languages (Shi et al., 2022a; Zhu et al., 2024b; Weysow
119 et al., 2024). This issue is mainly due to the dominance of English in both pretraining (Blevins &
120 Zettlemoyer, 2022) and instruction datasets (Wang et al., 2023). In this work, we focus on a core
121 capability of LLMs—their reasoning ability—and aim to push the boundaries of their performance
122 in multilingual reasoning tasks.
123

124 3 METHODOLOGY 125

126 This section describes the systematic method used to determine the appropriate language proportions
127 in the dataset and the required data volume for effective fine-tuning. Our method is divided into three
128 phases, as detailed below, with the experimental setup provided in Section 4.1.
129
130

131 3.1 PHASE ONE: LANGUAGE ALIGNMENT ANALYSIS 132

133 Directly determining the optimal ratio among 10 languages is also unrealistic, as directly exploring
134 the optimal proportions among 10 variables would result in a huge search space. To address this
135 issue and better generalize the reasoning capabilities of LLMs to different language contexts, we
136 first examined the correlation between LLMs’ alignment with English and other languages, and their
137 mathematical reasoning capabilities. Through this experiment, we aim to reduce the search space to
138 an acceptable size.
139

140 This alignment was assessed through a Non-English to English translation task. First, several
141 LLMs were fine-tuned on the English mathematical reasoning dataset, MetaMathQA-395k (Yu et al.,
142 2024), to evaluate their translation capabilities from Non-English to English, both before and after
143 fine-tuning. These models were then fine-tuned on translation data from non-English to English, and
144 their performance on multilingual mathematical reasoning tasks was evaluated. Our results showed
145 a positive correlation between LLMs’ reasoning abilities and their alignment between English and
146 French, German, and Russian, while no significant correlation was found with other languages. This
147 finding will greatly help the subsequent experiments.
148

149 3.2 PHASE TWO: LANGUAGE GROUP OPTIMIZATION 150

151 Building on the insights from Phase One, we constructed a multilingual mathematical reasoning
152 dataset encompassing 10 languages and grouped these 10 languages into three groups (as shown in
153 Table 3) based on the findings from Phase One, treating them as three variables for further explo-
154 ration. By systematically varying the proportions of data for each group, we utilized five models and
155 conducted over 600 experiments to identify the optimal distribution of three groups. The findings
156 indicate that, given a fixed total data volume, the optimal language ratio for enhancing the multilin-
157 gual reasoning capabilities of LLMs is group1:group2:group3 = 4:4:1 (en:ru:de:fr:es:ja:zh:bn:sw =
158 24:8:8:8:1:1:1:1:1). Additionally, we employed Gaussian Process Regression (GPR) to model the
159 relationship between data proportions and fine-tuning outcomes. This statistical approach allows us
160 to predict optimal language distributions based on experimental data, thereby guiding efficient data
161 allocation in subsequent phases. We also evaluated the contribution of each language to the overall
improvement of multilingual reasoning.

3.3 PHASE THREE: DATA VOLUME DETERMINATION

To determine the appropriate volume of fine-tuning data, we first conducted experiments with 10 languages and later extended the study to 25 languages. We tested varying data volumes and compared the results with previous approaches that used equal data volumes across languages. The experiments included fine-tuning data volumes up to 9.875M and scaling model parameters to 70B, demonstrating the clear advantages of our proposed methodology.

4 EXPERIMENT

4.1 OPTIMAL PROPORTION AND DATA VOLUME SELECTION FOR MULTILINGUAL FINE-TUNING

This section describes the experimental setup and key results for the three phases outlined in the methodology.

4.1.1 EXPERIMENTAL SETUP

Base LLMs

Different base LLMs were selected for each phase based on the research objectives. The specific models are listed in Table 1. We employed the FMT fine-tuning method by default, with training conducted for 3 epochs.

Table 1: Base LLMs used in different phases

Phase	Base LLMs
Phase 1	LLaMA3.1-8B, LLaMA3-8B, Phi3.5-mini-instruct, Mistral-7B-v0.3, Qwen2-7B
Phase 2	LLaMA3.1-8B, LLaMA3-8B, Phi3.5-mini-instruct, Qwen2-7B
Phase 3	LLaMA3-8B, LLaMA3-70B

Experimental Languages

The selection of test languages for each phase was guided by specific research objectives and computational constraints. In Phase One, the objective was to investigate the relationship between alignment with English and the mathematical reasoning capabilities of various non-English languages. To achieve this, We primarily selected several high-resource non-English languages. Phase Two employed the 10 languages from the MGSM dataset, while Phase Three aimed to assess the generalizability of our approach across 25 widely spoken languages. The languages chosen for each phase are detailed in Table 2.

Table 2: Experimental languages used in different phases

Phase	Experimental Languages
Phase 1	en, ru, de, fr, es, ja, zh, sw, th, bn, lt, cs, ka, ar
Phase 2	en, ru, de, fr, es, ja, zh, sw, th, bn
Phase 3	zh, en, es, hi, ar, pt, bn, ru, ja, pa, jv, de, ko, fr, te, vi, tr, ta, it, fa, ur, mr, sw, th, pl

We translated both the questions and answers in the multilingual chain-of-thought mathematical reasoning dataset. To ensure that the generated multilingual data effectively enhanced the model’s reasoning capabilities, we used the powerful open-source model DeepSeek-Chat-v2-236B (DeepSeek-AI, 2024) for high-quality translations, adjusting its hyperparameters for this task. To maintain translation quality and consistency, we implemented the following strategies:

- Mathematical formulas were preserved, and all numbers were converted to Arabic numerals to facilitate cross-linguistic prediction.
- To improve translation accuracy, we included two examples in the prompts for each language.

For the mathematical reasoning with code dataset, we only translated the question parts since the solution involves Python code.

