

Unlocking Multilingual Math Potential: Strategies Amidst Data Scarcity

Anonymous ACL submission

Abstract

Existing works studying behaviors of Large Language Models (LLMs) in multilingual settings focus mainly on general downstream tasks such as instruction following. To fill this gap, we perform an in-depth analysis of LLMs' math reasoning capacity under multilingual settings and propose to alleviate the shortage of high quality multilingual math reasoning post-training data by exploring whether prior English math knowledge and additional English data helps, and by observing the effects of multilingual synthetic data on performance. For models pre-trained with mostly English data, we find that prior English math knowledge helps and that scaling English data helps only when the training and evaluation data belong to similar distributions (human/machine translated). Additionally, we find that inclusion of multilingual synthetic data leads to improved performance on human-translated data but degraded performance on machine-translated data. Our findings shed light on effective finetuning of LLMs for better multilingual math reasoning performance given the shortage high-quality multilingual math reasoning data.

1 Introduction

In the era of modern Large Language Models (LLMs), the rising interest in multilingual capabilities of models has led to the extension of existing datasets into different languages via human translation (Shi et al., 2023), machine translation services or LLM-based translation (Chen et al., 2023b; Li et al., 2023; Lai et al., 2023; Chai et al., 2024; Lai and Nissim, 2024). Detailed analyses for behaviors of models on translated datasets exists for tasks like multilingual instruction following (Shaham et al., 2024), but are missing for multilingual math reasoning, which is an important task to measure the reasoning capabilities of LLMs. Moreover, the low availability of multilingual math training data, as compared to English, hinders improvements as

well. In this work, we study model behaviors and explore various strategies to improve multilingual math reasoning performance given the scarcity of high quality multilingual data.

We first study the variance of cross-lingual performance for a wide pool of LLMs when post-trained for a certain language on math reasoning task. We observe that finetuning multilingual LLMs (which are already pre-trained with target task multilingual data) with additional math reasoning data belonging to one language does not lead to highly varying performance difference between the training language and other languages, indicating that post-training in this scenario leads to saturating results. This is not the case for monolingual models (pretrained mostly with English data), which exhibit high performance variance across languages.

Despite their variance, monolingual models are suitable to study the impact of different training and data setups for multilingual tasks. We observe that monolingual models having math reasoning capabilities in English transfer better to multilingual setting compared to the base model, implying that intermediate training with additional English math reasoning data helps in overall performance. Further, we find that scaling English training data while keeping non-English data constant benefits performance on evaluation sets which align with the training data in-terms of dataset construction (human v/s machine-translated).

The above disparity between data distribution of evaluations datasets and the gaining popularity of synthetic data usage for post-training (Huang et al., 2023; Sun et al., 2024; Ri et al., 2024; Dubey et al., 2024), prompt us to examine the effect of introduction of multilingual synthetic data on available benchmarks. We find that scaling synthetic multilingual data obtained via self-training helps the model to consistently perform better on human-translated test set while leading to diminishing per-

