

From Unaligned to Aligned: Scaling Multilingual LLMs with Multi-Way Parallel Corpora

Anonymous ACL submission

Abstract

Continued pretraining and instruction tuning on large-scale multilingual data have proven to be effective in scaling large language models (LLMs) to low-resource languages. However, the unaligned nature of such data limits its ability to effectively capture cross-lingual semantics. In contrast, *multi-way parallel data*, where identical content is aligned across multiple languages, provides stronger cross-lingual consistency and offers greater potential for improving multilingual performance. In this paper, we introduce a large-scale, high-quality multi-way parallel corpus, TED2025, based on TED Talks. The corpus spans 113 languages, with up to 50 languages aligned in parallel, ensuring extensive multilingual coverage. Using this dataset, we investigate best practices for leveraging multi-way parallel data to enhance LLMs, including strategies for continued pretraining, instruction tuning, and the analysis of key influencing factors. Experiments on six multilingual benchmarks show that models trained on multi-way parallel data consistently outperform those trained on unaligned multilingual data.¹

1 Introduction

Large language models (LLMs) have demonstrated remarkable performance on various tasks in high-resource languages (Huang et al., 2024; Qin et al., 2024). However, their performance still lags behind in low-resource languages (Huang et al., 2023; Lai et al., 2023a). To bridge this gap, recent efforts have focused on continued pretraining (Ji et al., 2024; Groeneveld et al., 2024) and instruction tuning (Lai et al., 2024; Üstün et al., 2024), utilizing large-scale unaligned multilingual text. However, these methods do not take full advantage of the explicit many-to-many alignments present in *multi-way parallel corpora*, which have been shown to

improve cross-lingual representations in multilingual NLP (Qi et al., 2018; Freitag and Firat, 2020; Xu et al., 2022; Wu et al., 2024). More recently, Mu et al. (2024) demonstrated that multi-way parallel inputs can also enhance in-context learning (Dong et al., 2024). However, scaling multilingual LLMs using multi-way parallel data remains underexplored.

Existing multi-way parallel datasets typically cover only a limited number of languages, domains and levels of parallelism (see Table 1). In contrast, TED Translators², a global community translating TED talk transcripts into over 100 languages, provides consistently high-quality, human-verified translations and serves as an ideal source for a large-scale multi-way parallel corpus. However, the largest TED-based datasets (Qi et al., 2018; Reimers and Gurevych, 2020) have not been updated since 2020, limiting their utility for LLM training and potentially exacerbating hallucinations (Ji et al., 2023).

To address these limitations, we introduce TED2025, a large-scale, high-quality multi-way parallel corpus derived from the latest TED talks. TED2025 covers 113 languages, 352 domain labels, and supports up to 50-way parallelism. Compared to existing resources, it offers more frequent updates and significantly broader coverage across both languages and domains, thereby strengthening the data foundation for multilingual LLM training.

Utilizing TED2025, we investigate three key research questions:

RQ1: How does fine-tuning on multi-way parallel data compare to training on unaligned multilingual text in terms of zero-shot cross-lingual transfer and representation alignment?

RQ2: Which strategies for selecting parallelism in multi-way parallel data (e.g., degree of parallelism and language subsets) lead to the greatest

¹We will make the code and data publicly available upon acceptance.

²<https://www.ted.com/participate/translate>

Dataset	Source	#Langs	#Domains	Max Parallelism	Type	Collection Method	Date Range	Open-source?
Bible (Christodoulopoulos and Steedman, 2015)	Bible	100	1	55	Training	Human Translator	2015	✓
UN Corpus (Ziemski et al., 2016)	UN	6	1	6	Training	Human Translator	2016	✓
MMCR4NLP (Dabre and Kurohashi, 2017)	Mix	59	5	13	Training	Mix Collection	2017	✗
GCP Corpus (Imamura and Sumita, 2018)	Speech	10	4	10	Training	Machine Translation	2018	✗
FLORES-101 (Goyal et al., 2022)	Wiki	101	10	101	Evaluation	Human Translator	2022	✓
FLORES-200 (Costa-Jussà et al., 2022)	Wiki	204	10	204	Evaluation	Human Translator	2022	✓
BPCC (Gal et al., 2023)	Many	22	13	22	Training	Mix Collection	2023	✓
XDailyDialog (Liu et al., 2023b)	DailyDialog	4	/	4	Training	Human Translator	2023	✓
MWccMatrix (Thompson et al., 2024)	Common Crawl	90	/	/	Training	Crawl	2024	✓
TED2018 (Qi et al., 2018)	TED	58	/	58	Training	Human Translator	2018	✓
TED2020 (Reimers and Gurevych, 2020)	TED	108	/	/	Training	Human Translator	2020	✓
Ours (TED2025)	TED	113	352	50	Training	Human Translator	2025	✓

Table 1: Comparison of existing multi-way parallel corpora and our constructed TED2025, highlighting key attributes such as the data source, the number of languages (#Langs), the number of domains (#Domains), the maximum parallelism supported, the type of data (training or evaluation), the collection method, and whether the corpus is open-source.

improvements in multilingual LLM performance?

RQ3: Which instruction-tuning objectives can most effectively leverage the advantages of multi-way parallel data?