To ensure consistency across languages, we extracted all mathematical expressions from the translated questions and answers. If the calculations were correct and matched the English version, the translation was deemed accurate. If errors persisted across five consecutive translations, the corresponding data was discarded.

This approach ensures coherence and accuracy in the translation process, enabling comprehensive evaluation and application of the dataset in a multilingual context, while maintaining both linguistic and mathematical integrity.

Tasks and Datasets

In terms of test task selection, in Phase One, we used a translation task to assess the alignment between English and non-English languages within the LLM. In the subsequent phases, chain-of-thought mathematical reasoning tasks were used to evaluate the LLM’s reasoning abilities.

For the training datasets, Phase One used the MetaMathQA-395k and WMT datasets. In later phases, we translated the 395k entries from MetaMathQA into 25 languages, generating a total of 9.875M entries. In Phase Two, 200k entries from each language were selected, totaling 1M entries. Phase Three used the full 9.875M dataset.

For the test datasets, Phase One used the MGSM dataset to evaluate the LLMs’ multilingual mathematical reasoning capabilities, and 3,000 entries from the WMT dataset were selected to assess translation performance. The WMT dataset is a collection of parallel texts designed to advance machine translation systems. In the subsequent phases, the MGSM dataset was used as the primary test dataset.

4.1.2 MAIN RESULTS

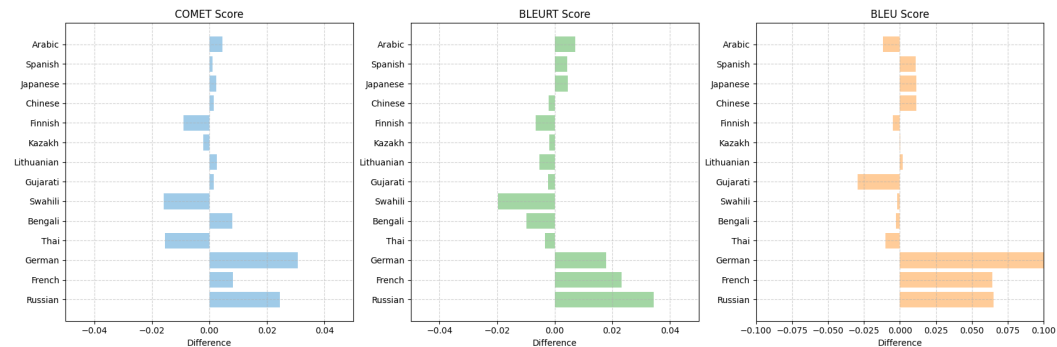


Figure 3: The figure shows the average performance across four LLMs; for detailed experimental data and descriptions of evaluation metrics, please refer to the appendix.

Phase 1

Figure 3 illustrates the changes in translation performance across different languages before and after fine-tuning the LLMs on the MetaMathQA-395k dataset. Significant improvements were observed for German, Russian, and French, while performance in other languages remained mostly unchanged.

Figure 4 presents the model’s performance on the MGSM dataset after fine-tuning with different non-English to English translation pairs. Since no additional measures were taken to enhance reasoning ability, mathematical reasoning performance declined post-fine-tuning. However, German, Russian, and French still outperformed other languages in reasoning tasks.

These results suggest a positive correlation between the LLM’s reasoning ability and its alignment with languages like French, German, and Russian.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

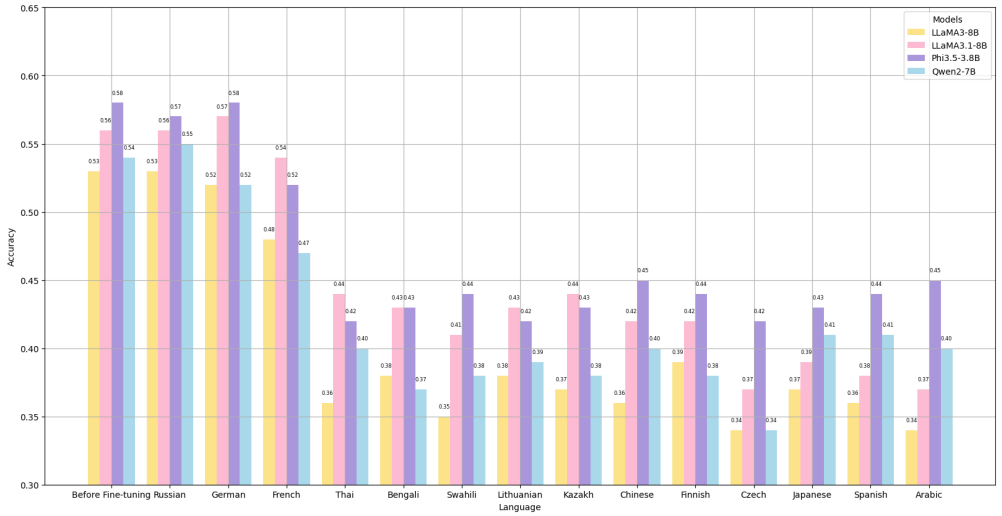


Figure 4: The score in the image represents the average of the LLM’s scores across the 10 languages.

Table 3: Language Groups for the Experiment

Group	Languages
Group 1	en
Group 2	de, ru, fr
Group 3	es, ja, zh, th, bn, sw

Phase 2

Based on the conclusions from the first phase, we divided the 10 languages selected for the second phase into three groups, as shown in Table 3. In each group, the amount of mathematical reasoning data for each language was kept consistent. By adjusting the overall data ratio among these three groups, we explored the optimal distribution of languages. As shown in Figure 5, although the pretraining corpora for each model differ, all five models demonstrated optimal performance when the ratio of the three language groups was around 24:24:6.(4:4:1)

To refine the optimal data ratio, we conducted further experiments to assess the relative importance of languages in the second and third groups. As shown in Figure 6, increasing the proportion of a single language did not lead to any clear improvement in prediction accuracy. Even when improvement occurred, the gains were minimal.