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082

083	formance on machine-translated test set. This result	133
084	highlights the importance of determining the opti-	134
085	mal mix of multilingual datasets (which are domi-	
086	nantly machine-translated) with synthetic data for	
087	obtaining maximal model performance.	
088	Overall, our experiments provide an in-depth un-	136
089	derstanding of the status of multilingual math rea-	137
090	soning, along with insights that can be leveraged to	138
091	train better models. More specifically, models with	139
092	prior English math knowledge and determining the	140
093	appropriate proportion of English and multilingual	141
094	synthetic data in the training data play are crucial	142
095	for multilingual math reasoning.	143
		144
		145
096	2 Related Work	146
		147
097	Multilingual Math Reasoning Measuring the	148
098	capability of models to solve mathematical prob-	
099	lems is directly correlated with their complex rea-	
100	soning proficiency. Due to the lack of abundant	
101	multilingual domain data, most works rely on ex-	
102	tending datasets using translation models or fron-	
103	tier LLMs (Chen et al., 2023a; Shi et al., 2023; Lai	
104	and Nissim, 2024). Works have also focused on	
105	leveraging auxiliary tasks to better align models for	
106	multilingual data (Zhu et al., 2024a,b). Another	
107	line of works emphasize on studying the reason-	
108	ing consistency of LLMs in different languages	
109	Shi et al. (2023); Chowdhery et al. (2023); Lai and	
110	Nissim (2024). However, none of the works study	
111	the differences between the behaviors of monolin-	
112	gual and multilingual models when finetuned for	
113	multilingual math.	
114	Influence of data attributes on multilingual per-	
115	formance Cross-lingual transfer is a well known	
116	phenomena, which has been shown to be effective	
117	in cases with limited multilingual data (Shaham	
118	et al., 2024; Chen et al., 2023b). Prior works have	
119	shown that multilingual performance benefits from	
120	intermediate multitask training in English (Phang	
121	et al., 2020) and from an imbalance proportion of	
122	languages in training data (Schäfer et al., 2024).	
123	While these works are aimed at studying the im-	
124	pact of data for multilingual tasks in general, our	
125	work is focused on exploring the effects specifically	
126	in math settings.	
127	The predominant focus of related works has been	
128	disjointly on multilingual math reasoning and the	
129	impact of data attributes for multilingual tasks in	
130	general. In this paper, we address the gap of an	
131	exhaustive comparison of monolingual and multi-	
132	lingual models for multilingual math reasoning and	
	explore techniques to maximize performance by	133
	leveraging English and multilingual synthetic data.	134
	3 Task Setup	135
	The task is framed as a supervised finetuning task,	136
	where given a math problem, the objective is to	137
	generate the steps to solve the problem and output	138
	the final answer. In addition to the question, the	139
	model is also prompted with the Chain of Thought	140
	prompt template (in the language corresponding	141
	to the question), which has been corroborated to	142
	enhance performance (Wei et al., 2022). We use	143
	the mCoT dataset (Lai and Nissim, 2024) as the	144
	training data for all our experiments and use answer	145
	exact match (accuracy) as our evaluation metric on	146
	the MGSM (Shi et al., 2023) and MSVAMP (Chen	147
	et al., 2023a) datasets.	148
	4 Improvement strategies under data	149
	constraint settings	150
	4.1 Cross-lingual performance variance for	151
	math/reasoning when trained on a single	152
	language	153
	We first explore the cross-lingual performance vari-	154
	ance in LLMs when trained on a single language.	155
	We perform full finetuning of the models for only	156
	one epoch to avoid any possibility of overfitting.	157
	As shown in Figure 2, we observe that multi-	158
	lingual models like Aya23 (Aryabumi et al., 2024)	159
	and Qwen2 (Yang et al., 2024), when trained on dif-	160
	ferent languages exhibit similar performance (low	161
	variance) for any particular language, though the	162
	performance range varies depending on the model	163
	and the resource level of the language. This implies	164
	that the models rely on their existing knowledge de-	165
	spite being post-trained specifically for a particular	166
	language. We hypothesize that the additional train-	167
	ing with data belonging to a particular language	168
	does not influence the model much since it has al-	169
	ready been trained on vast amounts of multilingual	170
	data during the pre- and post- training stages.	171
	These observations are contrary for models that	172
	have a low proportion of multilingual data in their	173
	pre-training corpus like Mistral-7B (Jiang et al.,	174
	2023a) and Llama3-8B (Dubey et al., 2024) (Fig-	175
	ure 3). Such models exhibit considerable variance	176
	when the evaluation data language is constant due	177
	to additional learning signals provided to them.	178
	Next, we study the correlations between lan-	179
	guages. Specifically, we tabulate the languages (in	180
	Figure 4) which when used to train models lead to	181

the top-k performance in a particular language. We consider only Mistral-7B-v0.1 Base and Llama3 Base for computing the results since their monolingual nature provides the best setup to study the influence of languages on performance. We consider evaluation metrics on the MGSM and MSVAMP datasets to bolster the confidence of the observations and set $k = 4$. Unsurprisingly, we observe that the target language benefits from training data in the same language and that European languages (English, Spanish, French, and German) help each other.

