

# How Well Does Mathematical Reasoning Transfer Across Languages in English-Centric LLMs?

Anonymous ACL submission

## Abstract

Recent advances in Reinforcement Post-Training (RPT) have substantially improved the mathematical reasoning capabilities of large language models (LLMs), but it remains unclear how well such gains transfer across languages. In this work, we study the cross-lingual transfer of mathematical reasoning in English-centric LLMs under controlled multilingual evaluation. We systematically evaluate English-centric reasoning models on multilingual mathematical reasoning benchmarks, and the results show that cross-lingual transfer is highly variable across initial model choice, target language, and training paradigm. Through controlled comparative experiments, we further find that stronger English-centric initialization does not necessarily lead to stronger relative transfer efficiency, even when it yields better multilingual reasoning accuracy. Finally, we show that the largest improvement comes from moving beyond English-only supervision: adding a single parallel language yields a substantial gain, while further gains from additional languages are smaller but consistent. Overall, our results suggest that multilingual mathematical reasoning should be evaluated directly rather than inferred from English benchmarks alone.

## 1 Introduction

Recent advances in Reinforcement Post-Training (RPT) (Jaech et al., 2024; Kimi et al., 2025; Qwen, 2025) have substantially improved the reasoning capabilities of large language models (LLMs), especially on mathematical benchmarks. In particular, Reinforcement Learning with Verifiable Rewards (RLVR) has led to strong gains on datasets such as MATH (Hendrycks et al., 2021) and AIME (Maxwell, 2024; Kaggle, 2025) (Lambert et al., 2024; Guo et al., 2025). As these models continue to improve on English reasoning benchmarks, an important question is whether such gains

transfer beyond English. Prior work has mainly studied transfer across tasks or modalities (Chu et al., 2025; Liu et al., 2025a; Hu et al., 2025a; Huan et al., 2025; Zhou et al., 2025), while cross-lingual transfer remains less understood.

This question matters for both multilingual use and model analysis. Strong English reasoning performance does not guarantee equally strong performance for multilingual users, especially when users require both correct answers and interpretable reasoning traces in their own language. More broadly, cross-lingual transfer offers a way to probe what English-centric post-training learns: transferable reasoning behavior, language-specific patterns, or some combination of both. This motivates the central question of this paper:

*(Q) To what extent does mathematical reasoning acquired through English-centric post-training transfer across languages under controlled multilingual evaluation?*

We answer this question in three stages. We first conduct an observational study over open-source English-centric reasoning models and find that cross-lingual transfer is highly variable across initial model choice, target language, and training paradigm. We then run controlled comparative studies on initialization type, model family, and model size. Across these settings, stronger English-centric initialization does not necessarily correspond to larger relative cross-lingual transfer gains. Finally, we study parallel multilingual supervision and find that it substantially improves transfer, with the largest gain appearing when moving from monolingual to bilingual training.

Our contributions are threefold. First, we provide a systematic empirical evaluation of cross-lingual transfer in English-centric LLMs on multilingual mathematical reasoning benchmarks. Second, we show through controlled comparative experiments that transfer varies across initialization

and model choices under matched post-training settings. Third, we identify the transition from monolingual to bilingual supervision as the most effective intervention in our setting, with further multilingual gains that are positive but smaller.

We organize the paper around the following research questions:

- **RQ1:** To what extent do English-centric LLMs transfer mathematical reasoning across languages? (Section 2)
- **RQ2:** What factors are associated with stronger or weaker cross-lingual transfer under controlled post-training settings? (Section 3)
- **RQ3:** How can cross-lingual transfer be improved effectively? (Section 4)

## 2 Observational Study

To answer **RQ1**—*To what extent do English-centric LLMs transfer mathematical reasoning across languages?*—we begin with an observational study of open-source reasoning models. This analysis establishes whether cross-lingual transfer is present and examines how it varies with model initialization, target language, and training paradigm.

### 2.1 Setup

**Models.** We evaluate a diverse set of open-source reasoning models further tuned with either Supervised Fine-Tuning (SFT) or Reinforcement Post-Training (RPT), including the Simple-Zoo (Zeng et al., 2025), s1 (Muennighoff et al., 2025), OpenThinker (Guha et al., 2025), OpenReasoner-Zero (Hu et al., 2025b), and DeepSeek-R1-Distill (Guo et al., 2025) series.

**Benchmarks.** We evaluate on multilingual mathematical reasoning benchmarks, including the multilingual versions of MATH500 (Hendrycks et al., 2021), AIME2024 (Maxwell, 2024), and AIME2025 (Kaggle, 2025) from XReasoning (Qi et al., 2025). We also include multilingual GPQA-Diamond from BenchMAX (Huang et al., 2025b) as an out-of-domain reasoning benchmark. The evaluation includes *English (En)*, *Spanish (es)*, *Russian (ru)*, *German (de)*, *French (fr)*, *Bengali (bn)*, *Swahili (sw)*, *Thai (th)*, *Japanese (ja)*, *Chinese (zh)*, and *Telugu (te)*—yielding eleven evaluation languages in total. More details are given in appendix C.1.

**Language-controlled prompting.** Following prior work (Luo et al., 2025; Wang et al., 2025b; Qi et al., 2025), we concatenate the language control instructions after the input prompts to make the model generate responses using the same language as the query. This setup reflects multilingual application scenarios where both the answer and the reasoning trace should be understandable in the user’s language (Yong et al., 2025; Wang et al., 2025a; Zhang et al., 2025). The full prompts are provided in Appendix G.3.

**Metrics.** We report reasoning accuracy (**Acc**) as the main task metric, and off-target rate (**Off-target**) to measure how often a model fails to respond in the requested language, using the LangDetect library. In all evaluation settings, correctness is determined from the final answer only. Intermediate reasoning traces are not scored. Therefore, a model does not receive higher accuracy simply for producing its reasoning trace in the target language.

To quantify relative transfer, we also use the *Multilingual Transferability Index* (MTI), following prior work (Huan et al., 2025; Huang et al., 2025a). MTI is intended to measure how much of the gain obtained in the training language transfers to unseen languages. Since base difficulty varies across languages and benchmarks, absolute accuracy alone is not sufficient for this purpose.

Let  $Acc_{b,l}^{\text{trained}}$  and  $Acc_{b,l}^{\text{base}}$  denote the accuracy of a trained model and its corresponding base model on benchmark  $b$  and language  $l$ . We define the relative gain on language  $l$  as

$$\Delta R_{b,l} = \frac{Acc_{b,l}^{\text{trained}} - Acc_{b,l}^{\text{base}}}{Acc_{b,l}^{\text{base}}}. \quad (1)$$

For a training language set  $\mathcal{L}_{\text{train}}$ , the average relative gain is

$$\Delta R_{b,\mathcal{L}_{\text{train}}} = \frac{1}{|\mathcal{L}_{\text{train}}|} \sum_{l \in \mathcal{L}_{\text{train}}} \Delta R_{b,l}. \quad (2)$$

We define the MTI for an unseen language  $l_{\text{unseen}}$  as the average benchmark-level relative transfer ratio:

$$\text{MTI}_{l_{\text{unseen}}} = \frac{1}{|B|} \sum_{b \in B} \frac{\Delta R_{b,l_{\text{unseen}}}}{\Delta R_{b,\mathcal{L}_{\text{train}}}}. \quad (3)$$

An MTI of 1 indicates that the relative gain on an unseen language matches the average gain on the training language(s); values below 1 indicate weaker transfer, and values above 1 indicate

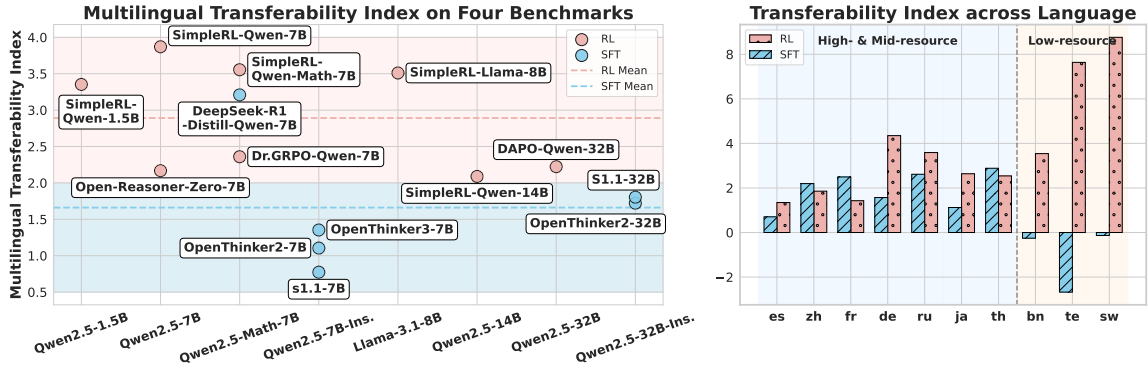


Figure 1: **Cross-lingual transfer across open-source reasoning models.** The *left* panel reports average MTI across four benchmarks and eleven languages. The *right* panel shows average transfer behavior on multilingual MATH500, comparing SFT-tuned and RL-tuned models.

stronger transfer. We use MTI only as an auxiliary measure of transfer efficiency and interpret it together with absolute accuracy and off-target rate. As a robustness check in Appendix E.1, we also repeated the main analyses with alternative transfer metrics, including normalized gain and raw gain ratio, and examined the sensitivity of MTI to low-base-accuracy cases. The numerical values differ, but the main qualitative conclusions remain unchanged.

## 2.2 Main Results

The main result of this section is that *cross-lingual transfer of mathematical reasoning is highly variable*. Figure 1 shows that strong English reasoning does not reliably translate into equally strong multilingual performance, and transfer varies across model initialization, training paradigm, and target language.

**Model choice matters.** Transfer varies substantially across models. Even among systems with similar training goals, multilingual performance differs markedly across languages. For example, SimpleRL-Qwen-7B shows slightly higher transferability than SimpleRL-Qwen-Math-7B despite similar training procedures. This variation is visible both across training paradigms and across initializations.

**English performance is not enough.** High English performance and strong multilingual transfer are related, but not identical. Some models that perform well in English do not maintain equally strong multilingual performance, while others with weaker English starting points show larger relative gains after post-training. English evaluation alone is therefore insufficient for characterizing multilin-

gual transfer.

**Training paradigm matters.** The clearest pattern is the difference between SFT and RL. RL-tuned models generally achieve stronger multilingual performance than SFT-tuned models, especially in lower-resource languages such as *bn*, *sw*, and *te*. In these languages, SFT often leads to weak or negative transfer, whereas RL more often yields positive gains.

**Target language matters.** Transfer is not uniform across languages. Higher-resource languages are generally more stable, whereas lower-resource languages show much larger variation across models and training paradigms. This suggests that cross-lingual transfer depends not only on English-side training, but also on properties of the target language.

## 3 Controlled Comparative Study

The observational study shows that cross-lingual transfer in English-centric reasoning models is highly variable, but open-source model comparisons alone cannot identify the main sources of this variation. Open-source models differ simultaneously in initialization, training data, optimization settings, and model family. To answer **RQ2**—*What factors are associated with stronger or weaker cross-lingual transfer under controlled post-training settings?*—we therefore conduct a set of controlled comparative studies.

Our goal is to test whether the variation observed in Section 2 persists under controlled post-training settings, and whether similar variation appears across initialization type, model family, and model size under matched post-training settings.

### 3.1 Setup

**Training data.** We use a compact training set of 1,000 examples sampled from the MATH training split, following prior work on small but carefully designed reasoning supervision sets (Ye et al., 2025; Muennighoff et al., 2025). All controlled studies in this section use the same training set unless otherwise noted. Details are provided in Appendix D.2.

**Post-training algorithm.** We use Group Relative Policy Optimization (GRPO) (Shao et al., 2024) as our RPT algorithm. Following prior work (Rastogi et al., 2025; Zhang et al., 2025; Liu et al., 2026), the policy is optimized with a composite reward that combines reasoning accuracy, response format, and language consistency:

$$R = \lambda_1 R_{\text{acc}} + \lambda_2 R_{\text{format}} + \lambda_3 R_{\text{lang}}. \quad (4)$$

Here,  $R_{\text{acc}}$  rewards correctness,  $R_{\text{format}}$  encourages well-formed reasoning traces, and  $R_{\text{lang}}$  encourages responses in the target language. Hyperparameter details are given in Appendix D.4 and E.3.1.

Because this objective includes formatting and language-consistency terms in addition to task accuracy, the resulting gains should not be interpreted as pure reasoning gains alone. We use the same objective across all settings so that the comparisons remain controlled.

**Controlled factors.** We study three factors: *initial model type*, *model family*, and *model size*. In each case, all training data and optimization settings are fixed.