Mathematical Modeling

To better model the patterns observed in Figure 5 and provide suggestions for data allocation in future studies, we used the Gaussian Process Regression method (GPR). GPR is employed to model the underlying function $f(x)$ from the observed data, assuming that any finite set of function values follows a multivariate Gaussian distribution. The model can be expressed as:

$$y = f(X) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I),$$

where X represents the input features, y is the observed output, and ϵ is Gaussian noise with variance σ^2 .

We experimented with various mathematical models and computational methods, with detailed analyses provided in the appendix. Ultimately, we finalized the following kernel:

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

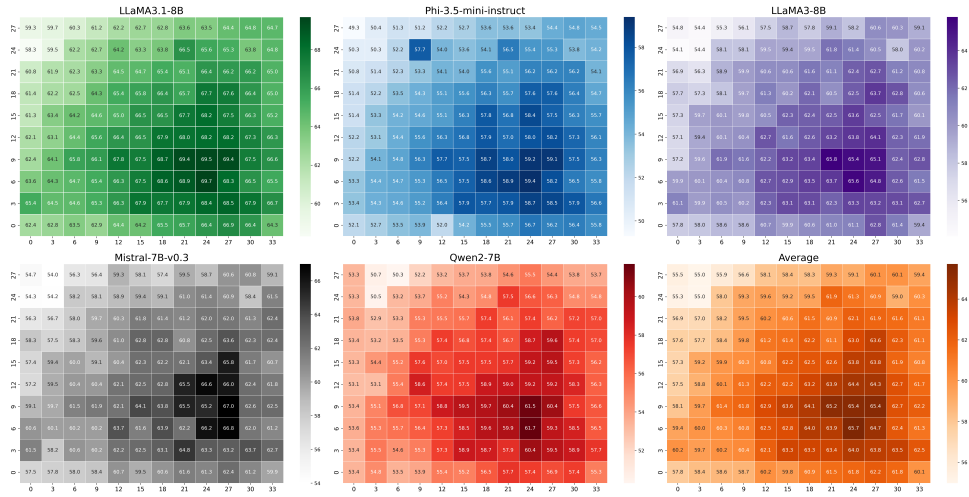


Figure 5: The relative proportion of English is fixed at 24. The horizontal axis represents the relative proportion of the second group of languages, while the vertical axis represents the relative proportion of the third group of languages. The results are averaged across the 10 languages.

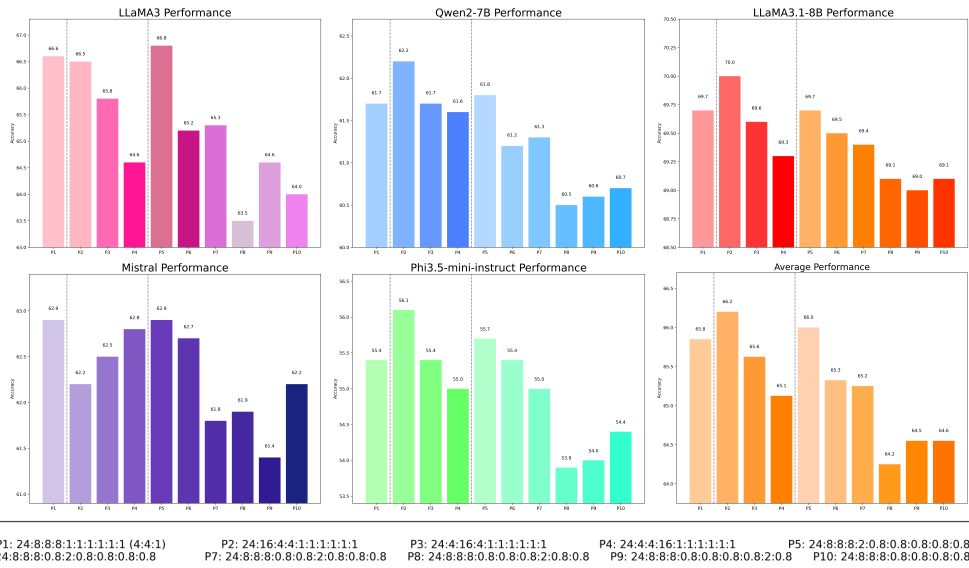


Figure 6: The proportions in the figure are represented as en:ru:de:fr:es:ja:zh:bn:sw:th in sequence.

$$k_{\text{RQ}}(X, X') = \sigma_f^2 \left(1 + \frac{\|X - X'\|^2}{2\alpha l^2} \right)^{-\alpha}$$

Phase 3

This phase of the experiment was divided into two groups. The first group used the original 10 languages to construct different sizes of the multilingual mathematical reasoning dataset, HighMath, with a language ratio of en:ru:de:fr:es:ja:zh:bn:sw:th = 24:8:8:8:1:1:1:1:1:1. The second group included 25 languages and constructed different sizes of the HighMath dataset, using a ratio of en:de:ru:fr:others = 24:8:8:8:12. It should be noted that our adjustment of proportions - increasing the ratios for German, French, and Russian - was qualitative in nature and may not represent the optimal distribution.

We conducted a series of experiments, with the results shown in Figure 7. HighMath-350k and HighMath-1.5M was identified as a relatively optimal size.