Based on the above results, monolingual models serve the optimal setting to study the influence of different training strategies and dataset mixtures since multilingual models might not be affected by them.

4.2 Does additional English math/reasoning data help in multilingual setting?

LLMs post-trained for a specific task or data might locate the model in a local minima, which leads to poor adaptation and generalization for new tasks/data. We intend to explore if this is the case when LLMs with math reasoning capabilities in English are adapted to multilingual data. To make fair comparisons, we experiment with WizardMath (Luo et al., 2023), which is a Mistral-7B-v0.1 Base model, trained for English Math problems using Reinforcement Learning, and Mathstral¹, another variant of the Mistral base model. We follow training hyperparameters similar to Lai and Nissim (2024) and train the models for two epochs.

We observe that both WizardMath and Mathstral outperform Mistral-7B-v0.1 (Table 1) with improvements on almost all languages in the evaluation data (Figure 5). We hypothesize two reasons for this gain in performance namely, additional English training data² and training strategies other than supervised finetuning (Uesato et al., 2022; Luo et al., 2023; Lightman et al., 2023; Gao et al., 2024)

Continuing our previous discussion and given that English data is easily available for most tasks, we design experiments in a controlled setting to verify if additional English training data leads to improved performance. We sample 25% of non-English data from our training set and gradually increase the proportion of English data. Thus, we sample 25%, 50%, 75% and 100% of English data

¹<https://mistral.ai/news/mathstral/>

²The models have benchmarks on GSM 8K only, which suggests the model being suitable mainly for English

Finetuned Model	MGSM	MSVAMP
Mistral-7B-v0.1	0.68	0.668
WizardMath	0.707	0.696
Mathstral	0.748	0.752

Table 1: Average accuracy of English Math Reasoning models in multilingual setting. **Takeaway:** English math knowledge improves performance

and mix it with the sampled non-English data.

Figure 1 shows the results of training a Mistral-7B-v0.1 model on these 4 data mixtures on the MGSM and MSVAMP datasets. We observe that scaling English data mostly improves performance in English with slight effects on non-English data, though the overall performance variations are negligible on MGSM, while non-trivial on MSVAMP. We hypothesize that this is due to data belonging to different distributions. Scaling English data whose Google translations constitute the multilingual component of the training set helps in improving performance on MSVAMP (Google Translated) and does not affect results on MGSM (human translated).

Thus, training models with prior math knowledge is beneficial for multilingual math reasoning, though sole reliance on machine-translated data must be avoided as it might only help improvements on machine-translated benchmarks and lead to plateauing performance on human-translated ones.

4.3 Impact of synthetic multilingual data

We create another controlled setting at a small scale to explore the influence of synthetic multilingual data. We sample 5K samples belonging to each of the 11 languages from the mCoT training data. We call this dataset *mCoT-55K*. We follow the procedure devised by Wang et al. (2024) to incorporate synthetic data into our training pipeline but restrict it to one iteration of synthetic data generation. First, we perform supervised finetuning (SFT) of Mistral-7B-v0.1 on *mCoT-55K*, which is treated as the baseline for the experiment. Then we sample 5 generations for each sample from the SFT model at a temperature of 0.7. Using the ground truth answer, we construct pair-wise data to train the SFT model further using Direct Policy Optimization (Rafailov et al., 2024) (DPO). The DPO model is then used to sample positive samples, which are then included with the *mCoT-55K* samples (we use the same sampling settings as earlier). The SFT model is then trained on this combination of origi-