### 3.2 Initial Model Type

We first study the effect of initial model type using three Qwen2.5-7B variants: a general base model, an instruction-tuned model, and a math-specialized model. Table 1 shows a trade-off between final multilingual performance and relative transfer efficiency. The instruction model achieves the highest multilingual accuracy after English-only post-training and the lowest off-target rate, but the base and math-specialized models obtain higher MTI. This suggests that stronger final multilingual performance and stronger relative transfer efficiency are not always aligned under our setup. Notably, this observation concerns relative transfer efficiency rather than final multilingual performance, which remains highest for the instruction-tuned model under our setup.

Model	Acc	Off-target	MTI
Qwen2.5-7B-Base	12.00	11.41	-
↪ GRPO on En Data	22.45	3.12	1.95
Qwen2.5-7B-Instruct	22.45	1.43	-
↪ GRPO on En Data	23.51	0.94	1.23
Qwen2.5-Math-7B	12.25	22.59	-
↪ GRPO on En Data	19.37	9.50	2.12

Table 1: **The Impact of Initial Model Type.** Average Accuracy (%), Off-target rate (%), and MTI across eleven languages on four benchmarks.

### 3.3 Model Family

We next compare Qwen2.5-7B-Instruct and Llama3.1-8B-Instruct under the same English-only GRPO setup to test whether the pattern above persists across model families.

Figure 2 shows that English-only GRPO improves performance in both English and unseen languages for both families. However, the relative gains differ substantially: although Llama3.1-8B-Instruct starts from weaker English performance than Qwen2.5-7B-Instruct, it exhibits larger relative gains across languages. This suggests that cross-lingual transfer is not simply determined by English performance alone. At the same time, this comparison does not isolate a pure family effect, since differences in tokenization, pretraining data, and architecture remain entangled. Notably, this comparison concerns relative gains rather than final multilingual accuracy, and should be interpreted accordingly.

### 3.4 Model Size

We finally compare Qwen2.5-1.5B-Instruct and Qwen2.5-7B-Instruct under the same English-only GRPO setup to study how the observed pattern changes with model scale.

Figure 3 shows that the effect of model scale is benchmark-dependent. On multilingual MATH500, the 1.5B model gains more than the 7B model, suggesting greater headroom on the in-domain task. On multilingual AIME24/25, however, the 7B model attains stronger final performance on the harder benchmarks despite smaller gains on MATH500. On multilingual GPQA-Diamond, the smaller model again shows larger improvements, while the 7B model improves only marginally and can even degrade in English. One possible explanation is that the small math-only training set over-specializes the 7B model toward mathematical reasoning, which can hurt out-of-domain scientific

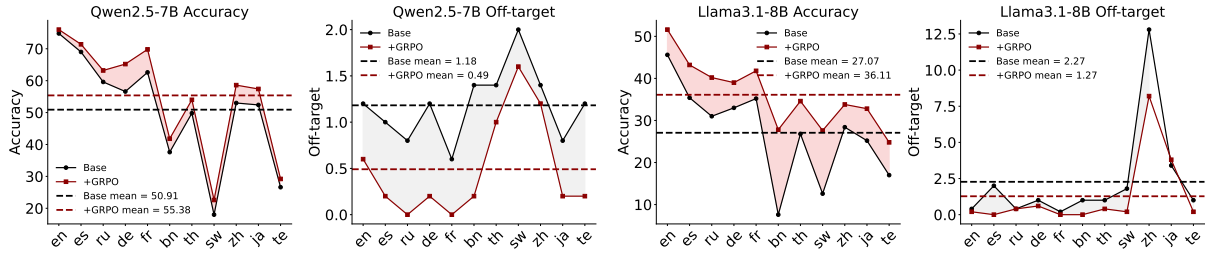


Figure 2: **Effect of model family under controlled post-training.** Multilingual reasoning performance on MATH500. *Base* denotes the initial model, and *+GRPO* the model after English-only GRPO.

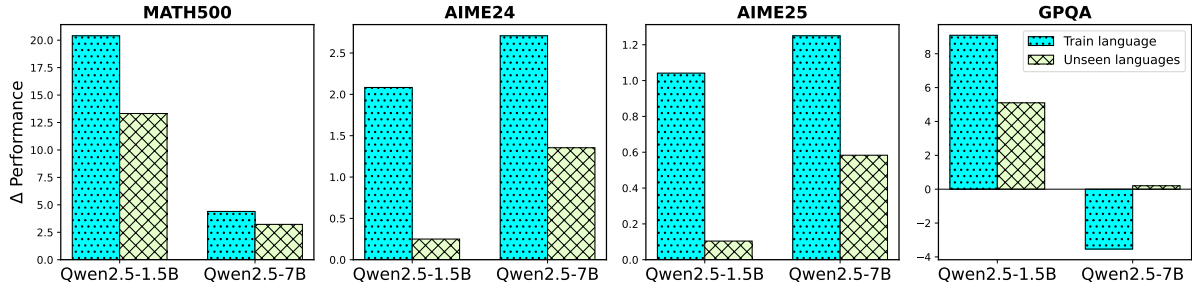


Figure 3: **Effect of model size under controlled post-training.**  $\Delta$ Performance denotes the average accuracy gain over the initial model, averaged across the training language and unseen languages.

QA. Overall, the observed transfer pattern depends on the interaction between initial model strength and benchmark difficulty, rather than on model size alone.

## 4 Parallel Training Study

The controlled comparative results suggest that cross-lingual transfer in English-centric post-training is unstable and sensitive to initialization and alignment choices. We therefore ask *RQ3*—*How can cross-lingual transfer be improved effectively?* In this section, we study a simple intervention: adding multilingual parallel supervision during post-training.

Our main finding is that the transition from monolingual to bilingual training produces the largest improvement in cross-lingual transfer. Further gains from adding more languages are positive but smaller, indicating that the key step is to move beyond English-only supervision.

### 4.1 Setup

We use Qwen2.5-7B-Instruct as the main initial model and post-train it with GRPO on English together with increasing numbers of parallel languages. The starting point is the same 1,000 English examples used in Section 3. For multilingual supervision, we use parallel problem sets drawn from the multilingual MATH dataset described in Section C.1, covering seven additional languages:

*es, ru, de, fr, bn, th, zh.* Following prior work, we further perform an LLM-based validation with Qwen3-32B to screen for translation errors in the multilingual dataset; details and summary statistics are provided in Appendix C.1. By varying the number of added languages from one to seven, we obtain a sequence of models that differ only in the amount of multilingual supervision. The full language configurations are given in Table 16. We evaluate these models on multilingual MATH500 across eleven languages, using absolute accuracy as the primary metric and MTI only as an auxiliary measure of relative transfer efficiency. Unless otherwise stated, the analyses in this section focus on multilingual MATH500, which provides the clearest view of the monolingual-to-bilingual transition under our setup.

**The Power of Bilingual Training** The most striking result in Figure 4 is that most of the improvement comes from the first parallel language. Moving from English-only to bilingual training yields a much larger gain than adding further languages afterward. Using Qwen2.5-7B-Instruct as the initial model, average accuracy rises from 54.24 to 57.87, while MTI increases from 1.16 to 2.50. By comparison, extending training from one parallel language to seven yields smaller incremental improvements.

We refer to this phenomenon as the *First-Parallel Leap*. In our setting, the key step is to move beyond English-only supervision: adding

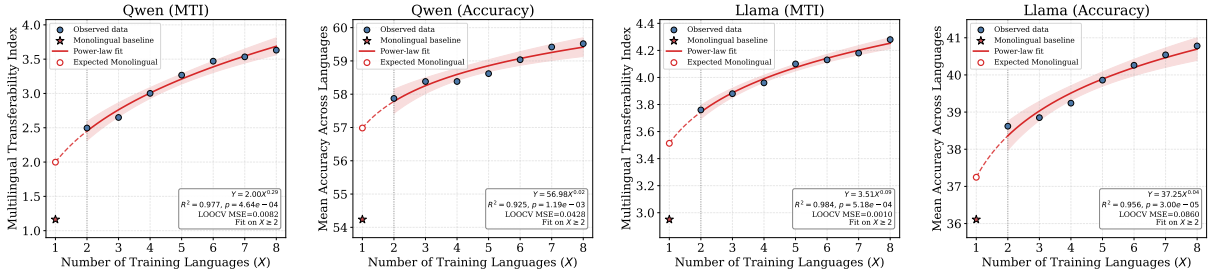


Figure 4: **Transfer improves most at the first parallel language.** Each panel shows performance as a function of the number of training languages. Blue points denote results obtained by fine-tuning on English together with additional parallel languages. The solid red curve shows a power-law-like fit over configurations with  $X \geq 2$ , and the shaded region indicates the 95% confidence interval. The dashed segment extrapolates the multilingual trend to the monolingual setting ( $X = 1$ ), which is excluded from fitting. The red star marks the observed monolingual baseline, and the hollow red circle denotes the extrapolated value at  $X = 1$ .

one parallel language matters more than increasing the number of languages from an already multilingual starting point.

**Beyond Bilingual: Positive but Diminishing Returns** Beyond the bilingual setting, performance continues to improve as more parallel languages are added, but the gains are smaller. Figure 4 shows a smooth sub-linear trend in both MTI and average accuracy. To summarize this pattern, we fit a power-law-like function

$$f(X) = \alpha X^\beta, \quad (5)$$

where  $X$  is the number of training languages. We fit only configurations with  $X \geq 2$  to separate the monolingual and multilingual regimes. The resulting fits are:

$$\begin{aligned} \text{Accuracy: } f(X) &= 56.98X^{0.02}, \\ \text{Transferability: } f(X) &= 2.00X^{0.29}. \end{aligned} \quad (6)$$

We use this *Parallel Scaling Trend* only as a descriptive summary of the bilingual-to-multilingual regime. It is not intended as a universal scaling law. The sub-linear exponents indicate diminishing returns as more languages are added.

The fitted trend also highlights a *Monolingual Extrapolation Gap*: when extrapolated back to  $X = 1$ , the observed English-only baseline falls below the value predicted from the multilingual regime in both MTI and accuracy. This again suggests that English-only post-training behaves differently from even weakly multilingual training.

## 4.2 Additional Analysis and Discussion

**Fixed training budget.** A possible alternative explanation is that parallel training helps simply because it uses more data. To test this, we fix the

total number of training examples and vary only the composition. Specifically, we compare scaling English-only data (e.g.,  $2N$  English examples) with replacing part of it by parallel multilingual data (e.g.,  $N$  English +  $N$  Russian). Table 2 shows that, under a sample-equivalent budget, multilingual mixtures yield stronger cross-lingual transfer than monolingual scaling alone.

Budget	Type	Acc	Off-target	MTI
1×	1×En	54.24	0.49	1.16
2×	2×En	55.82	0.58	1.92
2×	1×En + 1×Ru	57.87	0.20	2.50
3×	3×En	57.13	0.42	2.38
3×	1×En + 1×Ru + 1×Fr	58.38	0.24	2.65

Table 2: The Fixed-Budget ablation experiment based on Qwen2.5-7B-Instruct.

**Parallel vs. non-parallel data.** We also compare parallel bilingual supervision against non-parallel multilingual data. Parallel data provides aligned versions of the same problems, while non-parallel data only increases language exposure. Figure 6 shows that parallel data yields stronger gains, suggesting that explicit cross-lingual alignment is an important part of the bilingual gain.

**Does the chosen language matter?** A practical question is whether the gain depends strongly on which bilingual language is added. Figure 7 suggests that the answer is no at a coarse level: adding one parallel language consistently improves multilingual transfer, and the differences across languages such as Russian, Bengali, German, and Chinese are modest relative to the monolingual-to-bilingual jump. Some language-specific effects remain, but the overall benefit of parallel supervi-

sion is robust to the language choice.

**Scaling the number of training data** We also study whether the benefit of parallel training persists as the amount of training data increases. Figure 5 compares the 1K and 3K training data for each language under English-only, bilingual, and trilingual training settings. Increasing the training data improves both MTI and average accuracy, while the relative pattern remains stable: the largest gain comes from moving from monolingual to bilingual training, and adding a second parallel language yields a smaller additional improvement.

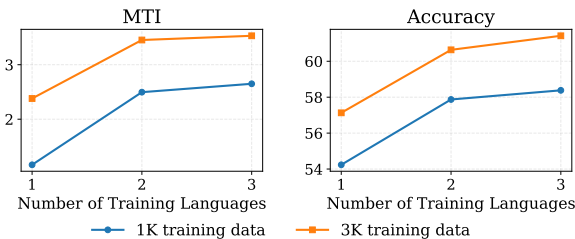


Figure 5: Scaling training data based on Qwen2.5-7B-Instruct.

**Prompt ablation.** A possible concern is that the bilingual gains in the main setting are driven by explicit language-control prompting. To test this, we compare two evaluation protocols: *language-controlled prompting*, used in the main experiments, and *no language control*, where all explicit target-language instructions are removed. We evaluate two Qwen2.5-7B-Instruct-based models: an *English-only* GRPO model and an *English+Ru* parallel-trained GRPO model. Table 3 shows that the bilingual gains are not reducible to test-time prompt control alone. Without language-control prompts, the English-only model exhibits severe language drift, whereas the parallel-trained model maintains near-zero off-target behavior and slightly higher accuracy. Under language-controlled evaluation, the parallel-trained model also remains stronger in accuracy.