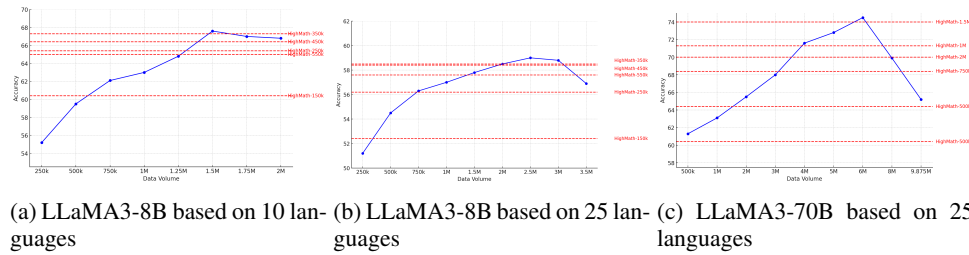


Figure 7: The x-axis represents the size of the dataset, while the line graphs show the results using equal amounts of data for each language. The horizontal lines indicate the fine-tuning results after optimizing the data proportions across different languages. The vertical axis shows the mean score across all languages. For detailed experimental data, please refer to the Appendix.

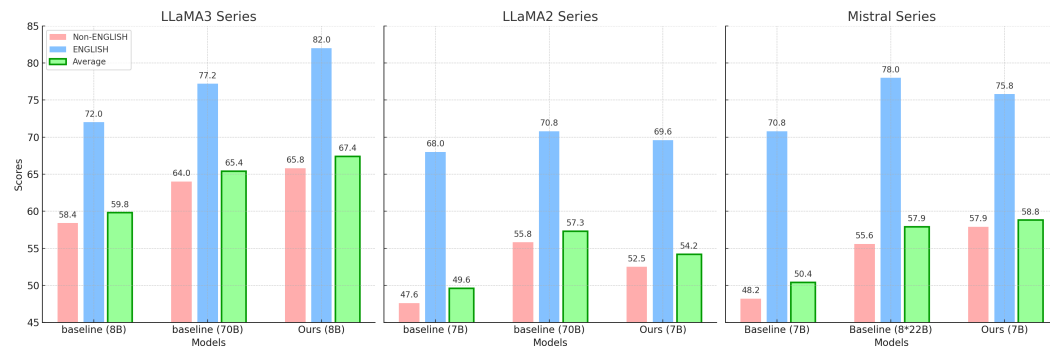


Figure 8: Evaluation results on the MGSM dataset. The score in the image represents the average of the LLM’s scores across the 10 languages.

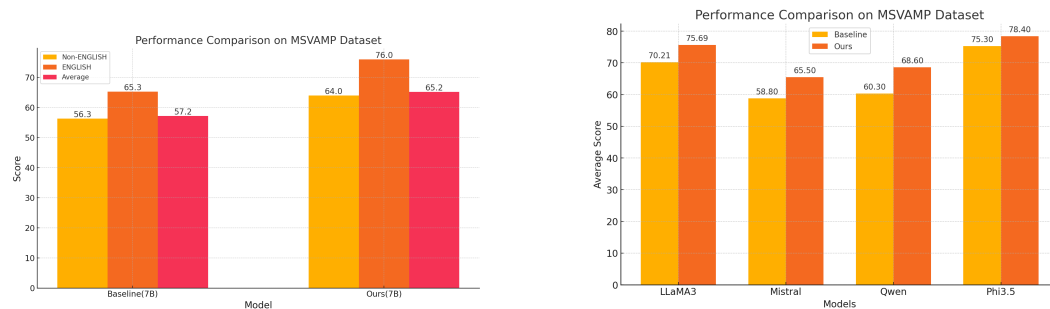


Figure 9: In the left image, the baseline represents (Zhu et al., 2024a), and non-English refers to the average performance across 9 non-English languages. In the right image, the baseline represents the fine-tuning results using an equal amount of data, but with the proportion of different languages in the dataset being balanced. The scores are averaged across 10 languages.

4.2 MODEL EVALUATION

In this section, we fine-tune several LLMs on the HighMath-350k dataset and compare their performance against baseline models. We also extend our approach to solving mathematical problems using python code, demonstrating the generalizability and effectiveness of our method.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

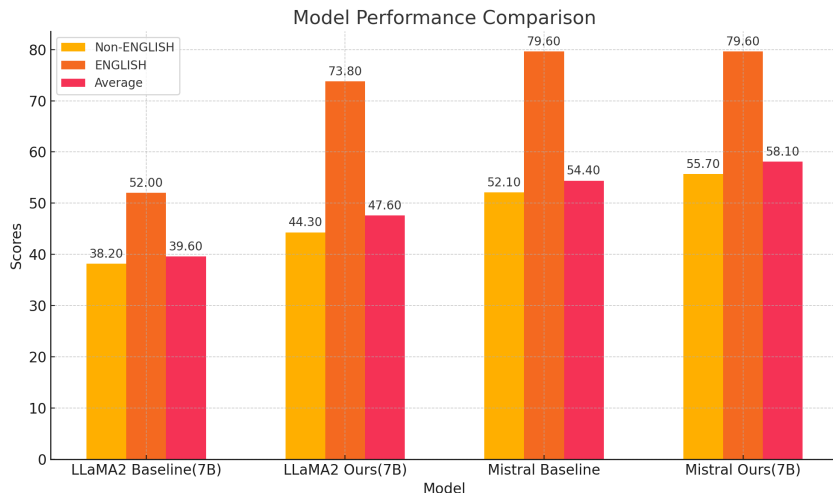


Figure 10: The score in the image represents the average of the LLM’s scores across the 10 languages. For our baseline, we adopted the state-of-the-art model proposed by Zhu et al. (2024a).

4.2.1 MATHEMATICAL REASONING

Baseline Models For the MGSM dataset, we selected state-of-the-art fine-tuned baselines from (Zhu et al., 2024b;a). Additionally, for the MSVAMP evaluation, we included baselines derived from fine-tuning LLaMA2 series models. Where direct comparisons were not feasible, we established an additional baseline using a dataset of 350k samples evenly distributed across 10 languages to ensure consistency.