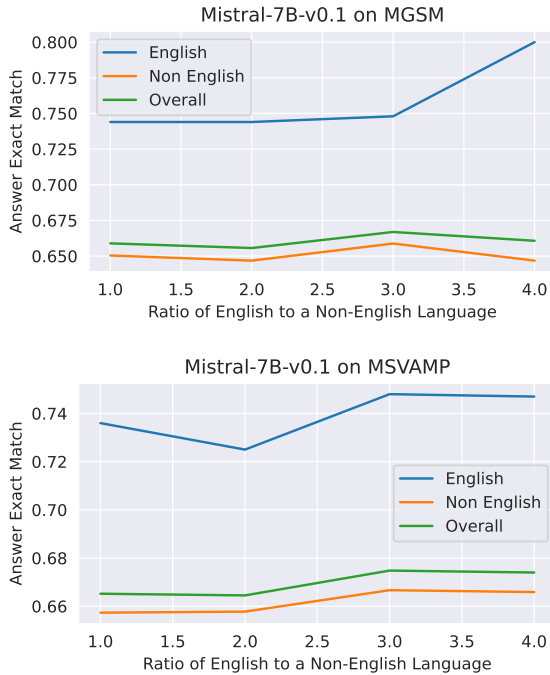


Figure 1: Mistral-7B-v0.1 on constant non-English samples and varying proportion of English samples. **Takeaway:** Scaling proportion of English helps when training and test data are created in a similar manner

nal training data and synthetic data to get the final model.

We ablate on the proportion of synthetic data used in the final training. More specifically, we gradually add 25%, 50%, 75% and 100% of the generated synthetic SFT data ensuring a balance in the samples per language (entire synthetic data has 3100 samples per language).

We find that the addition of synthetic data always achieves better performance than the baseline model on the human-translated MGSM test set, whereas the performance on the Google-translated MSVAMP is always poor compared to the baseline (Table 2). For MGSM, we do not find any particular trend as the synthetic data is scaled, which can be attributed to different synthetic samples having different quality levels. For MSVAMP, we find that the performance decreases as the proportion of synthetic data increases. We hypothesize that the addition of synthetic data adds diversity to the Google-Translated training dataset and shift the overall data distribution away from Google-Translated distribution, thus leading to degrading performance on MSVAMP and improved results on MGSM.

The above observations signify that the an appropriate mixture of machine-translated data and

multilingual synthetic data is mandatory to optimize performance on both human-translated and machine-translated evaluation benchmarks and additional work is needed to devise methods and provide recommendations for the same.

Dataset	MGSM	MSVAMP
<i>mCoT-55K</i>	0.469	0.615
+ 25% synthetic	0.532	0.596
+ 50% synthetic	0.515	0.577
+ 75% synthetic	0.517	0.553
+ 100% synthetic	0.497	0.525

Table 2: Results for multilingual synthetic data scaling. **Takeaway:** Synthetic data is similar to human-generated data since it leads to consistent improvements on human-translated test bed (MGSM), but degradation on Google-Translated test bed (MSVAMP)

5 Conclusion

We study various aspects revolving around multilingual math settings regarding the performance variance of monolingual and multilingual LLMs when trained for only one language, the transfer of English math reasoning models to multilingual data, and the impact of multilingual synthetic data. Given that the majority of multilingual task data is machine-translated, determining the optimal proportion of multilingual synthetic data becomes crucial to align models toward human generated data. Overall, our findings help us understand the behaviors and training dynamics of various models in different training data settings and thus help in designing better experimental setups for improving multilingual math reasoning performance.

6 Limitations

First, our study is limited to 11 languages, for which original English datasets were translated. Moreover, amongst the 11 languages, very few of them belong to the low-resource category. Second, the training dataset used is constructed entirely using machine-translation (Google Translate). Prior works have highlighted the difference between machine-translated data and human-translated data (Jiang et al., 2023b; Luo et al., 2024). Further, we experiment with only one technique of synthetic data generation and restrict our study to models with parameter count in the 7-8 B range due to computation constraints.

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
378
379

380
381
382
383
384

385
386
387
388

References

Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Kelly Marchisio, Sebastian Ruder, et al. 2024. Aya 23: Open weight releases to further multilingual progress. *arXiv preprint arXiv:2405.15032*.