Evaluation	Training	Acc	Off-target
No language control	English-only	57.82	75.76
No language control	English+Ru	58.43	0.36
Language-controlled	English-only	54.20	0.50
Language-controlled	English+Ru	57.90	0.20

Table 3: **Prompt ablation.** Average accuracy and off-target rate on multilingual MATH500 under two evaluation protocols.

**Extension across architectures and source languages.** We repeat the analysis on Llama3.1-8B-Instruct and observe the same qualitative pattern in Figure 4: the largest gain comes from moving beyond English-only training. We further test a Chinese-centric variant in Table 4. The same trend appears there as well: adding one parallel language yields the largest improvement, while adding another language provides a smaller additional gain. These results suggest that the bilingual gain is not specific to a single model family or source language.

Training Data	Acc	Off-target	MTI
Zh	55.00	0.38	0.86
Zh + Ru	58.13 <sub>(+3.13)</sub>	0.26 <sub>(-0.12)</sub>	1.69 <sub>(+0.83)</sub>
Zh + Ru + Fr	58.76 <sub>(+3.76)</sub>	0.23 <sub>(-0.15)</sub>	1.87 <sub>(+1.01)</sub>

Table 4: **Chinese-centric parallel training results** on multilingual MATH500.

**Robustness to transfer metric choice.** We further examine whether the observed monolingual-to-bilingual gain depends on the specific formulation of MTI. To this end, we repeat the main parallel-training analysis using two alternative transfer metrics, normalized-gain ratio and raw-gain ratio in Table 6, and additionally test the sensitivity of the results to low-base-accuracy cases by excluding unseen-language cases below different base-accuracy thresholds in Table 7. While the absolute values differ across metrics and filters, the main qualitative result remains unchanged: the largest improvement still comes from moving from monolingual to bilingual training. This suggests that the bilingual gain is not specific to the original MTI definition (Appendix E.1).

## 5 Related Work

**Reasoning Transfer** Recent large reasoning models (LRMs) have been driven by reinforcement-learning-based post-training. OpenAI’s O1 (Jaech et al., 2024) and DeepSeek-R1 (Guo et al., 2025) show that RL can substantially improve mathematical reasoning, especially when combined with verifiable or rule-based rewards (Shao et al., 2024). As these models improve, an important question is how well the acquired reasoning behavior transfers beyond the training setting. Prior work has mainly studied this problem across tasks, domains, or modalities. Hu et al. (2025a) show that RL improves structured reasoning but transfers less ef-

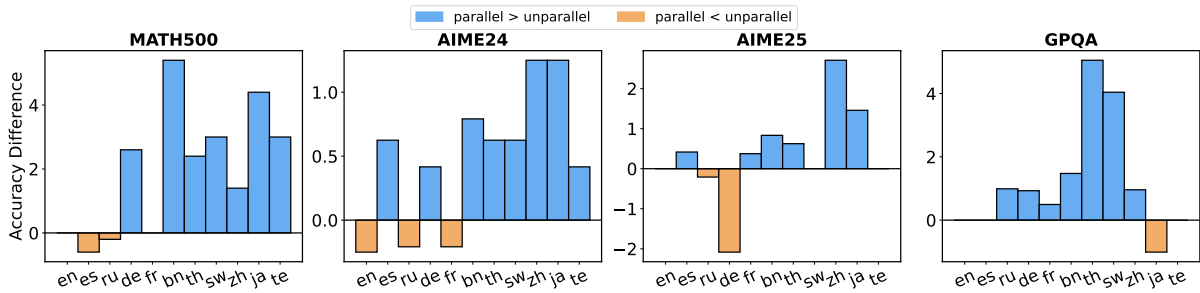


Figure 6: **Parallel vs. non-parallel multilingual supervision.** Accuracy differences under parallel and non-parallel data training based on Qwen2.5-7B-Instruct.

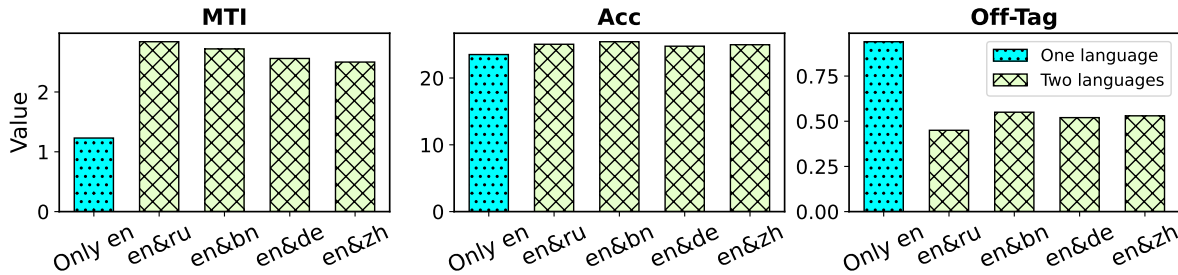


Figure 7: **The impact of the selected parallel language.** Multilingual reasoning performance across different single-language parallel training settings based on Qwen2.5-7B-Instruct.

512 fectively to unstructured tasks. Huan et al. (2025);  
 513 Chu et al. (2025) compare RL and SFT in terms  
 514 of broader transfer rather than narrow task-specific  
 515 gains, and X-REASONER (Liu et al., 2025a) stud-  
 516 ies transfer across domains and modalities. Our  
 517 work is complementary to this line of research:  
 518 instead of cross-task or cross-modal transfer, we  
 519 focus on cross-lingual transfer of mathematical rea-  
 520 soning.

521 **Cross-Lingual Transfer** Cross-lingual transfer  
 522 in English-centric LLMs has been widely studied  
 523 in multilingual NLP and, more recently, in align-  
 524 ment and reasoning settings. Prior work shows that  
 525 English-trained reward models (Wu et al., 2024;  
 526 Hong et al., 2025), preference alignment (Yang  
 527 et al., 2025b,a), and minimal multilingual adap-  
 528 tation (Li et al., 2024; Chirkova and Nikoulina,  
 529 2024) can often generalize beyond English. In  
 530 multilingual reasoning, Bandarkar et al. (2025)  
 531 transfers mathematical capability by combining  
 532 components from math-specialized and multilin-  
 533 gual models, Yong et al. (2025) shows that cross-  
 534 lingual test-time scaling improves multilingual rea-  
 535 soning, and Qi et al. (2025) highlights the role of  
 536 response-language control in multilingual evalua-  
 537 tion. Parallel multilingual data has also been shown  
 538 to improve multilingual capability in prompting  
 539 and pretraining settings (Mu et al., 2024; Qorib  
 540 et al., 2025). Huang et al. (2025a) shows that RL

541 training on non-English data yields better overall  
 542 performance and generalization than training on  
 543 English data. Our work differs from these studies  
 544 in two ways. First, we study cross-lingual transfer  
 545 in reasoning-oriented post-training rather than pre-  
 546 training or prompting alone. Second, beyond show-  
 547 ing that multilingual supervision helps, we identify  
 548 a sharper empirical pattern: the largest gain comes  
 549 from moving from monolingual to bilingual super-  
 550 vision, while additional languages yield smaller but  
 551 consistent improvements.

## 552 6 Conclusion

553 We study the cross-lingual transfer of mathemat-  
 554 ical reasoning in English-centric LLMs. Our re-  
 555 sults show that such transfer is real but highly  
 556 variable, and that strong English reasoning per-  
 557 formance does not reliably predict equally strong  
 558 transfer across languages. Through controlled com-  
 559 parative experiments, we further find that stronger  
 560 English-centric initialization does not necessarily  
 561 yield stronger relative transfer efficiency. We then  
 562 show that the strongest intervention is to move  
 563 beyond English-only supervision; adding a single  
 564 parallel language already yields most of the gain.  
 565 Overall, our findings suggest that multilingual rea-  
 566 soning transfer should be studied and evaluated  
 567 directly, rather than inferred from English bench-  
 568 marks alone.

## 569 Limitations

570 Our study has several limitations. First, although  
571 we extend the training data from 1K to 3K exam-  
572 ples and observe the same qualitative pattern, our  
573 controlled experiments remain small-scale com-  
574 pared with production-level reasoning post-training.  
575 The reported results should therefore be interpreted  
576 as evidence under this controlled setup rather than  
577 as claims about large-scale post-training dynamics.  
578 Second, our multilingual evaluation relies on trans-  
579 lated benchmarks, which may introduce language-  
580 dependent artifacts such as wording shifts, notation  
581 mismatch, or uneven difficulty across languages.  
582 Third, evaluation is conducted under language-  
583 controlled prompting, so the reported gains may re-  
584 flect not only task-solving ability, but also language  
585 compliance and response-format stability. In addi-  
586 tion, off-target rate is measured with LangDetect  
587 and should be interpreted as a proxy for response-  
588 language control rather than a perfect measure  
589 of multilingual reasoning validity, especially for  
590 mixed-language mathematical outputs. Finally,  
591 MTI is used only as an auxiliary metric, since it  
592 can be unstable in low-accuracy regimes, and our  
593 controlled results establish empirical associations  
594 rather than a definitive causal mechanism.

## 595 Ethical Considerations

596 This work studies the cross-lingual transfer of math-  
597 ematical reasoning in English-centric LLMs. Our  
598 contributions are primarily empirical and method-  
599 ological: we analyze transfer behavior under con-  
600 trolled multilingual evaluation and study training  
601 interventions that may improve multilingual rea-  
602 soning performance. We use publicly available  
603 benchmarks that do not contain sensitive personal  
604 information to the best of our knowledge. At the  
605 same time, multilingual evaluation based on trans-  
606 lated benchmarks may overstate or understate per-  
607 formance in some languages if translation artifacts,  
608 notation mismatch, or mixed-language mathemat-  
609 ical responses are not carefully accounted for. Im-  
610 provements in response-language control should  
611 therefore not be interpreted as complete evidence  
612 of robust multilingual reasoning in real-world set-  
613 tings.

## 614 References

615 Lucas Bandarkar, Benjamin Muller, Pritish Yuvraj, Rui  
616 Hou, Nayan Singhal, Hongjiang Lv, and Bing Liu.

2025. Layer swapping for zero-shot cross-lingual  
transfer in large language models. In *The Thirteenth  
International Conference on Learning Representa-  
tions*. 617  
618  
619  
620

Nadezhda Chirkova and Vassilina Nikoulina. 2024. 621  
Zero-shot cross-lingual transfer in instruction tun- 622  
ing of large language models. In *Proceedings of 623  
the 17th International Natural Language Generation 624  
Conference*, pages 695–708. 625

Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang 626  
Tong, Saining Xie, Dale Schuurmans, Quoc V Le, 627  
Sergey Levine, and Yi Ma. 2025. Sft memorizes, 628  
rl generalizes: A comparative study of foundation 629  
model post-training. In *Forty-second International 630  
Conference on Machine Learning*. 631

Etash Guha, Ryan Marten, Sedrick Keh, Negin Raouf, 632  
Georgios Smyrnis, Hritik Bansal, Marianna Nezhurina, 633  
Jean Mercat, Trung Vu, Zayne Sprague, et al. 634  
2025. Openthoughts: Data recipes for reasoning 635  
models. *arXiv preprint arXiv:2506.04178*. 636

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, 637  
Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, 638  
Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: In- 639  
centivizing reasoning capability in llms via reinforce- 640  
ment learning. *arXiv preprint arXiv:2501.12948*. 641

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul 642  
Arora, Steven Basart, Eric Tang, Dawn Song, and 643  
Jacob Steinhardt. 2021. Measuring mathematical 644  
problem solving with the math dataset. In *Thirty- 645  
fifth Conference on Neural Information Processing 646  
Systems Datasets and Benchmarks Track (Round 2)*. 647

Jiwoo Hong, Noah Lee, Rodrigo Martínez-Castaño, 648  
César Rodríguez, and James Thorne. 2025. Cross- 649  
lingual transfer of reward models in multilingual 650  
alignment. In *Proceedings of the 2025 Conference 651  
of the Nations of the Americas Chapter of the Asso- 652  
ciation for Computational Linguistics: Human Lan- 653  
guage Technologies (Volume 2: Short Papers)*, pages 654  
82–94. 655

Chuxuan Hu, Yuxuan Zhu, Antony Kellermann, Caleb 656  
Biddulph, Suppakit Waiwitlikhit, Jason Benn, and 657  
Daniel Kang. 2025a. Breaking barriers: Do reinforce- 658  
ment post training gains transfer to unseen domains? 659  
*arXiv preprint arXiv:2506.19733*. 660

Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, 661  
Xiangyu Zhang, and Heung-Yeung Shum. 2025b. 662  
Open-reasoner-zero: An open source approach to 663  
scaling up reinforcement learning on the base model. 664  
*arXiv preprint arXiv:2503.24290*. 665