Main Results The results, depicted in Figures 8 and 9, demonstrate that models fine-tuned on the HighMath-350k dataset achieved an average accuracy improvement of 12% on the MGSM evaluation dataset compared to the baselines. Furthermore, these models outperformed larger parameter models from the LLaMA3 and Mistral series by 8%, despite using only 70% of the training data volume required by baseline models. Furthermore, compared to methods that uniformly translate English data into multiple languages, our approach saves over 70% in translation costs.

4.2.2 SOLVING MATHEMATICAL REASONING TASKS USING PYTHON CODE

Using the HighCode-350k dataset, which focuses on solving mathematical problems through Python code, our models consistently outperformed the baselines across multiple base models. This demonstrates that optimized language proportions enhance both natural and code-based reasoning, highlighting the versatility of our fine-tuning approach.

5 CONCLUSION

For scenarios that require simultaneously enhancing LLMs’ reasoning capabilities across multiple languages through fine-tuning methods, we propose an innovative three-phase approach and conduct extensive experiments to investigate the impact of different language data proportions in multilingual training datasets on LLM fine-tuning for multilingual mathematical reasoning tasks. Additionally, we establish a mathematical model to illustrate the relationship between these factors. We also perform large-scale experiments to explore the correlation between fine-tuning effectiveness and the amount of training data in multilingual mathematical reasoning datasets. Using an optimal data scale, we construct the multilingual chain-of-thought reasoning dataset HighMath-350k and the HighCode-350k dataset, designed for solving mathematical problems using Python code in a multilingual context. Leveraging these datasets, we fine-tuned several LLMs to achieve outstanding performance, providing insights for future research.

REFERENCES

- 486
487
488 Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. Large language models
489 for mathematical reasoning: Progresses and challenges, 2024. URL <https://arxiv.org/abs/2402.00157>.
490
- 491 Keiiti Aki. Scaling law of seismic spectrum. *Journal of geophysical research*, 72(4):1217–1231,
492 1967.
493
- 494 Terra Blevins and Luke Zettlemoyer. Language contamination helps explain the cross-lingual capa-
495 bilities of english pretrained models, 2022. URL <https://arxiv.org/abs/2204.08110>.
- 496 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
497 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,
498 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
499 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz
500 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
501 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL
502 <https://arxiv.org/abs/2005.14165>.
- 503 Nuo Chen, Zinan Zheng, Ning Wu, Ming Gong, Yangqiu Song, Dongmei Zhang, and Jia Li. Break-
504 ing language barriers in multilingual mathematical reasoning: Insights and observations, 2023.
505 URL <https://arxiv.org/abs/2310.20246>.
506
- 507 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
508 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
509 Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):
510 1–113, 2023.
- 511 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
512 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
513 Schulman. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
514
- 515 DeepSeek-AI. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language
516 model, 2024.
517
- 518 Luciano Floridi and Massimo Chiriatti. Gpt-3: Its nature, scope, limits, and consequences. *Minds
519 and Machines*, 30:681–694, 2020.
520
- 521 Sherry G Hendrickson, Tracie C Harrison, Nora A Lopez, Aurea G Zegarra-Coronado, and Tiffany
522 Ricks. Translation cost, quality, and adequacy. *Journal of Nursing Scholarship*, 45(2):185–191,
523 2013.
- 524 Peng Hu, Sizhe Liu, Changjiang Gao, Xin Huang, Xue Han, Junlan Feng, Chao Deng, and Shujian
525 Huang. Large language models are cross-lingual knowledge-free reasoners, 2024. URL <https://arxiv.org/abs/2406.16655>.
526
- 527 Danqing Huang, Shuming Shi, Chin-Yew Lin, Jian Yin, and Wei-Ying Ma. How well do computers
528 solve math word problems? large-scale dataset construction and evaluation. In *Proceedings of the
529 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
530 pp. 887–896, 2016.
531
- 532 Haoyang Huang, Tianyi Tang, Dongdong Zhang, Wayne Xin Zhao, Ting Song, Yan Xia, and Furu
533 Wei. Not all languages are created equal in llms: Improving multilingual capability by cross-
534 lingual-thought prompting, 2023a. URL <https://arxiv.org/abs/2305.07004>.
- 535 Haoyang Huang, Tianyi Tang, Dongdong Zhang, Wayne Xin Zhao, Ting Song, Yan Xia, and Furu
536 Wei. Not all languages are created equal in llms: Improving multilingual capability by cross-
537 lingual-thought prompting. *arXiv preprint arXiv:2305.07004*, 2023b.
538
- 539 Yunlong Huo and Ghassan S Kassab. A scaling law of vascular volume. *Biophysical journal*, 96(2):
347–353, 2009.