Linzheng Chai, Jian Yang, Tao Sun, Hongcheng Guo, Jiaheng Liu, Bing Wang, Xiannian Liang, Jiaqi Bai, Tongliang Li, Qiyao Peng, et al. 2024. xcot: Cross-lingual instruction tuning for cross-lingual chain-of-thought reasoning. *arXiv preprint arXiv:2401.07037*.

Nuo Chen, Zinan Zheng, Ning Wu, Linjun Shou, Ming Gong, Yangqiu Song, Dongmei Zhang, and Jia Li. 2023a. Breaking language barriers in multilingual mathematical reasoning: Insights and observations. *arXiv preprint arXiv:2310.20246*.

Pinzhen Chen, Shaoxiong Ji, Nikolay Bogoychev, Barry Haddow, and Kenneth Heafield. 2023b. Monolingual or multilingual instruction tuning: Which makes a better alpaca. *arXiv preprint arXiv:2309.08958*.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.

Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Bofei Gao, Zefan Cai, Runxin Xu, Peiyi Wang, Ce Zheng, Runji Lin, Keming Lu, Junyang Lin, Chang Zhou, Tianyu Liu, et al. 2024. The reason behind good or bad: Towards a better mathematical verifier with natural language feedback. *arXiv preprint arXiv:2406.14024*.

Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2023. **Large language models can self-improve**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1051–1068, Singapore. Association for Computational Linguistics.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023a. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Zhaokun Jiang, Qianxi Lv, and Ziyin Zhang. 2023b. Distinguishing translations by human, nmt, and chatgpt: A linguistic and statistical approach. *arXiv preprint arXiv:2312.10750*.

Huiyuan Lai and Malvina Nissim. 2024. mcot: Multilingual instruction tuning for reasoning consistency in language models. *arXiv preprint arXiv:2406.02301*.

Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan A Rossi, and Thien Huu Nguyen. 2023. Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback. *arXiv preprint arXiv:2307.16039*.

Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. 2023. Bactrian-x: Multilingual replicable instruction-following models with low-rank adaptation. *arXiv preprint arXiv:2305.15011*.

Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*.

Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*.

Jiaming Luo, Colin Cherry, and George Foster. 2024. To diverge or not to diverge: A morphosyntactic perspective on machine translation vs human translation. *Transactions of the Association for Computational Linguistics*, 12:355–371.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.

Jason Phang, Iacer Calixto, Phu Mon Htut, Yada Pruk-sachatkun, Haokun Liu, Clara Vania, Katharina Kann, and Samuel R. Bowman. 2020. **English intermediate-task training improves zero-shot cross-lingual transfer too**. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 557–575, Suzhou, China. Association for Computational Linguistics.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

Ryokan Ri, Shun Kiyono, and Sho Takase. 2024. Self-translate-train: A simple but strong baseline for cross-lingual transfer of large language models. *arXiv preprint arXiv:2407.00454*.