Maggie Huan, Yuetai Li, Tuney Zheng, Xiaoyu Xu, 666  
Seungone Kim, Minxin Du, Radha Poovendran, Gra- 667  
ham Neubig, and Xiang Yue. 2025. Does math rea- 668  
soning improve general llm capabilities? understand- 669  
ing transferability of llm reasoning. *arXiv preprint 670  
arXiv:2507.00432*. 671

672	Shulin Huang, Yiran Ding, Junshu Pan, and Yue Zhang. 2025a. Beyond english-centric training: How reinforcement learning improves cross-lingual reasoning in llms. <i>arXiv preprint arXiv:2509.23657</i> .		
673			
674			
675			
676	Xu Huang, Wenhao Zhu, Hanxu Hu, Conghui He, Lei Li, Shujian Huang, and Fei Yuan. 2025b. Benchmax: A comprehensive multilingual evaluation suite for large language models. <i>arXiv preprint arXiv:2502.07346</i> .		
677			
678			
679			
680			
681	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. 2024. Openai o1 system card. <i>arXiv preprint arXiv:2412.16720</i> .		
682			
683			
684			
685			
686	Kaggle. 2025. <a href="#">Aime2025</a> .		
687	Team Kimi, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. 2025. Kimi k1. 5: Scaling reinforcement learning with llms. <i>arXiv preprint arXiv:2501.12599</i> .		
688			
689			
690			
691			
692	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .		
693			
694			
695			
696			
697			
698			
699	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. 2024. Tulu 3: Pushing frontiers in open language model post-training. <i>arXiv preprint arXiv:2411.15124</i> .		
700			
701			
702			
703			
704			
705	Chong Li, Wen Yang, Jiajun Zhang, Jinliang Lu, Shaonan Wang, and Chengqing Zong. 2024. X-instruction: Aligning language model in low-resource languages with self-curated cross-lingual instructions. In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 546–566.		
706			
707			
708			
709			
710			
711	Junxiao Liu, Zhijun Wang, Yixiao Li, Zhejian Lai, Liqian Huang, Xin Huang, Xue Han, Junlan Feng, and Shujian Huang. 2026. Self-improving multilingual long reasoning via translation-reasoning integrated training. <i>arXiv preprint arXiv:2602.05940</i> .		
712			
713			
714			
715			
716	Qianchu Liu, Sheng Zhang, Guanghui Qin, Timothy Ossowski, Yu Gu, Ying Jin, Sid Kiblawi, Sam Preston, Mu Wei, Paul Vozila, et al. 2025a. X-reasoner: Towards generalizable reasoning across modalities and domains. <i>arXiv preprint arXiv:2505.03981</i> .		
717			
718			
719			
720			
721	Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. 2025b. Understanding r1-zero-like training: A critical perspective. <i>arXiv preprint arXiv:2503.20783</i> .		
722			
723			
724			
725	Wenyang Luo, Wayne Xin Zhao, Jing Sha, Shijin Wang, and Ji-Rong Wen. 2025. Mmath: A multilingual benchmark for mathematical reasoning. <i>Preprint</i> .		
726			
727			
	Jia Maxwell. 2024. <a href="#">Aime2024</a> .		728
	Yongyu Mu, Peinan Feng, Zhiquan Cao, Yuzhang Wu, Bei Li, Chenglong Wang, Tong Xiao, Kai Song, Tongran Liu, Chunliang Zhang, et al. 2024. Revealing the parallel multilingual learning within large language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 6976–6997.		729 730 731 732 733 734 735
	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling. <i>arXiv preprint arXiv:2501.19393</i> .		736 737 738 739 740
	Jirui Qi, Shan Chen, Zidi Xiong, Raquel Fernández, Danielle S Bitterman, and Arianna Bisazza. 2025. When models reason in your language: Controlling thinking trace language comes at the cost of accuracy. <i>arXiv preprint arXiv:2505.22888</i> .		741 742 743 744 745
	Muhammad Reza Qorib, Junyi Li, and Hwee Tou Ng. 2025. Just go parallel: Improving the multilingual capabilities of large language models. <i>arXiv preprint arXiv:2506.13044</i> .		746 747 748 749
	Team Qwen. 2025. Qwq-32b: Embracing the power of reinforcement learning.		750 751
	Abhinav Rastogi, Albert Q Jiang, Andy Lo, Gabrielle Berrada, Guillaume Lample, Jason Rute, Joep Barmantlo, Karmesh Yadav, Kartik Khandelwal, Khyathi Raghavi Chandu, et al. 2025. Magistral. <i>arXiv preprint arXiv:2506.10910</i> .		752 753 754 755 756
	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2024. Gpqa: A graduate-level google-proof q&a benchmark. In <i>First Conference on Language Modeling</i> .		757 758 759 760 761
	John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust region policy optimization. In <i>International conference on machine learning</i> , pages 1889–1897. PMLR.		762 763 764 765
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .		766 767 768 769 770
	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. Hybridflow: A flexible and efficient rlhf framework. <i>arXiv preprint arXiv:2409.19256</i> .		771 772 773 774 775
	Mingyang Wang, Lukas Lange, Heike Adel, Yunpu Ma, Jannik Strötgen, and Hinrich Schütze. 2025a. Language mixing in reasoning language models: Patterns, impact, and internal causes. <i>arXiv preprint arXiv:2505.14815</i> .		776 777 778 779 780

- 781 Yiming Wang, Pei Zhang, Jialong Tang, Haoran Wei,  
782 Baosong Yang, Rui Wang, Chenshu Sun, Feitong  
783 Sun, Jiran Zhang, Junxuan Wu, et al. 2025b. Poly-  
784 math: Evaluating mathematical reasoning in multi-  
785 lingual contexts. *arXiv preprint arXiv:2504.18428*.
- 786 Zhaofeng Wu, Ananth Balashankar, Yoon Kim, Jacob  
787 Eisenstein, and Ahmad Beirami. 2024. Reuse your  
788 rewards: Reward model transfer for zero-shot cross-  
789 lingual alignment. In *Proceedings of the 2024 Con-  
790 ference on Empirical Methods in Natural Language  
791 Processing*, pages 1332–1353.
- 792 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,  
793 Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong  
794 Tu, Jingren Zhou, Junyang Lin, et al. 2024. Qwen2.  
795 5-math technical report: Toward mathematical ex-  
796 pert model via self-improvement. *arXiv preprint  
797 arXiv:2409.12122*.
- 798 Wen Yang, Junhong Wu, Chen Wang, Chengqing Zong,  
799 and Jiajun Zhang. 2025a. Implicit cross-lingual re-  
800 warding for efficient multilingual preference align-  
801 ment. *arXiv preprint arXiv:2503.04647*.
- 802 Wen Yang, Junhong Wu, Chen Wang, Chengqing Zong,  
803 and Jiajun Zhang. 2025b. Language imbalance  
804 driven rewarding for multilingual self-improving. In  
805 *The Thirteenth International Conference on Learning  
806 Representations*.
- 807 Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie  
808 Xia, and Pengfei Liu. 2025. Limo: Less is more for  
809 reasoning. *arXiv preprint arXiv:2502.03387*.
- 810 Zheng-Xin Yong, M Farid Adilazuarda, Jonibek  
811 Mansurov, Ruochen Zhang, Niklas Muennighoff,  
812 Carsten Eickhoff, Genta Indra Winata, Julia Kreutzer,  
813 Stephen H Bach, and Alham Fikri Aji. 2025.  
814 Crosslingual reasoning through test-time scaling.  
815 *arXiv preprint arXiv:2505.05408*.
- 816 Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,  
817 Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu,  
818 Lingjun Liu, Xin Liu, et al. 2025. Dapo: An open-  
819 source llm reinforcement learning system at scale.  
820 *arXiv preprint arXiv:2503.14476*.
- 821 Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Ke-  
822 qing He, Zejun Ma, and Junxian He. 2025. Simplerl-  
823 zoo: Investigating and taming zero reinforcement  
824 learning for open base models in the wild. *arXiv  
825 preprint arXiv:2503.18892*.
- 826 Xue Zhang, Yunlong Liang, Fandong Meng, Songming  
827 Zhang, Kaiyu Huang, Yufeng Chen, Jinan Xu, and Jie  
828 Zhou. 2025. Think natively: Unlocking multilingual  
829 reasoning with consistency-enhanced reinforcement  
830 learning. *arXiv preprint arXiv:2510.07300*.
- 831 Ruochen Zhou, Minrui Xu, Shiqi Chen, Junteng Liu,  
832 Yunqi Li, Xinxin Lin, Zhengyu Chen, and Junxian He.  
833 2025. Does learning mathematical problem-solving  
834 generalize to broader reasoning? *arXiv preprint  
835 arXiv:2507.04391*.

836

## Appendix

837

**A Reproducibility Statement** 13

838

**B The Usage of AI Assistants** 13

839

**C Evaluation Details and Setup** 13

840

C.1 Multilingual Reasoning Benchmarks . . . . . 13

841

C.2 An Overview of Open-source LRMs . . . . . 13

842

**D Implementation Details** 13

843

D.1 GRPO Algorithm . . . . . 13

844

D.2 Training Dataset . . . . . 14

845

D.3 Experimental Environments . . . . . 14

846

D.4 Hyperparameters . . . . . 14

847

**E Detailed Results and Analysis** 15

848

E.1 Robustness of Transfer Metrics . . . . . 15

849

E.2 Observational Study . . . . . 15

850

E.2.1 Template Choice for Base Models . . . . . 15

851

E.2.2 Detailed Results for Initial Model Types . . . . . 16

852

E.2.3 Detailed Results for Open-source Models . . . . . 16

853

E.3 Controlled Comparative Study . . . . . 16

854

E.3.1 Reward Hyperparameter Sensitivity . . . . . 16

855

E.3.2 Detailed Results for Initial Model Types . . . . . 16

856

E.3.3 Detailed Results for Model Family . . . . . 16

857

E.3.4 Detailed Results for Model Size . . . . . 16

858

E.4 Parallel Training Results . . . . . 20

859

E.4.1 Parallel Training Language Settings . . . . . 20

860

E.4.2 Detailed Results of Parallel Training . . . . . 20

861

E.4.3 Effect of the Selected Parallel Language . . . . . 20

862

E.5 Parallel Training from a Chinese-Centric Perspective . . . . . 20

863

**F Dataset License** 20

864

**G Prompts Template** 23

865

G.1 Templates for Base Model . . . . . 23

866

G.2 Multilingual Reasoning Instruction . . . . . 23

867

G.3 Prompt hacking to force response language . . . . . 23

868

G.4 Template for R1-like Reasoning . . . . . 23

## A Reproducibility Statement

Codes and model weights will be released after review to facilitate future research. For evaluation, we follow prior works and report averaged results over 16 sampled generations per question on data-scarce benchmarks. All evaluations are conducted with temperature set to 0.6 and top-p to 0.95, with the random seed fixed to ensure deterministic outputs across runs. Note that minor variations in inference results may still occur due to differences in hardware or the version of the inference framework.

## B The Usage of AI Assistants

We declare that the AI Assistants (ChatGPT and Gemini) were only used to refine the fluency of certain sentences during the writing of this paper. Every sentence polished with the LLM was carefully reviewed and approved by the authors. The LLM was not used for any other part of this research.

## C Evaluation Details and Setup

### C.1 Multilingual Reasoning Benchmarks

We use the multilingual version of these four reasoning benchmarks provided in (Qi et al., 2025), which use GPT-4o (Jaech et al., 2024) to translate all questions into the other ten languages *Spanish (es)*, *Russian (ru)*, *German (de)*, *French (fr)*, *Bengali (bn)*, *Swahili (sw)*, *Thai (th)*, *Japanese (ja)*, *Chinese (zh)*, and *Telugu (te)*, resulting in a total of eleven languages for evaluation.

**MATH500** The MATH500 (Hendrycks et al., 2021) benchmark assesses the mathematical reasoning and problem-solving abilities of language models, addressing the need for more challenging evaluations as their general capabilities advance. It consists of 500 problems across five core mathematical domains: algebra, combinatorics, geometry, number theory, and precalculus. Each problem is designed to test multi-step reasoning and complex problem-solving skills, going beyond simple calculations or factual recall.

**AIME24&25** The AIME24 (Maxwell, 2024) and AIME25 (Kaggle, 2025) datasets contain problems from the American Invitational Mathematics Examination (AIME) for 2024 and 2025, respectively. AIME is a prestigious high school mathematics

competition renowned for its challenging problems, consisting of 30 questions.

**GPQA-Diamond** GPQA-Diamond (Rein et al., 2024) consists of 198 multiple-choice questions across biology, chemistry, and physics, with difficulty levels ranging from challenging undergraduate to postgraduate. It is the highest quality subset, which includes only questions where both experts answer correctly and the majority of non-experts answer incorrectly.