- 540 Katikapalli Subramanyam Kalyan. A survey of gpt-3 family large language models including chat-
541 gpt and gpt-4. *Natural Language Processing Journal*, pp. 100048, 2023.
- 542
- 543 Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi.
544 Mawps: A math word problem repository. In *Proceedings of the 2016 conference of the north
545 american chapter of the association for computational linguistics: human language technologies*,
546 pp. 1152–1157, 2016.
- 547
- 548 Klemens Lagler, Michael Schindelegger, Johannes Böhm, Hana Krásná, and Tobias Nilsson. Gpt2:
549 Empirical slant delay model for radio space geodetic techniques. *Geophysical research letters*,
550 40(6):1069–1073, 2013.
- 551
- 552 Adam Dahlgren Lindström and Savitha Sam Abraham. Clevr-math: A dataset for compositional
553 language, visual and mathematical reasoning. *arXiv preprint arXiv:2208.05358*, 2022.
- 554
- 555 Xiangyang Liu, Tianqi Pang, and Chenyou Fan. Federated prompting and chain-of-thought reason-
556 ing for improving llms answering. In *International Conference on Knowledge Science, Engineer-
557 ing and Management*, pp. 3–11. Springer, 2023.
- 558
- 559 Xuan-Phi Nguyen, Wenxuan Zhang, Xin Li, Mahani Aljunied, Zhiqiang Hu, Chenhui Shen,
560 Yew Ken Chia, Xingxuan Li, Jianyu Wang, Qingyu Tan, Liying Cheng, Guanzheng Chen, Yue
561 Deng, Sen Yang, Chaoqun Liu, Hang Zhang, and Lidong Bing. Seallms – large language models
562 for southeast asia, 2024. URL <https://arxiv.org/abs/2312.00738>.
- 563
- 564 OpenAI. Gpt-4 technical report, 2023.
- 565
- 566 Libo Qin, Qiguang Chen, Fuxuan Wei, Shijue Huang, and Wanxiang Che. Cross-lingual prompt-
567 ing: Improving zero-shot chain-of-thought reasoning across languages. In Houda Bouamor,
568 Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Meth-
569 ods in Natural Language Processing*, pp. 2695–2709, Singapore, December 2023a. Associa-
570 tion for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.163. URL <https://aclanthology.org/2023.emnlp-main.163>.
- 571
- 572 Libo Qin, Qiguang Chen, Fuxuan Wei, Shijue Huang, and Wanxiang Che. Cross-lingual
573 prompting: Improving zero-shot chain-of-thought reasoning across languages. *arXiv preprint
574 arXiv:2310.14799*, 2023b.
- 575
- 576 Yaakov Rosenfeld. A quasi-universal scaling law for atomic transport in simple fluids. *Journal of
577 Physics: Condensed Matter*, 11(28):5415, 1999.
- 578
- 579 Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi,
580 Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Lan-
581 guage models are multilingual chain-of-thought reasoners, 2022a. URL <https://arxiv.org/abs/2210.03057>.
- 582
- 583 Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi,
584 Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, et al. Language models are multi-
585 lingual chain-of-thought reasoners. *arXiv preprint arXiv:2210.03057*, 2022b.
- 586
- 587 Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Git-
588 man. Openmathinstruct-1: A 1.8 million math instruction tuning dataset, 2024. URL <https://arxiv.org/abs/2402.10176>.
- 589
- 590 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
591 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open
592 and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 593
- 594 J Wang, HL Duan, ZP Huang, and Bhushan Lal Karihaloo. A scaling law for properties of nano-
595 structured materials. *Proceedings of the Royal Society A: Mathematical, Physical and Engineer-
596 ing Sciences*, 462(2069):1355–1363, 2006.

- 594 Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi
595 Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi.
596 How far can camels go? exploring the state of instruction tuning on open resources, 2023. URL
597 <https://arxiv.org/abs/2306.04751>.
- 598 Martin Weysow, Aton Kamanda, and Houari Sahraoui. Codeultrafeedback: An llm-as-a-
599 judge dataset for aligning large language models to coding preferences. *arXiv preprint*
600 *arXiv:2403.09032*, 2024.
- 602 BigScience Workshop, :, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić,
603 Daniel Hesslow, Roman Castagné, and Alexandra Sasha Luccioni et al. Bloom: A 176b-parameter
604 open-access multilingual language model, 2023. URL <https://arxiv.org/abs/2211.05100>.
- 606 Dongkeun Yoon, Joel Jang, Sungdong Kim, Seungone Kim, Sheikh Shafayat, and Minjoon Seo.
607 Langbridge: Multilingual reasoning without multilingual supervision, 2024. URL [https://](https://arxiv.org/abs/2401.10695)
608 arxiv.org/abs/2401.10695.
- 610 Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhen-
611 guo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions
612 for large language models. In *The Twelfth International Conference on Learning Representations*,
613 2024. URL <https://openreview.net/forum?id=N8N0hgNDRT>.
- 614 Boning Zhang, Chengxi Li, and Kai Fan. Mario eval: Evaluate your math llm with your math llm—a
615 mathematical dataset evaluation toolkit. *arXiv preprint arXiv:2404.13925*, 2024a.
- 616 Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou,
617 Pan Lu, Kai-Wei Chang, Peng Gao, et al. Mathverse: Does your multi-modal llm truly see the
618 diagrams in visual math problems? *arXiv preprint arXiv:2403.14624*, 2024b.
- 620 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in
621 large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- 622 Wei Zhao, Mingyue Shang, Yang Liu, Liang Wang, and Jingming Liu. Ape210k: A large-scale and
623 template-rich dataset of math word problems. *arXiv preprint arXiv:2009.11506*, 2020.
- 625 Wenhao Zhu, Shujian Huang, Fei Yuan, Cheng Chen, Jiajun Chen, and Alexandra Birch. The
626 power of question translation training in multilingual reasoning: Broadened scope and deepened
627 insights. *arXiv preprint arXiv:2405.01345*, 2024a.
- 628 Wenhao Zhu, Shujian Huang, Fei Yuan, Shuaijie She, Jiajun Chen, and Alexandra Birch. Question
629 translation training for better multilingual reasoning. *arXiv preprint arXiv:2401.07817*, 2024b.
- 630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

APPENDIX: GAUSSIAN PROCESS REGRESSION MODEL FITTING AND EVALUATION

INTRODUCTION TO GAUSSIAN PROCESS REGRESSION

A GAUSSIAN PROCESS REGRESSION: MATHEMATICAL FOUNDATION AND KERNEL SELECTION

A.1 MATHEMATICAL FOUNDATION OF GAUSSIAN PROCESS REGRESSION

Gaussian Process Regression (GPR) is a non-parametric Bayesian method used for regression problems. It assumes that the observed data can be generated by a Gaussian process, which is a collection of random variables such that any finite subset has a joint Gaussian distribution.