444	Anton Schäfer, Shauli Ravfogel, Thomas Hofmann,	Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo,	498
445	Tiago Pimentel, and Imanol Schlag. 2024. Lan-	Chien-Chin Huang, Min Xu, Less Wright, Hamid	499
446	guage imbalance can boost cross-lingual generali-	Shojanazeri, Myle Ott, Sam Shleifer, et al. 2023. Py-	500
447	sation. <i>arXiv preprint arXiv:2404.07982</i> .	torch fsdp: experiences on scaling fully sharded data	501
		parallel. <i>arXiv preprint arXiv:2304.11277</i> .	502
448	Uri Shaham, Jonathan Herzig, Roei Aharoni, Idan	Wenhao Zhu, Shujian Huang, Fei Yuan, Cheng Chen,	503
449	Szpektor, Reut Tsarfaty, and Matan Eyal. 2024. Mul-	Jiajun Chen, and Alexandra Birch. 2024a. The power	504
450	tilingual instruction tuning with just a pinch of multi-	of question translation training in multilingual reason-	505
451	linguality. <i>arXiv preprint arXiv:2401.01854</i> .	ing: Broadened scope and deepened insights. <i>arXiv</i>	506
		<i>preprint arXiv:2405.01345</i> .	507
452	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang,	Wenhao Zhu, Shujian Huang, Fei Yuan, Shuaijie She,	508
453	Suraj Srivats, Soroush Vosoughi, Hyung Won Chung,	Jiajun Chen, and Alexandra Birch. 2024b. Question	509
454	Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das,	translation training for better multilingual reasoning.	510
455	and Jason Wei. 2023. Language models are multi-	<i>arXiv preprint arXiv:2401.07817</i> .	511
456	lingual chain-of-thought reasoners . In <i>The Eleventh</i>		
457	<i>International Conference on Learning Representa-</i>		
458	<i>tions</i> .		
459	Zhiqing Sun, Yikang Shen, Qinhong Zhou, Hongxin		
460	Zhang, Zhenfang Chen, David Cox, Yiming Yang,		
461	and Chuang Gan. 2024. Principle-driven self-		
462	alignment of language models from scratch with		
463	minimal human supervision. <i>Advances in Neural</i>		
464	<i>Information Processing Systems</i> , 36.		
465	Jonathan Uesato, Nate Kushman, Ramana Kumar, Fran-		
466	cis Song, Noah Siegel, Lisa Wang, Antonia Creswell,		
467	Geoffrey Irving, and Irina Higgins. 2022. Solv-		
468	ing math word problems with process-and outcome-		
469	based feedback. <i>arXiv preprint arXiv:2211.14275</i> .		
470	Tianduo Wang, Shichen Li, and Wei Lu. 2024. Self-		
471	training with direct preference optimization im-		
472	proves chain-of-thought reasoning. <i>arXiv preprint</i>		
473	<i>arXiv:2407.18248</i> .		
474	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten		
475	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,		
476	et al. 2022. Chain-of-thought prompting elicits rea-		
477	soning in large language models. <i>Advances in neural</i>		
478	<i>information processing systems</i> , 35:24824–24837.		
479	Andrea W Wen-Yi, Unso Eun Seo Jo, Lu Jia Lin, and		
480	David Mimno. 2024. How chinese are chinese lan-		
481	guage models? the puzzling lack of language policy		
482	in china’s llms. <i>arXiv preprint arXiv:2407.09652</i> .		
483	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien		
484	Chaumond, Clement Delangue, Anthony Moi, Pier-		
485	ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,		
486	et al. 2020. Transformers: State-of-the-art natural		
487	language processing. In <i>Proceedings of the 2020 con-</i>		
488	<i>ference on empirical methods in natural language</i>		
489	<i>processing: system demonstrations</i> , pages 38–45.		
490	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,		
491	Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan		
492	Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2		
493	technical report. <i>arXiv preprint arXiv:2407.10671</i> .		
494	Jun Zhao, Zhihao Zhang, Qi Zhang, Tao Gui, and Xu-		
495	anjing Huang. 2024. Llama beyond english: An em-		
496	pirical study on language capability transfer. <i>arXiv</i>		
497	<i>preprint arXiv:2401.01055</i> .		

A Experimental Settings

We train our models using a learning rate of $5e - 6$ and a batch size of 8 per GPU with 16 steps of gradient accumulation. Cosine learning rate scheduler is used with a linear warm up of 3% training steps. The model is trained for a maximum sequence length of 1024 and the number of training epochs is set to 1 and 2 for the language-wise and all language training. We use Pytorch (Paszke et al., 2019), Huggingface’s transformers (Wolf et al., 2020) and native Pytorch Fully Sharded Data Parallel (FSDP) (Zhao et al., 2023) on 8 x H100s for training our models.