**Translation Validation.** To assess the quality of the translated multilingual benchmarks, we perform a full-coverage LLM-based validation following prior work (Liu et al., 2026). For each benchmark and each target language, we compare every translated question against its English source using Qwen3-32B as a judge. This validation is conducted per question rather than by sampling. The judge scores each translation along four dimensions: semantic fidelity, mathematical consistency, solvability preservation, and fluency/naturalness. Based on these scores, it assigns an overall label of *Valid*, *Minor issue*, or *Major issue*, together with a short rationale and an error type when applicable (e.g., semantic shift, omission, mathematical inconsistency, or ambiguity). Across all checked translations, 99.2% are labeled as *Valid*, 0.7% as *Minor issue*, and 0.1% as *Major issue*. We use this procedure to screen for obvious translation errors and to verify that the translated benchmarks preserve the original mathematical problems at scale.

### C.2 An Overview of Open-source LRMs

Table 5 provides an overview of the various open-source LLMs evaluated in our observational study. These models, which include the DeepSeek-R1-Distill-Qwen-7B (Guo et al., 2025), OpenThinker series (Guha et al., 2025), Simple-RL-Zoo series (Zeng et al., 2025), s1 series (Muennighoff et al., 2025), and DAPO-Qwen-32B (Yu et al., 2025), range in size from 1.5B to 32B.

## D Implementation Details

### D.1 GRPO Algorithm

GRPO is a simplified PPO-based algorithm that significantly reduces training costs by eliminating the need for a value model. It operates by sampling  $G$  rollouts  $\{o_1, \dots, o_G\}$  from the current policy for a given input, calculating their cumulative rewards  $R = \{R_1, \dots, R_G\}$ , and then using these rewards

Model	Initial Model	Size	Training Paradigm
DeepSeek-R1-Distill-Qwen-7B	Qwen2.5-Math-7B-Base	7B	SFT
Open-Reasoner-Zero-7B	Qwen2.5-7B-Base	7B	RL
OpenThinker2-7B	Qwen2.5-7B-Instruct	7B	SFT
OpenThinker3-7B	Qwen2.5-7B-Instruct	7B	SFT
Qwen-2.5-1.5B-SimpleRL-Zoo	Qwen2.5-1.5B-Base	1.5B	RL
Qwen-2.5-7B-SimpleRL-Zoo	Qwen2.5-7B-Base	7B	RL
Qwen-2.5-14B-SimpleRL-Zoo	Qwen2.5-14B-Base	14B	RL
Qwen-2.5-Math-7B-SimpleRL-Zoo	Qwen2.5-Math-7B-Base	7B	RL
Qwen2.5-Math-7B-Dr.GRPO	Qwen2.5-Math-7B-Base	7B	RL
s1.1-7B	Qwen2.5-7B-Instruct	7B	SFT
DAPO-Qwen-32B	Qwen2.5-32B-Base	32B	RL
OpenThinker2-32B	Qwen2.5-32B-Instruct	32B	SFT
s1.1-32B	Qwen2.5-32B-Instruct	32B	SFT

Table 5: **The Overview of the Open-source LLMs Used in Observational Study**, including their initial model, parameter size, and training paradigm.

to estimate advantages  $\hat{A}_{i,t}$  to guide policy updates. The optimization objective for GRPO is defined as follows:

$$L_{\text{GRPO}}(\theta) = \mathbb{E}_{q \sim \mathcal{D}, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|q)} \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left( \mathcal{L}_{i,t}^{\text{clip}}(\theta) - \beta D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) \right) \right] \quad (7)$$

where

$$\mathcal{L}_{i,t}^{\text{clip}}(\theta) = \min \left( r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip} \left( r_{i,t}(\theta), 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_{i,t} \right) \\ r_{i,t}(\theta) = \frac{\pi_{\theta}(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t} | q, o_{i,<t})}, \hat{A}_{i,t} = \frac{R_i - \text{mean}(R)}{\text{std}(R)} \quad (8)$$

The clipping term with ratio  $\varepsilon$  (Schulman et al., 2015) keeps the new policy close to the old one, improving training stability.

## D.2 Training Dataset

**The Distribution of Parallel Questions** Figure 8a shows the type and level distributions of the 1,000 English training questions sampled from the MATH dataset (Hendrycks et al., 2021). The type distribution is relatively balanced, and the number of questions increases steadily from Level 1 to Level 5.

### The Distribution of non-parallel Questions

Moreover, Figure 8b presents the type and level distributions of 1,000 Russian questions used for a non-parallel training analysis experiment, which form a separate, non-overlapping set from the 1000 English questions. The distributions of both type and level closely match those of the English training questions. This indicates that, in the analysis comparing parallel and non-parallel training, the

performance drop observed in non-parallel training is not due to distributional differences between the non-parallel and English datasets.

## D.3 Experimental Environments

Training and inference were conducted on Ubuntu 22.04 with  $8 \times$  NVIDIA A800 GPUs. We use VeRL (Sheng et al., 2024) for RL training, vLLM 0.8.5 (Kwon et al., 2023) for inference, and Qwen’s Math evaluation codebase (Yang et al., 2024) following prior work (Zeng et al., 2025; Liu et al., 2025b).

## D.4 Hyperparameters

**RL Training** The maximum generation length was set to 4096 tokens, and the maximum prompt length to 1024 tokens, such that their sum matches the model’s maximum context length. The learning rate was fixed at  $1 \times 10^{-6}$ . Training was performed with a batch size of 128 questions. For each question,  $G = 16$  rollouts were sampled, using a sampling temperature of 1.0.  $\lambda_1 = 0.8$ ,  $\lambda_2 = 0.1$ , and  $\lambda_3 = 0.1$ .

**Inference** In the evaluation setup, we used a temperature of 0.6, a top- $p$  value of 0.95, and a maximum generation length of 8912 tokens for all models in the 1.5B–14B series. For 32B models, we used the same temperature (0.6) and top- $p$  value (0.95), but set the maximum generation length to 32,768 tokens, except for DAPO-Qwen-32B, which followed the official recommended settings: a temperature of 1.0, a top- $p$  value of 0.7, and a maximum generation length of 20,480 tokens. For AIME2024 and AIME2025, we report accuracy by

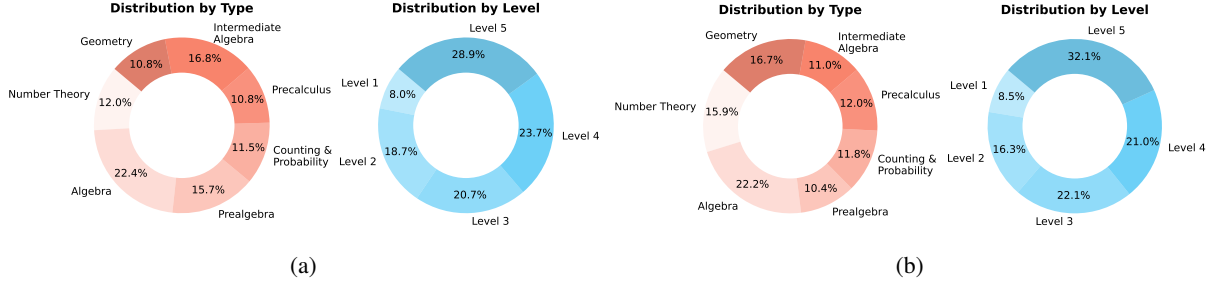


Figure 8: Distribution of question difficulty. (a) The 1,000 English questions were utilized in the controlled comparative study and for parallel training. (b) The 1,000 Russian questions for non-parallel training, comprising a separate and non-overlapping set from the questions in (a).

averaging over 16 sampled generations per question, while for MATH500 and GPQA, accuracy is computed using a single sampled generation per question.

## E Detailed Results and Analysis

### E.1 Robustness of Transfer Metrics

To examine whether our conclusions depend on the specific choice of MTI, we repeated the main analyses using two alternative transfer metrics.

First, we consider a normalized-gain variant:

$$\Delta N_{b,l} = \frac{Acc_{b,l}^{\text{trained}} - Acc_{b,l}^{\text{base}}}{1 - Acc_{b,l}^{\text{base}}}, \quad (9)$$

and define the corresponding transfer ratio by replacing  $\Delta R$  with  $\Delta N$  in the MTI computation.

Second, we consider a raw-gain ratio:

$$\Delta G_{b,l} = Acc_{b,l}^{\text{trained}} - Acc_{b,l}^{\text{base}}, \quad (10)$$

and compare the gain on an unseen language with the average gain on the training language(s):

$$RG_{b,l_{\text{unseen}}} = \frac{\Delta G_{b,l_{\text{unseen}}}}{\Delta G_{b,\mathcal{L}_{\text{train}}}}. \quad (11)$$

Table 6 summarizes the corresponding results on the main parallel-training comparison. Although the absolute values differ across metrics, the main qualitative pattern remains unchanged: moving from English-only to bilingual training yields the largest gain, and adding another parallel language provides a smaller additional improvement. This supports the main-text conclusion that the First-Parallel Leap is not specific to the original MTI definition.

Setting	MTI	NGR	RGR	Interpretation
Only English	1.163	0.396	0.732	<i>Baseline</i>
w. One parallel	2.496	0.963	1.480	<i>Large bilingual jump</i>
w. Two parallel	2.650	1.061	1.601	<i>Smaller additional gain</i>

Table 6: **Robustness of transfer metrics on the main parallel-training comparison.** We compare the original MTI with two alternative transfer metrics: normalized-gain ratio (NGR) and raw-gain ratio (RGR).

**Sensitivity to base accuracy.** Because MTI is based on relative gain, it can become sensitive when base accuracy is low. To assess whether the main result is driven by such cases, we recompute the transfer scores on multilingual MATH500 after excluding unseen-language cases whose base accuracy falls below different thresholds. Table 7 shows that the absolute values vary across metrics and filters, especially for the original relative-gain MTI. However, the main qualitative result remains unchanged: the transition from monolingual to bilingual training yields a clear improvement under all transfer metrics and filtering thresholds. This suggests that the bilingual gain is not an artifact of the lowest-base-accuracy cases alone.

### E.2 Observational Study

#### E.2.1 Template Choice for Base Models

To evaluate base models consistently, we compare three template settings for *Qwen2.5-7B-Base* and *Qwen2.5-Math-7B-Base*: Qwen-Math Template, Qwen-Instruct Template, and No Template. As shown in Table 8, the Qwen-Instruct Template yields the best overall reasoning accuracy and target-language consistency on multilingual MATH500. We therefore use it as the default template for evaluating all general-base and math-base models.

Setting	All	$Acc^{base} > 20$	$Acc^{base} > 30$
<i>Relative-gain MTI</i>			
Only English	1.16	1.16	1.23
+ One parallel	2.50	2.42	1.88
+ Two parallel	2.65	2.62	2.13
+ Three parallel	3.00	3.16	2.58
+ Four parallel	3.28	3.25	2.34
+ Five parallel	3.48	3.41	2.02
<i>Normalized-gain ratio</i>			
Only English	0.40	0.40	0.48
+ One parallel	0.96	0.95	1.05
+ Two parallel	1.06	1.08	1.22
+ Three parallel	1.15	1.19	1.37
+ Four parallel	1.18	1.20	1.33
+ Five parallel	1.63	1.57	1.85
<i>Raw-gain ratio</i>			
Only English	0.84	0.84	0.85
+ One parallel	1.45	1.47	1.46
+ Two parallel	1.61	1.67	1.68
+ Three parallel	1.95	2.04	1.97
+ Four parallel	1.85	1.88	1.81
+ Five parallel	2.06	2.11	2.08

Table 7: **Sensitivity of transfer scores to low-base-accuracy cases on multilingual MATH500.** We recompute the main transfer scores after excluding unseen-language cases whose base accuracy falls below different thresholds.

## E.2.2 Detailed Results for Initial Model Types

Table 9 reports the full multilingual accuracy and off-target results for the initial models used in our analysis. Two patterns are worth noting. First, within the Qwen2.5-7B series, the instruction-tuned model has the lowest off-target rate and the highest final multilingual accuracy. Second, within the Qwen2.5-Base series, larger models generally achieve higher multilingual accuracy and lower off-target rates. These results support the main-text observation that stronger final multilingual performance and stronger relative transfer efficiency are not always aligned.

## E.2.3 Detailed Results for Open-source Models

Tables 10 and 11 provide the full transfer, accuracy, and off-target results for the open-source models in the observational study. These tables complement Figure 1 in the main text by showing the benchmark-level and language-level variation behind the averaged results.

## E.3 Controlled Comparative Study

### E.3.1 Reward Hyperparameter Sensitivity

To examine the role of different reward components, we vary the weights on the accuracy, format, and language-consistency rewards. Table 12 reports the corresponding results. Performance is most sensitive to the accuracy reward and the format reward. Removing the language-consistency reward substantially increases the off-target rate, while removing the format reward leads to the largest overall degradation. We use this analysis as a sensitivity check on the reward design; it does not fully disentangle reasoning improvement from formatting or language compliance.