Given a training dataset $\{(X_i, y_i)\}_{i=1}^n$, where X_i represents the input variables (in this case, two-dimensional coordinates) and y_i represents the corresponding observations, we assume:

$$y_i = f(X_i) + \varepsilon_i$$

where $f(X)$ is the latent true function, and $\varepsilon_i \sim \mathcal{N}(0, \sigma_n^2)$ is independent and identically distributed Gaussian noise.

A Gaussian process prior is placed on $f(X)$:

$$f(X) \sim \mathcal{GP}(m(X), k(X, X'))$$

where $m(X)$ is the mean function, typically set to zero, $m(X) = 0$, and $k(X, X')$ is the covariance function (kernel) that measures the similarity between input points.

A.1.1 PREDICTIVE DISTRIBUTION

For a new input point X_* , the predictive distribution is also a Gaussian distribution with mean and variance given by:

$$\mu_* = k_*^\top (K + \sigma_n^2 I)^{-1} y, \quad \sigma_*^2 = k(X_*, X_*) - k_*^\top (K + \sigma_n^2 I)^{-1} k_*$$

where:

- $y = [y_1, y_2, \dots, y_n]^\top$ is the vector of observed values.
- K is the kernel matrix, $K_{ij} = k(X_i, X_j)$.
- $k_* = [k(X_*, X_1), k(X_*, X_2), \dots, k(X_*, X_n)]^\top$ is the covariance vector between the new input point and the training data.
- I is the identity matrix.

A.2 KERNEL FUNCTION SELECTION AND PARAMETER TUNING

Choosing an appropriate kernel function $k(X, X')$ is crucial for the effectiveness of Gaussian Process Regression. Several kernel functions were considered:

Radial Basis Function (RBF) Kernel / Gaussian Kernel:

$$k_{\text{RBF}}(X, X') = \sigma_f^2 \exp\left(-\frac{\|X - X'\|^2}{2l^2}\right)$$

Reason: The RBF kernel is smooth and infinitely differentiable, making it suitable for predicting smoothly varying functions. It is one of the most commonly used kernels.

702 **Matern Kernel (e.g., $\nu = \frac{3}{2}$ or $\nu = \frac{5}{2}$):**

$$703 \quad k_{\text{Matern}}(X, X') = \sigma_f^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|X - X'\|}{l} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu} \|X - X'\|}{l} \right)$$

704 where K_ν is the modified Bessel function.

705 **Reason:** The Matern kernel is more flexible than the RBF kernel, allowing control over the smoothness of the function, making it suitable for handling data with varying smoothness.

706 **Rational Quadratic Kernel:**

$$707 \quad k_{\text{RQ}}(X, X') = \sigma_f^2 \left(1 + \frac{\|X - X'\|^2}{2\alpha l^2} \right)^{-\alpha}$$

708 **Reason:** The Rational Quadratic kernel is a weighted sum of RBF kernels with different length scales, making it suitable for data with multi-scale features.

709 A.2.1 PARAMETER TUNING

710 The hyperparameters of the kernel functions were adjusted as follows:

- 711 • **Length Scale l :** Values tested: $l = 0.1, 1, 10$.
 712 **Reason:** The length scale controls the rate of change of the function. Smaller l values allow greater changes over shorter distances, suitable for capturing local features; larger l values result in smoother functions.
- 713 • **Signal Variance σ_f^2 :** Values tested: $\sigma_f^2 = 1, 5, 10$.
 714 **Reason:** The signal variance determines the amplitude of function variation. Adjusting σ_f^2 helps match the overall variability level of the data.
- 715 • **Noise Variance σ_n^2 :** Values tested: $\sigma_n^2 = 0.01, 0.1, 1$.
 716 **Reason:** The noise variance reflects the level of noise in the observed data. Adjusting σ_n^2 according to the actual data helps prevent overfitting or underfitting.

717 **Combining Kernels:** Multiple kernels were also combined, for example, a linear kernel plus an RBF kernel:

$$718 \quad k(X, X') = k_{\text{Linear}}(X, X') + k_{\text{RBF}}(X, X')$$

719 **Reason:** Combined kernels can capture different features in the data, such as global trends and local variations.

720 A.3 RESULTS EVALUATION METHODS

721 To evaluate the performance of the GPR models with different kernel functions, several metrics were used:

722 **Mean Squared Error (MSE):**

$$723 \quad \text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \mu_i)^2$$

724 **Purpose:** Measures the average squared difference between the predicted and actual values. The smaller the MSE, the better the model performance.

725 **Root Mean Squared Error (RMSE):**

$$726 \quad \text{RMSE} = \sqrt{\text{MSE}}$$

727 **Purpose:** RMSE is in the same units as the original data, making it easier to interpret the magnitude of the error.

Mean Absolute Error (MAE):

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_i - \mu_i|$$

Purpose: Measures the average absolute difference between the predicted and actual values, less sensitive to outliers than MSE.

R-squared (Coefficient of Determination):

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \mu_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

Purpose: Evaluates the proportion of variance in the dependent variable that is predictable from the independent variables. The closer R^2 is to 1, the better the model fits the data.

A.4 EXPERIMENTAL RESULTS AND ANALYSIS

The results for different kernel functions are summarized in Table 4. These results indicate the performance of each kernel function in terms of MSE, RMSE, MAE, and R^2 score.