B Comparison of models when trained on a single language

Figures 2, 3 and 4 show the results for Section 4.1. An interesting observation for Llama3-8B Base is that training on Chinese (high resource) and Thai (low resource) does not lead to the highest performance. Many works revolve around adapting Llama-based models on Chinese data (Zhao et al., 2024; Cui et al., 2023). On the other hand, Wen-Yi et al. (2024) show that the behavior of international LLMs is similar to Chinese LLMs. Thus, future work is needed to investigate if this trend exhibited by Llama3 holds for other domains and tasks.

C Language Wise Results for Mistral variants

Figure 5 shows the results for Section 4.2

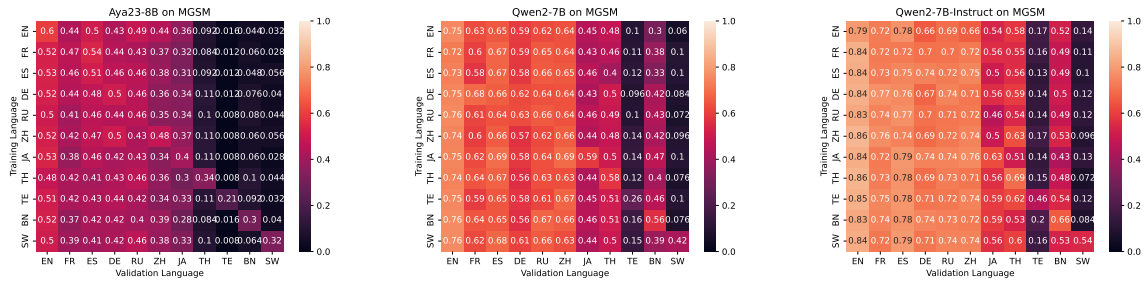


Figure 2: Language-wise performance matrix for multilingual models on MGSM. **Takeaway:** Performance on a fixed evaluation language is consistent for multilingual models despite a language specific post-training

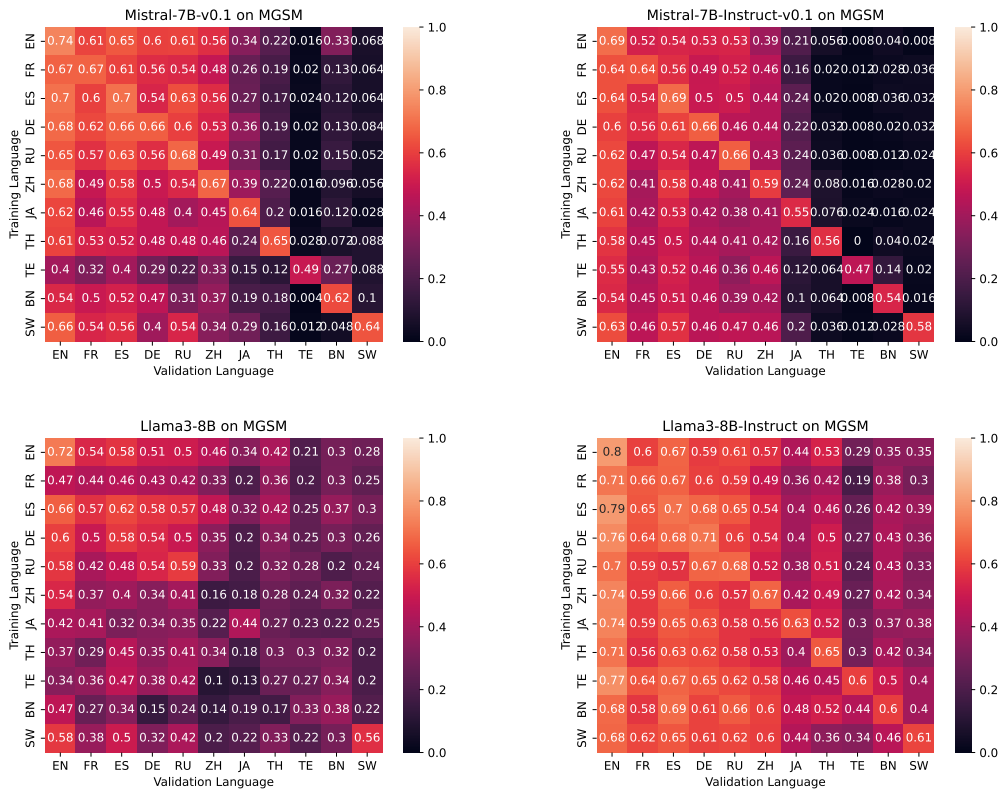


Figure 3: Language-wise performance matrix for monolingual models on MGSM. **Takeaway:** Performance on a fixed evaluation language varies significantly as compared to multilingual models