### E.3.2 Detailed Results for Initial Model Types

Table 13 reports the full results for the initial model type comparison in the controlled comparative study.

### E.3.3 Detailed Results for Model Family

Figure 9 provides the full benchmark-level results for the model family comparison between Qwen2.5-7B-Instruct and Llama3.1-8B-Instruct. In both families, English-only GRPO improves performance across languages, but the relative gains differ substantially. As discussed in the main text, this comparison should be interpreted as a controlled comparison under matched post-training settings rather than as a pure estimate of family effects, since tokenization, pretraining data, and architecture remain entangled.

**Cross-lingual Transfer within the Sino-Tibetan Family.** To complement the model family comparison, we also evaluate Qwen2.5-7B-Instruct and Meta-Llama-3.1-8B-Instruct on several Sino-Tibetan languages after English-only GRPO. Table 14 shows that Qwen retains higher absolute accuracy in Chinese, while Llama exhibits larger relative gains on lower-performing languages such as Tibetan and Myanmar. We report this result as an additional case study consistent with the main-text observation that stronger starting performance and larger relative gains do not always coincide.

### E.3.4 Detailed Results for Model Size

Table 15 reports the full results for the model size comparison between *Qwen2.5-1.5B-Instruct* and *Qwen2.5-7B-Instruct*. The smaller model shows larger gains on multilingual MATH500 and GPQA-Diamond, whereas the larger model attains stronger

Settings	Accuracy per language											Average	
	en	es	ru	de	fr	bn	th	sw	zh	ja	te	Acc	Off-target
<i>Qwen2.5-7B-Base</i>													
Qwen-Math Template	49.2	31.8	25.2	28.2	30.0	5.8	23.0	2.4	27.0	15.2	1.4	21.7	16.6
Qwen-Instruct Template	50.6	38.0	30.0	33.2	38.4	10.4	26.8	2.4	30.0	27.8	4.4	26.5	15.7
No Template	44.4	38.2	28.2	28.4	35.6	6.0	17.0	0.2	29.2	19.8	1.2	22.6	18.5
<i>Qwen2.5-Math-7B-Base</i>													
Qwen-Math Template	43.4	36.6	2.8	21.8	21.4	26.2	15.0	2.2	36.6	9.4	1.2	19.7	30.5
Qwen-Instruct Template	56.6	46.4	11.2	33.4	36.8	31.4	28.2	4.2	44.4	25.2	3.2	29.2	18.0
No Template	37.8	33.0	4.2	29.8	37.4	29.0	12.4	0.0	36.2	17.2	3.0	21.8	31.5

Table 8: **The Performance of Base Models with Different Template Settings.** Accuracy (%) and Off-target rate (%) across languages for different template settings on multilingual MATH500 benchmark.

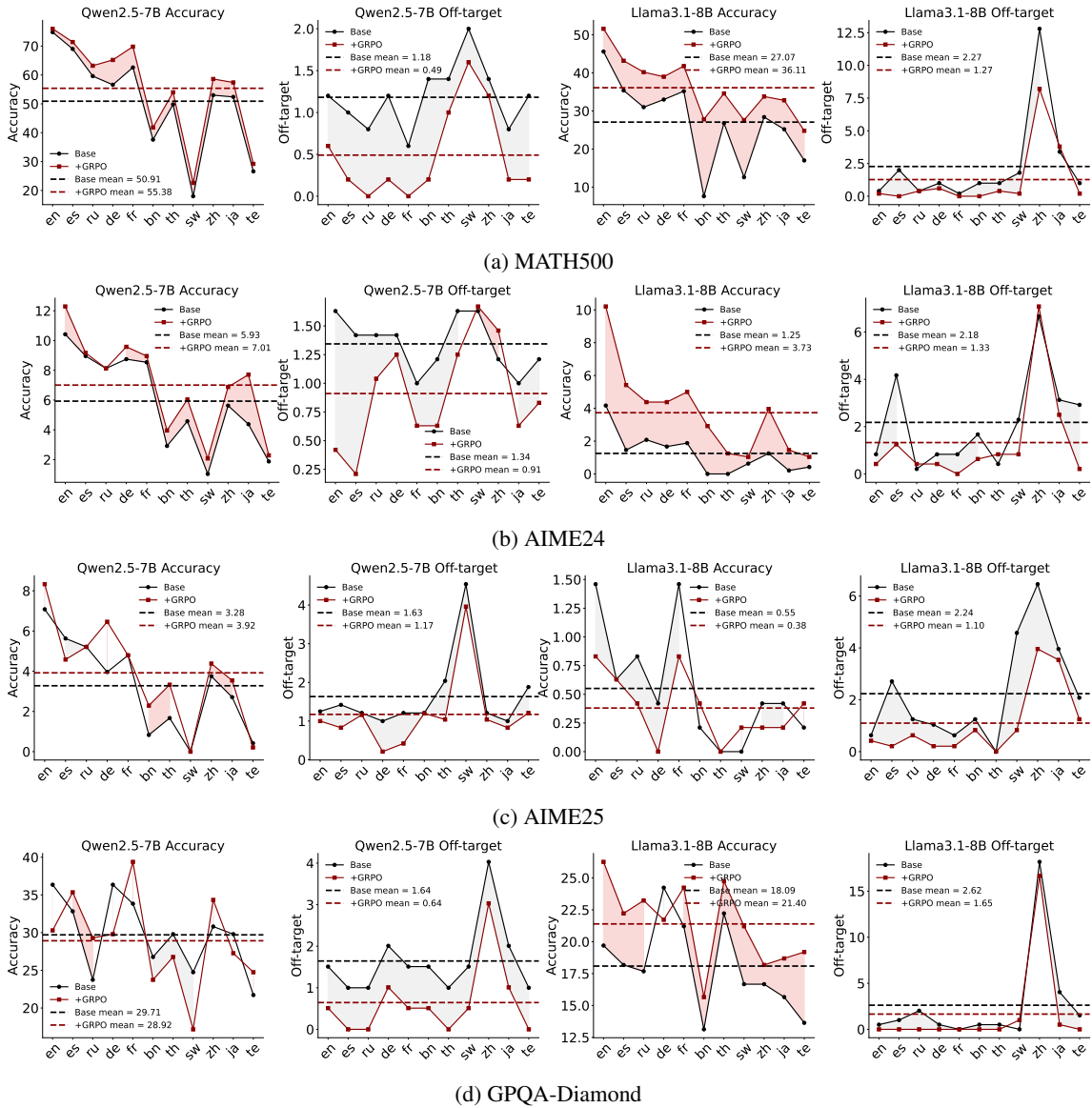


Figure 9: **The Impact of Different Model Families in Controlled Comparative Study.** Multilingual reasoning performance across languages, comparing the influence of model family using Qwen2.5-7B-Instruct and Llama3.1-8B-Instruct as initial models.

Settings	Accuracy per language											Average	
	en	es	ru	de	fr	bn	th	sw	zh	ja	te	Acc	Off-target
<i>Multilingual MATH500</i>													
Qwen2.5-1.5B	19.60	10.60	7.80	1.00	9.80	1.00	3.00	0.00	8.60	2.20	0.00	5.78	21.02
Qwen2.5-Math-7B	56.60	46.40	11.20	33.40	36.80	31.40	28.20	4.20	44.40	25.20	3.20	29.18	17.96
Qwen2.5-7B-Instruct	74.80	69.00	59.60	56.60	62.60	37.60	49.80	18.00	53.00	52.40	26.60	50.91	0.18
Qwen2.5-7B	50.60	38.00	30.00	33.20	38.40	10.40	26.80	2.40	30.00	27.80	4.40	26.55	15.69
Qwen2.5-14B	42.20	40.40	36.00	29.60	36.20	22.60	26.60	5.60	18.80	25.20	5.60	26.25	15.55
Qwen2.5-32B	54.00	50.80	42.00	37.80	46.80	24.60	33.00	13.60	28.00	42.60	11.60	34.98	5.56
Qwen2.5-32B-Instruct	78.60	73.60	68.00	68.40	69.40	53.20	60.60	37.40	61.40	65.00	43.40	61.73	0.13
<i>Multilingual AIME24</i>													
Qwen2.5-1.5B	0.21	0.42	0.00	0.00	0.00	0.00	0.21	0.00	0.63	0.00	0.00	0.13	19.26
Qwen2.5-Math-7B	13.75	6.46	1.67	3.33	4.79	2.71	3.33	0.00	6.88	1.67	0.63	4.11	21.76
Qwen2.5-7B-Instruct	10.42	8.96	8.13	8.75	8.54	2.92	4.58	1.04	5.63	4.38	1.88	5.93	0.34
Qwen2.5-7B	2.29	1.67	2.08	2.71	2.08	0.21	0.63	0.00	1.25	0.63	0.00	1.23	9.89
Qwen2.5-14B	2.50	2.50	2.29	2.50	2.71	0.63	0.21	0.00	1.04	0.83	0.21	1.40	16.02
Qwen2.5-32B	2.71	3.13	2.08	3.33	2.71	0.21	1.04	0.00	1.25	2.08	0.00	1.69	4.07
Qwen2.5-32B-Instruct	15.63	12.71	11.25	11.67	12.08	5.42	7.50	2.92	7.29	10.63	2.50	9.05	0.51
<i>Multilingual AIME25</i>													
Qwen2.5-1.5B	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.21	0.00	0.00	0.02	61.14
Qwen2.5-Math-7B	6.04	3.13	0.83	1.46	2.29	0.42	0.83	0.00	4.79	0.42	0.42	1.88	23.47
Qwen2.5-7B-Instruct	7.08	5.63	5.21	3.96	4.79	0.83	1.67	0.00	3.75	2.71	0.42	3.28	0.55
Qwen2.5-7B	0.83	0.83	0.21	1.25	0.63	0.21	0.00	0.00	0.42	0.21	0.00	0.42	10.38
Qwen2.5-14B	1.25	2.08	2.50	1.67	1.46	0.00	0.21	0.00	0.42	1.46	0.21	1.02	16.99
Qwen2.5-32B	1.04	1.04	1.46	0.21	0.83	0.21	0.00	0.00	1.04	1.04	0.21	0.64	4.02
Qwen2.5-32B-Instruct	11.25	7.29	7.29	6.04	6.67	1.25	2.71	0.00	5.42	2.71	0.42	4.64	0.30
<i>Multilingual GPQA-Diamond</i>													
Qwen2.5-1.5B	15.66	14.65	15.15	16.67	11.62	7.58	15.15	13.64	15.15	6.06	15.15	13.31	20.02
Qwen2.5-Math-7B	16.16	13.64	3.54	17.17	15.66	17.68	17.68	11.62	22.22	0.51	16.16	13.82	27.18
Qwen2.5-7B-Instruct	36.36	32.83	23.74	36.36	33.84	26.77	29.80	24.75	30.81	29.80	21.72	29.71	0.64
Qwen2.5-7B	28.79	24.24	20.71	22.22	18.18	11.11	21.72	17.68	23.23	17.17	12.63	19.79	9.69
Qwen2.5-14B	26.26	15.15	20.71	21.21	27.78	20.20	24.24	22.22	12.12	10.61	16.16	19.70	20.98
Qwen2.5-32B	28.28	30.30	28.28	33.33	23.74	15.66	24.75	17.68	27.27	27.78	17.17	24.93	9.00
Qwen2.5-32B-Instruct	45.45	41.92	38.89	41.41	44.44	29.80	38.38	32.83	36.36	38.38	27.27	37.74	0.28

Table 9: **The Performance of Initial Models.** Accuracy (%) and Off-target rate (%) across languages for different Initial models.

Models	Multilingual Reasoning Benchmarks				MTI		
	<i>MATH500</i>	<i>AIME24</i>	<i>AIME25</i>	<i>GPQA-D</i>	<i>ID</i>	<i>OOD</i>	<i>Avg</i>
DeepSeek-R1-Distill-Qwen-7B	3.493	2.312	2.864	4.168	3.493	3.115	3.209
Open-Reasoner-Zero-7B	3.195	2.677	1.479	1.320	3.195	1.825	2.168
OpenThinker2-7B	0.093	0.876	1.843	1.604	0.093	1.441	1.104
OpenThinker3-7B	0.157	1.502	2.434	1.318	0.157	1.752	1.353
Qwen-2.5-1.5B-SimpleRL-Zoo	5.322	2.173	0.856	1.383	5.322	1.383	3.353
Qwen-2.5-7B-SimpleRL-Zoo	4.543	3.189	1.217	6.531	4.543	3.646	3.870
Qwen-2.5-14B-SimpleRL-Zoo	2.381	3.360	0.959	1.655	2.381	1.991	2.089
Qwen-2.5-Math-7B-SimpleRL-Zoo	3.920	2.884	3.079	4.335	3.920	3.433	3.555
Qwen2.5-Math-7B-Dr.GRPO	2.807	1.324	3.158	2.149	2.807	2.210	2.359
s1.1-7B	0.310	0.920	1.192	0.671	0.310	0.928	0.773
DAPO-Qwen-32B	3.634	2.337	2.066	0.854	3.634	1.752	2.223
OpenThinker2-32B	0.936	1.513	4.235	0.201	0.936	1.983	1.721
S1.1-32B	1.382	1.583	3.429	0.821	1.382	1.944	1.804

Table 10: **The Performance of Various Open-source Models. Part 1:** Multilingual Transferability Index (MTI) of various models across benchmarks. The columns ID, OOD, and Avg refer to the MTI on in-domain (MATH500), out-of-domain (AIME24, AIME25, GPQA-Diamond), and all tasks, respectively.