Table 4: Performance of Different Kernels

Kernel	MSE	RMSE	MAE	R^2
RBF	0.6929	0.8324	0.6692	0.8764
Matern ($\nu = 1.5$)	0.3211	0.5666	0.4659	0.9427
Matern ($\nu = 2.5$)	0.4152	0.6443	0.5265	0.9259
Rational Quadratic	0.3199	0.5656	0.4639	0.9429

Discussion:

- The **RBF kernel** showed relatively poor performance with the highest MSE and RMSE, indicating it might not be flexible enough to handle the non-smooth data in this case.
- The **Matern kernel** ($\nu = 1.5$) performed very well with a low MSE of 0.3211 and a high R^2 of 0.9427, suggesting a good balance between model complexity and generalization ability.
- The **Matern kernel** ($\nu = 2.5$) also performed well, but slightly worse than $\nu = 1.5$, indicating its smoothness might not fully match the data characteristics.
- The **Rational Quadratic kernel** achieved the best performance with the lowest MSE (0.3199), RMSE (0.5656), and MAE (0.4639), indicating it is very well-suited to the multi-scale characteristics of the data.

A.5 CONCLUSION

Based on the experimental results, the **Rational Quadratic kernel** is selected as the optimal kernel function for the Gaussian Process Regression model. It shows the best fitting performance in terms of MSE, RMSE, and MAE metrics. The **Matern kernel** ($\nu = 1.5$) is also a recommended choice when a slightly smoother function is needed.

ISO 639-1 LANGUAGE CODES AND THEIR CORRESPONDING FULL NAMES.

Table 5: ISO 639-1 Language Codes and Their Corresponding Full Names. This table presents a comprehensive list of 25 major world languages, including their standardized two-letter codes (ISO 639-1) and full English names. Languages are arranged alphabetically by their ISO codes within each column.

Code	Full Name	Code	Full Name	Code	Full Name
zh	Chinese (Mandarin)	ja	Japanese	ta	Tamil
en	English	pa	Punjabi	it	Italian
es	Spanish	jv	Javanese	fa	Persian (Farsi)
hi	Hindi	de	German	ur	Urdu
ar	Arabic	ko	Korean	mr	Marathi
pt	Portuguese	fr	French	sw	Swahili
bn	Bengali	te	Telugu	th	Thai
ru	Russian	vi	Vietnamese	pl	Polish
		tr	Turkish		

DETAILED RESULTS ON TRANSLATION TASKS

BLEU (Bilingual Evaluation Understudy) is a traditional reference-based metric that measures the overlap of n-grams between the machine translation output and human reference translations. It calculates precision for different n-gram lengths (typically 1-4) and combines them using a geometric mean, along with a brevity penalty to penalize overly short translations. Despite its limitations, BLEU remains widely used as a baseline metric due to its simplicity and interpretability.

BLEURT (Bilingual Evaluation Understudy with Representations from Transformers) is a learned metric that leverages pre-trained language models to evaluate translation quality. It is trained on synthetic data and human judgments, using BERT-based architecture to capture semantic similarities between translations and references. BLEURT demonstrates stronger correlation with human judgments compared to traditional metrics, particularly in capturing meaning-preserving variations and grammatical nuances.

COMET (Crosslingual Optimized Metric for Evaluation of Translation) is a neural-based metric that utilizes multilingual pre-trained models to evaluate translation quality. It can operate in both reference-based and reference-free modes, making it particularly versatile. COMET is trained on human judgment data and employs a cross-lingual encoder to assess translation adequacy and fluency. It has shown superior performance in WMT Metrics Shared Tasks, demonstrating strong correlation with human evaluations across diverse language pairs.

Below are the performance trends across different models. We present the results for German (de), French (fr), and Russian (ru), which show significant improvements. Results for other languages did not demonstrate clear trends, and thus are not displayed.

Table 6: Performance Changes After Fine-tuning on LLaMA3-8B(Difference = After - Before)

Language	COMET	BLEURT	BLEU
Russian	0.0246	0.0345	0.0652
French	0.0082	0.0233	0.0640
German	0.0307	0.0178	0.1240

Table 7: Performance Changes After Fine-tuning on Phi3.5(Difference = After - Before)

Language	COMET	BLEURT	BLEU
Russian	0.0243	0.077	0.0645
French	0.064	0.0222	0.0686
German	0.034	0.0164	0.1411

Table 8: Performance Changes After Fine-tuning on LLaMA3.1-8B(Difference = After - Before)

Language	COMET	BLEURT	BLEU
Russian	0.0386	0.0297	0.0557
French	0.0166	0.0224	0.0648
German	0.0312	0.0153	0.1243

Table 9: Performance Changes After Fine-tuning on Qwen2-7B(Difference = After - Before)

Language	COMET	BLEURT	BLEU
Russian	0.0379	0.0333	0.0664
French	0.0135	0.0235	0.0745
German	0.0357	0.0136	0.1198

DETAILED DATA CORRESPONDING TO FIGURE 7

We present the improvements achieved by the LLM across 25 languages illustrated in Figure 7, with particular attention to low-resource languages. The results demonstrate that the model maintains significant enhancement in reasoning capabilities even in low-resource language contexts.

Table 10: LLaMA3-8B

hi	pt	bn	pa	jv	de	ko	th	pl
+19.5	+21.9	+21.6	+14.4	+21.7	+26.0	+21.8	+19.1	18.6
te	vi	tr	ta	it	fa	ur	mr	sw
+22.1	+23.4	+15.6	+11.5	+16.3	+15.7	+15.4	+23.5	+17.2

Table 11: LLaMA3-70B

hi	pt	bn	pa	jv	de	ko	th	pl
+17.5	+20.9	+23.6	+15.4	+20.7	+21.0	+21.4	+17.4	17.9
te	vi	tr	ta	it	fa	ur	mr	sw
+21.1	+20.4	+13.6	+14.5	+13.7	+16.3	+15.4	+20.7	+15.8