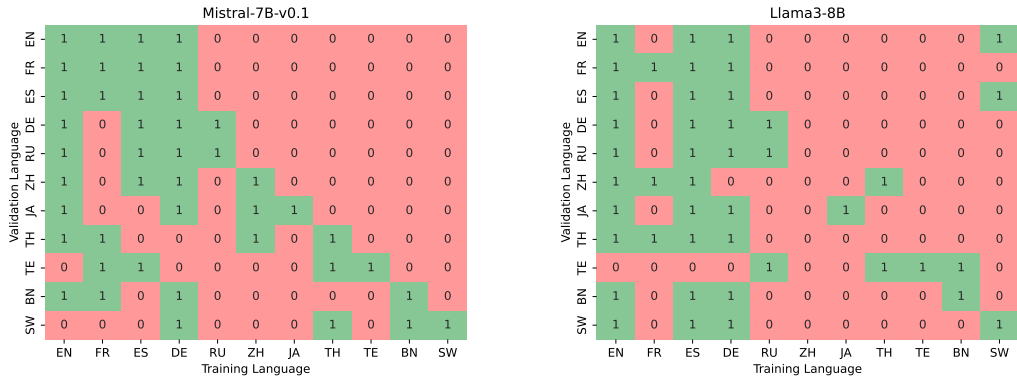


Figure 4: Training Languages which lead to highest performances while validating for a particular language. **Takeaway:** Same training and validation language work better in most cases; some exceptions are found for Llama3-8B Base

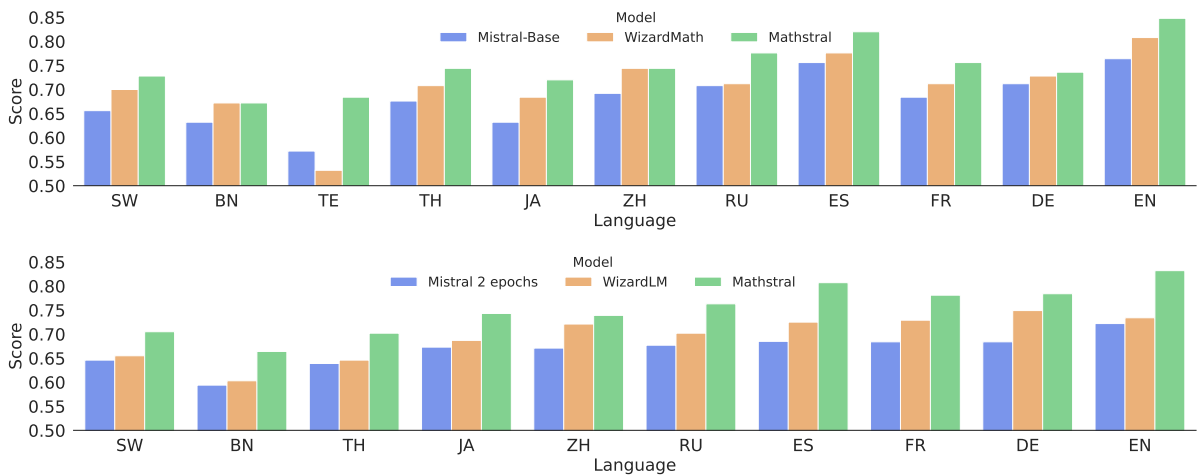


Figure 5: Evaluation of Mistral variants on MGSM (top) and MSVAMP (bottom) **Takeaway:** English math knowledge improves performance