Settings	Accuracy per language										Average		
	en	es	ru	de	fr	bn	th	sw	zh	ja	te	Acc	Off-target
<i>Multilingual MATH500</i>													
DeepSeek-R1-Distill-Qwen-7B	86.20	76.20	67.00	66.40	69.80	40.20	53.80	18.40	70.00	53.20	17.60	56.25	2.69
Open-Reasoner-Zero-7B	81.60	76.00	70.40	69.80	71.20	45.20	63.20	15.00	62.80	61.80	17.60	57.69	5.20
OpenThinker2-7B	86.00	74.80	64.80	63.00	73.00	34.40	58.80	12.60	65.80	70.00	8.40	55.60	36.13
OpenThinker3-7B	85.80	81.80	75.00	69.20	76.40	17.00	48.80	17.40	69.60	65.00	10.40	56.04	60.75
Qwen-2.5-1.5B-SimpleRL-Zoo	57.60	41.40	36.80	38.20	42.40	11.80	28.80	14.60	33.60	31.00	7.40	31.24	21.51
Qwen-2.5-7B-SimpleRL-Zoo	77.60	72.00	65.00	64.20	68.00	42.20	60.40	24.20	62.00	59.80	25.80	56.47	4.31
Qwen-2.5-14B-SimpleRL-Zoo	82.40	74.40	71.00	69.40	73.60	54.80	69.00	33.20	65.60	68.80	41.00	63.93	0.47
Qwen-2.5-Math-7B-SimpleRL-Zoo	80.40	72.60	66.00	68.60	70.60	45.80	57.40	15.00	61.80	54.40	14.20	55.16	8.33
Qwen2.5-Math-7B-Oat-Zero	79.80	72.40	32.80	55.60	50.40	47.60	49.40	16.80	58.00	43.40	11.80	47.09	15.11
s1.1-7B	75.80	68.20	60.80	59.00	69.60	37.00	57.40	16.40	57.20	51.60	20.40	52.13	12.76
DAPO-Qwen-32B	68.80	65.00	58.80	60.20	63.00	52.80	58.40	44.80	54.20	56.80	44.80	57.05	11.85
OpenThinker2-32B	96.00	88.60	85.20	84.20	85.00	73.00	80.60	35.40	77.40	75.60	47.20	75.29	13.02
S1.1-32B	95.40	91.20	85.00	83.60	88.00	73.00	81.20	57.00	77.40	80.80	53.60	78.75	27.40
<i>Multilingual AIME24</i>													
DeepSeek-R1-Distill-Qwen-7B	40.63	27.71	26.25	23.13	27.50	6.46	9.79	2.50	30.42	9.79	0.83	18.64	7.77
Open-Reasoner-Zero-7B	16.25	18.13	17.29	15.21	17.71	9.58	14.79	1.67	14.79	14.79	1.25	12.86	10.78
OpenThinker2-7B	37.08	18.33	17.08	13.75	20.83	11.04	20.21	2.50	25.63	27.29	5.42	18.11	39.72
OpenThinker3-7B	26.25	32.08	23.54	26.46	29.38	5.63	16.25	3.54	26.46	19.38	3.54	19.32	63.28
Qwen-2.5-1.5B-SimpleRL-Zoo	0.00	0.00	0.42	0.00	0.21	0.00	0.21	0.42	1.25	0.21	0.00	0.25	60.42
Qwen-2.5-7B-SimpleRL-Zoo	6.25	7.29	6.25	7.29	8.33	3.75	5.42	2.08	5.42	4.38	2.50	5.36	77.16
Qwen-2.5-14B-SimpleRL-Zoo	12.71	13.13	13.13	10.42	13.33	9.17	9.79	3.54	10.42	11.46	5.63	10.25	0.42
Qwen-2.5-Math-7B-SimpleRL-Zoo	25.83	15.42	13.96	11.25	13.33	5.00	8.13	1.67	10.21	8.54	2.50	10.53	11.29
Qwen2.5-Math-7B-Oat-Zero	28.33	12.92	5.42	8.54	11.67	6.25	8.75	1.04	12.29	6.67	0.42	9.30	20.57
s1.1-7B	14.38	10.42	10.42	10.21	11.46	5.21	7.92	1.67	8.75	7.71	0.21	8.03	7.95
DAPO-Qwen-32B	54.58	50.00	51.67	46.04	50.00	42.50	36.04	19.17	40.83	45.42	27.29	42.14	5.91
OpenThinker2-32B	74.17	61.88	55.42	56.67	55.42	59.17	49.17	13.96	56.88	37.71	34.58	50.45	22.08
S1.1-32B	58.75	55.21	49.17	51.25	53.33	36.04	41.25	19.38	44.58	46.88	17.08	42.99	7.80
<i>Multilingual AIME25</i>													
DeepSeek-R1-Distill-Qwen-7B	29.58	20.21	21.25	22.29	19.79	5.63	8.75	0.42	26.67	10.00	0.00	14.96	6.97
Open-Reasoner-Zero-7B	14.58	13.33	11.88	9.58	11.04	1.67	9.38	0.00	10.21	9.79	0.21	8.33	10.04
OpenThinker2-7B	28.33	21.67	20.63	17.08	21.46	9.38	17.08	2.71	25.00	24.38	2.08	17.25	39.77
OpenThinker3-7B	22.50	27.71	23.33	20.00	27.50	6.67	14.79	3.13	28.54	21.67	1.67	17.95	63.28
Qwen-2.5-1.5B-SimpleRL-Zoo	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.21	0.00	0.00	0.02	61.14
Qwen-2.5-7B-SimpleRL-Zoo	4.58	3.33	2.08	3.96	5.42	0.83	3.13	0.63	1.46	2.50	0.21	2.56	79.51
Qwen-2.5-14B-SimpleRL-Zoo	13.96	11.67	10.42	10.83	10.00	3.54	6.46	1.88	6.67	8.54	2.08	7.82	0.61
Qwen-2.5-Math-7B-SimpleRL-Zoo	13.75	9.58	6.04	5.42	9.79	2.71	5.63	1.04	6.25	3.75	1.04	5.91	12.12
Qwen2.5-Math-7B-Oat-Zero	10.00	9.38	2.29	6.67	4.38	1.46	2.08	1.04	6.67	2.92	0.42	4.30	22.12
s1.1-7B	13.96	11.67	9.58	6.88	11.88	2.08	7.08	0.21	9.79	5.21	0.00	7.12	6.17
DAPO-Qwen-32B	38.13	38.54	37.29	36.25	34.58	30.83	32.71	18.33	31.67	34.17	22.29	32.25	4.56
OpenThinker2-32B	57.29	50.00	48.13	52.29	43.96	45.42	41.88	12.50	52.50	36.04	25.63	42.33	22.65
S1.1-32B	50.00	43.54	38.33	43.33	42.71	29.38	31.88	16.04	38.75	35.42	14.58	34.91	8.05
<i>Multilingual GPQA-Diamond</i>													
DeepSeek-R1-Distill-Qwen-7B	32.32	33.33	33.33	35.35	35.35	18.18	21.21	23.74	29.80	14.65	14.14	26.49	6.11
Open-Reasoner-Zero-7B	37.37	26.77	31.82	33.33	33.33	24.24	32.83	12.63	33.33	26.77	7.07	27.23	6.20
OpenThinker2-7B	28.79	17.68	17.17	16.67	22.73	22.22	25.76	14.14	22.22	18.69	14.14	20.02	38.15
OpenThinker3-7B	23.23	18.69	22.22	16.67	24.24	5.05	14.65	12.63	21.72	10.10	7.07	16.02	59.23
Qwen-2.5-1.5B-SimpleRL-Zoo	20.71	15.66	23.74	16.67	24.75	10.10	18.18	13.13	17.17	21.21	8.59	17.26	14.69
Qwen-2.5-7B-SimpleRL-Zoo	30.30	31.82	29.80	31.82	33.84	20.71	23.74	17.17	29.80	23.74	10.10	25.71	3.49
Qwen-2.5-14B-SimpleRL-Zoo	41.92	40.40	34.85	40.91	39.39	27.78	34.85	29.80	39.39	33.33	26.26	35.35	2.62
Qwen-2.5-Math-7B-SimpleRL-Zoo	30.81	26.26	22.73	27.78	28.28	18.18	17.17	9.09	27.27	16.67	8.08	21.12	12.26
Qwen2.5-Math-7B-Oat-Zero	25.76	17.17	7.07	21.21	15.66	19.19	21.72	9.09	30.81	6.06	12.63	16.94	19.74
s1.1-7B	17.68	14.14	20.20	22.22	29.29	9.09	17.17	16.16	24.75	16.67	18.69	18.73	11.98
DAPO-Qwen-32B	52.50	44.44	40.91	48.99	41.92	37.37	42.93	31.82	46.97	47.98	30.81	42.42	5.88
OpenThinker2-32B	62.63	57.58	58.08	59.09	58.59	50.51	47.47	21.72	56.57	0.00	0.00	42.93	22.91
S1.1-32B	64.65	57.58	57.58	59.60	56.57	41.41	48.48	36.36	56.57	53.03	32.83	51.33	11.85

Table 11: **The Performance of Various Open-source Models. Part 3:** Accuracy (%) and Off-target rate (%) across languages for various open-source models.

ID	Reward Weights			Purpose	Accuracy (Acc) per Language										Avg. Acc	Avg. Off-target	MTI	
	$\lambda_1$	$\lambda_2$	$\lambda_3$		en	es	ru	de	fr	bn	th	sw	zh	ja				te
C1	0.8	0.1	0.1	Balanced	78.4	73.4	66.0	65.6	67.2	48.8	57.4	26.2	57.8	62.0	33.8	<b>57.87</b>	<b>0.20</b>	<b>2.50</b>
C2	0.9	0.1	0.0	w/o LCR	77.8	71.4	66.2	63.4	68.0	44.2	56.4	26.2	59.6	57.8	31.8	56.62	1.24	2.10
C3	0.7	0.1	0.2	Strong LCR	77.2	70.2	65.8	65.8	67.2	43.6	55.4	25.6	60.2	56.2	32.6	56.35	0.38	2.25
C4	0.9	0.0	0.1	w/o FR	75.2	60.8	62.2	56.0	60.2	36.2	43.8	18.0	54.4	31.8	24.8	47.58	0.31	-3.45
C5	0.7	0.2	0.1	Strong FR	78.2	71.6	65.2	62.6	67.8	45.0	60.0	24.8	60.0	60.8	29.2	56.84	0.35	2.23

Table 12: Sensitivity analysis of reward hyperparameters  $\lambda$ , trained on English (En) with parallel Russian (Ru) data. **Acc**, **Off**, and **MTI** denote Accuracy, Off-target rate, and Multilingual Transfer Index, respectively.

Model	Average accuracy across all languages				Avg	Off-target	MTI
	MATH500	AIME24	AIME25	GPQA			
Qwen2.5-7B-Base	26.55	1.23	0.42	19.79	12.00	11.41	-
↗ GRPO on En Data	52.16	7.10	3.35	27.18	22.45	3.12	1.95
Qwen2.5-7B-Instruct	50.91	5.93	3.28	29.71	22.45	1.43	-
↗ GRPO on En Data	54.24	7.41	3.92	28.47	23.51	0.94	1.23
Qwen2.5-Math-7B	29.18	4.11	1.88	13.82	12.25	22.59	-
↗ GRPO on En Data	45.73	8.84	3.96	18.96	19.37	9.50	2.12

Table 13: **The Impact of Initial Model Type on Interventional Study.** Accuracy (%), Off-target rate (%) and MTI across different initial model types.

Model	Tibetan (bo)	Myanmar (my)	S. Chinese (zh)
Qwen2.5-7B-Instruct	10.40	9.20	53.00
+ GRPO (En)	13.60 (+3.20)	11.60 (+2.40)	58.60 (+5.60)
Llama-3.1-8B-Instruct	0.80	2.40	28.40
+ GRPO (En)	15.00 (+14.20)	8.40 (+6.00)	33.80 (+5.40)

Table 14: Cross-lingual transfer performance (Accuracy %) on Sino-Tibetan languages before and after English-only RPT (GRPO).

final performance on the harder AIME benchmarks. These results support the main-text observation that the effect of model scale depends on the interaction between initial model strength and benchmark difficulty.

## E.4 Parallel Training Results

### E.4.1 Parallel Training Language Settings

Table 16 lists the language configurations used in the parallel training study, where the number of training languages increases from 1 to 7.

### E.4.2 Detailed Results of Parallel Training

Table 17 reports the full accuracy results across languages for different numbers of training languages. These tables complement Figure 4 in the main text by providing the full language-level breakdown behind the averaged trends.

### E.4.3 Effect of the Selected Parallel Language

Figures 10 report benchmark-level results for different choices of the added parallel language. Across settings, the same qualitative pattern remains: adding one parallel language consistently improves transfer, and the differences across language

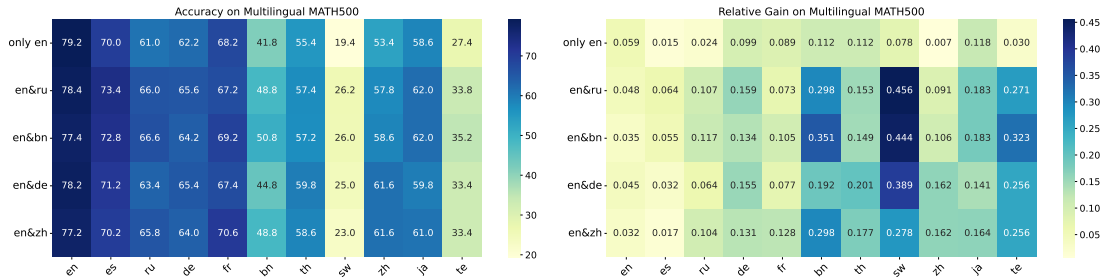
choices are modest relative to the monolingual-to-bilingual jump. Some language-specific effects remain, especially in lower-resource languages, but the main bilingual gain is robust to the choice of the added language.

## E.5 Parallel Training from a Chinese-Centric Perspective

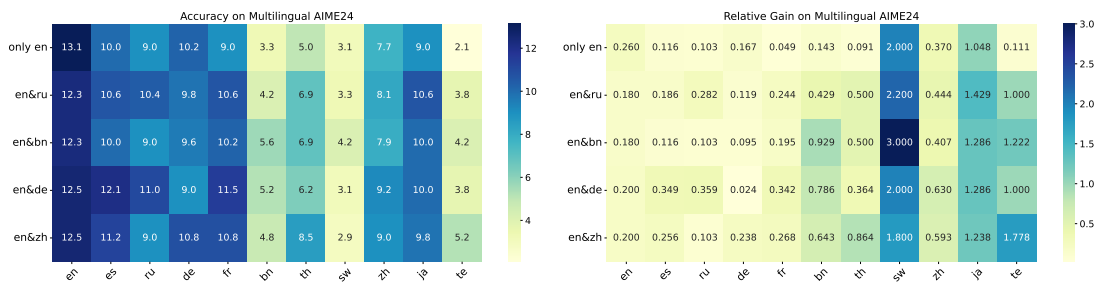
Table 18 reports the full results for three settings: monolingual Zh training, bilingual Zh+Ru training, and trilingual Zh+Ru+Fr training. A similar qualitative pattern appears: adding one parallel language yields the largest improvement in average accuracy and MTI, while adding a second language provides a smaller additional gain. We report this result as supporting evidence that the bilingual gain is not unique to English-only post-training.

## F Dataset License

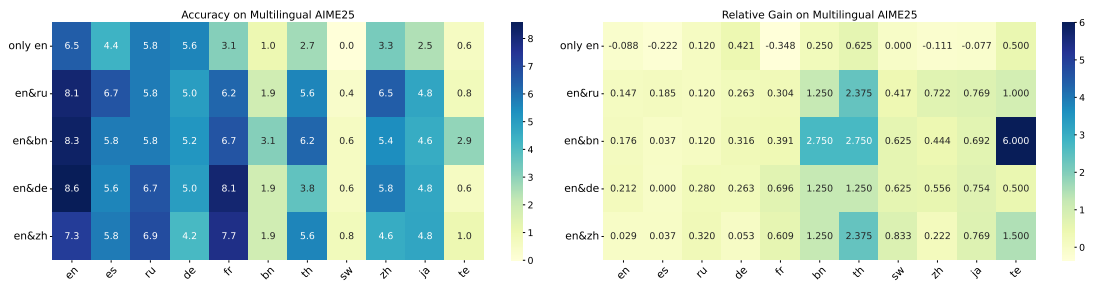
We use publicly available models and benchmarks, including the multilingual versions of MATH500 (Hendrycks et al., 2021), AIME2024 (Maxwell, 2024), AIME2025 (Kaggle, 2025) from XReasoning (Qi et al., 2025), and GPQA-Diamond (Rein et al., 2024) from BenchMAX (Huang et al., 2025b), and follow their corresponding usage guidelines.



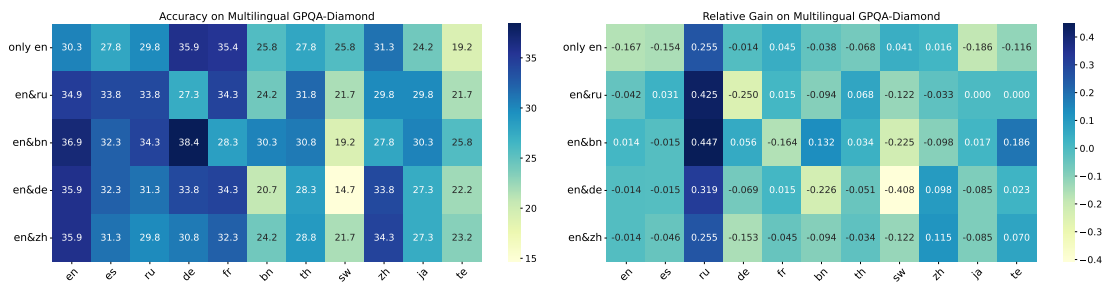
(a) MATH500



(b) AIME24



(c) AIME25



(d) GPQA-Diamond

Figure 10: **Analysis across Selected Parallel Languages.** The accuracy and relative gain across various benchmarks with different parallel languages. “Only en” denotes only fine-tuned on English data. “en&LANGUAGE” indicates the model was fine-tuned on English and a parallel language, with LANGUAGE representing *ru*, *bn*, *de*, *zh*, respectively.

Settings	$\Delta$ Performance											Average across languages	
	en	es	ru	de	fr	bn	th	sw	zh	ja	te	Training	Untraining
<i>Qwen2.5-1.5B-Instruct with GRPO on En Data</i>													
MATH500	20.40	27.20	17.40	14.40	17.40	6.20	7.80	4.40	16.80	14.20	7.40	20.40	13.32
AIME24	2.08	0.42	-0.21	1.25	0.42	0.21	0.21	0.00	-0.21	0.21	0.21	2.08	0.25
AIME25	1.04	0.21	0.00	0.21	0.00	0.00	0.00	-0.21	0.42	0.21	0.21	1.04	0.10
GPQA-Diamond	9.09	20.71	-0.51	5.05	12.12	1.52	-2.53	2.02	8.59	3.54	0.51	9.09	5.10
<i>Qwen2.5-7B-Instruct with GRPO on En Data</i>													
MATH500	4.40	1.00	1.40	5.60	5.60	4.20	5.60	1.40	0.40	6.20	0.80	4.40	3.22
AIME24	2.71	1.04	0.83	1.46	0.42	0.42	0.42	2.08	2.08	4.58	0.21	2.71	1.35
AIME25	1.25	-1.04	0.00	2.50	0.00	1.46	1.67	0.00	0.63	0.83	-0.21	1.25	0.58
GPQA-Diamond	-3.54	4.55	0.00	-1.01	7.07	-2.02	-1.01	-7.07	4.04	-3.03	0.51	-3.54	0.20

Table 15: **The Impact of Model Size in Interventional Study.**  $\Delta$  Performance on various benchmarks across *Qwen2.5-1.5B-Instruct* and *Qwen2.5-7B-Instruct*.

Settings	Training Parallel Languages
Only English	<i>en</i>
w. One parallel	<i>en, ru</i>
w. Two parallel	<i>en, ru, fr</i>
w. Three parallel	<i>en, ru, fr, es</i>
w. Four parallel	<i>en, ru, fr, es, de</i>
w. Five parallel	<i>en, ru, fr, es, de, bn</i>
w. Six parallel	<i>en, ru, fr, es, de, bn, th</i>
w. Seven parallel	<i>en, ru, fr, es, de, bn, th, zh</i>

Table 16: **The Language Settings in Parallel Scaling Law.**

Settings	Accuracy per language											Average		
	en	es	ru	de	fr	bn	th	sw	zh	ja	te	Acc	Off-target	MTI
<i>Multilingual MATH500</i>														
Only English	79.2	70.0	61.0	62.2	68.2	41.8	55.4	19.4	53.4	58.6	27.4	54.2	0.5	1.163
w. One parallel	78.4	73.4	66.0	65.6	67.2	48.8	57.4	26.2	57.8	62.0	33.8	57.9	0.2	2.496
w. Two parallel	79.0	73.4	64.4	67.4	69.2	45.2	60.2	26.2	63.0	61.6	32.6	58.4	0.2	2.650
w. Three parallel	77.8	73.6	64.6	68.4	69.8	46.0	60.8	24.0	62.2	60.8	34.2	58.4	0.4	3.002
w. Four parallel	77.2	71.2	66.2	66.8	68.0	47.6	61.8	28.6	62.0	60.2	35.2	58.6	0.6	3.282
w. Five parallel	77.4	71.4	62.2	66.2	66.0	48.6	62.0	32.4	62.2	63.8	37.0	59.0	0.4	3.475
w. Six parallel	76.4	70.8	63.8	65.6	66.8	48.6	61.8	34.6	63.4	63.4	38.4	59.4	0.5	3.534
w. Seven parallel	76.6	71.2	63.6	66.2	66.2	49.4	62.6	33.8	63.5	63.4	38.2	59.5	0.2	3.631

Table 17: **The Detailed Results in Parallel Scaling Law.**

Training Data	en	es	ru	de	fr	bn	th	sw	zh	ja	te	Avg. Acc	Avg. Off-target	MTI
Qwen2.5-7B-Instruct	74.8	69.0	59.6	56.6	62.6	37.6	49.8	18.0	53.0	52.4	26.6	50.91	1.18	-
Zh	76.8	70.4	61.8	62.2	66.8	43.0	55.2	21.8	59.4	56.6	31.0	55.00	0.38	0.86
Zh + Ru	78.8	72.0	64.4	65.0	70.0	44.4	60.0	28.6	61.4	60.4	33.4	58.13	0.26	1.69
Zh + Ru + Fr	78.0	72.2	63.6	66.4	69.8	45.4	61.0	30.0	62.8	64.0	33.2	<b>58.76</b>	<b>0.23</b>	<b>1.87</b>

Table 18: Parallel training experiments with Chinese (Zh) as the source language. Chinese-centric parallel training results on multilingual MATH500. A similar qualitative pattern appears: adding one parallel language yields the largest gain, while adding a second language provides a smaller additional improvement.

## G Prompts Template

1192

### G.1 Templates for Base Model

1193

#### Qwen-Instruct Template:

```
<|im_start|>system\n
You are Qwen, created by Alibaba Cloud. You are a helpful assistant.
<|im_end|>\n
<|im_start|>user\n{instruction}<|im_end|>\n
<|im_start|>assistant\n
```

#### Qwen-Math Template:

```
<|im_start|>system\n
Please reason step by step, and put your final answer within \boxed{ }.
<|im_end|>\n
<|im_start|>user\n{instruction}<|im_end|>\n
<|im_start|>assistant\n
```

#### No Template:

```
{instruction}
```

1194

### G.2 Multilingual Reasoning Instruction

1195

#### The Instruction Used in Multilingual Reasoning Prompt

Please always think in [LANGUAGE].

Solve the following mathematics problem step by step. At the end, provide your final answer enclosed in \boxed{ }.

Problem: {}

1196

### G.3 Prompt hacking to force response language

1197

#### The Prefixes Used in Prompt Hacking. Note that we list seven out of eleven languages.

- **English:** By request, I will start thinking in English.
- **Japanese:** 要求があれば、日本語で考え始めます。
- **Chinese:** 应要求，我将开始用中文思考。
- **Spanish:** A petición, empezaré a pensar en español.
- **French:** Sur demande, je commencerai à penser en français.
- **German:** Auf Anfrage werde ich anfangen, in Deutsch zu denken.
- **Swahili:** Kwa ombi, nitaanza kufikiria kwa Kiswahili.

1198

### G.4 Template for R1-like Reasoning

1199

#### The Template for R1-like Reasoning

You are a helpful AI Assistant that provides well-reasoned and detailed responses. You first think about the reasoning process as an internal monologue and then provide the user with the answer. The final answer must be put in \boxed{ }. Respond in the following format: <think>\n...\n</think>\n<answer>\n...\n</answer>

1200