We perform a comprehensive evaluation across six multilingual benchmarks to assess the benefits of using multi-way parallel data for scaling multilingual LLMs. Our results reveal that, at an equal data scale, fine-tuning on multi-way parallel data consistently outperforms training on unaligned multilingual text for both low-resource and high-resource languages (Section 4). Additionally, we identify the most effective configurations of parallelism (Section 5). Furthermore, we investigate how different instruction-tuning objectives impact LLM performance and cross-domain robustness (Section 6).

In summary, our contributions are as follows: **(i)** We construct TED2025, a 50-way parallel corpus derived from recent TED talk transcripts, covering 113 languages and 352 domains. **(ii)** To the best of our knowledge, this is the first work to leverage multi-way parallel data for scaling multilingual LLMs. We present a systematic comparison of multilingual LLM fine-tuning using multi-way versus unaligned data, analyzing their effects on zero-shot transfer and cross-lingual representation alignment. **(iii)** We explore instruction-tuning objectives specifically designed for multi-way parallel data and provide practical recommendations for optimizing multilingual LLM performance.

2 Related Work

Multi-Way Parallel Corpora. Datasets containing the same content across multiple languages (typically more than two) are known as *multi-way parallel corpora*. These corpora have demonstrated substantial benefits for machine translation (Freitag and Firat, 2020; Xu et al., 2022; Wu

et al., 2024) and cross-lingual representation alignment (Tran et al., 2020). Existing methods for constructing such corpora include mining comparable texts (Resnik et al., 1999), aligning independently collected monolingual corpora via translation pivots (Thompson et al., 2024), extracting multi-way subsets from large bilingual collections (Ramesh et al., 2022), and harvesting multilingual web crawls (Resnik and Smith, 2003; Qi et al., 2018). However, many of these resources are limited in terms of language and domain coverage. In contrast, we construct a multi-way parallel corpus derived from recent TED Talk transcripts, offering a broader and more diverse set of languages and domains.

Scaling Multilingual LLMs. The multilingual capabilities of LLMs have been significantly enhanced through two complementary strategies: continued pretraining on diverse multilingual corpora and multilingual instruction tuning. Continued pretraining on unaligned multilingual data has improved both in-language fluency and cross-lingual transfer (Ji et al., 2024; Groeneveld et al., 2024), while instruction tuning with human-curated multilingual prompts has boosted task performance across a wide range of languages (Lai et al., 2024; Üstün et al., 2024). More recently, Mu et al. (2024) demonstrated that incorporating multi-way parallel examples into in-context prompts leads to further gains in zero-shot transfer. Building on these insights, we systematically fine-tune multilingual LLMs on large-scale multi-way parallel data and quantify their impact compared to conventional unaligned approaches.

3 Experimental Setup

TED2025. We introduce TED2025, a new multi-way parallel corpus derived from the latest TED

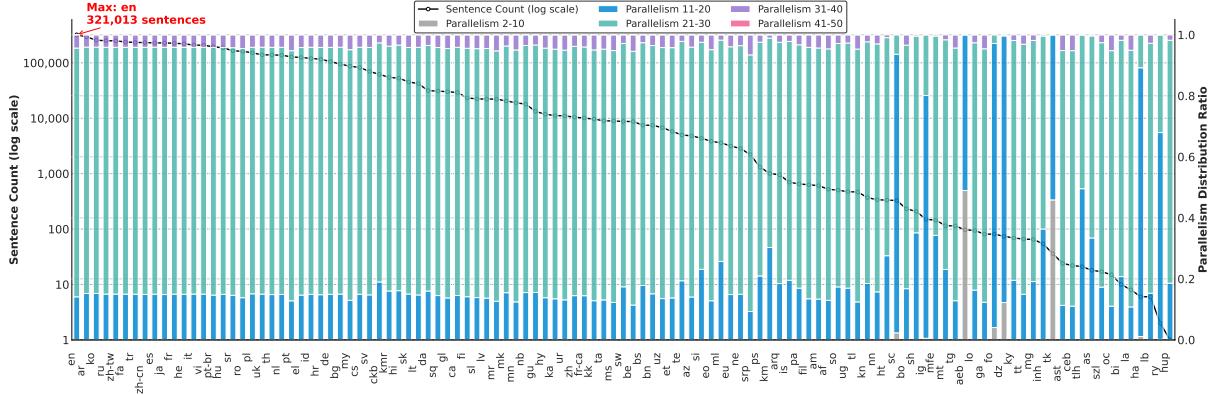


Figure 1: Distribution of sentence counts (line chart, left y-axis, log scale) and parallelism spans (bar chart, right y-axis, ratio) across languages (x-axis) in the TED2025 corpus. The parallelism spans, with a notable concentration between 21 and 30 languages, and high range even for low-resource languages.

Talk transcripts. It encompasses 113 languages with up to 50-way parallelism, making it one of the largest and most diverse resources for multilingual fine-tuning. Figure 1 illustrates the total number of sentences and the distribution of parallelism spans across languages in TED2025. Figure 2 compares translation quality, as measured by COMET-QE (Rei et al., 2020), across TED2025 (which includes 4,765 language pairs³) and other existing multi-way datasets: TED2018 (Qi et al., 2018), TED2020 (Reimers and Gurevych, 2020), and MWccMatrix (Thompson et al., 2024). Additional dataset statistics and construction details are provided in Appendix A.

We observe that: (1) the most common parallelism span in TED2025 ranges from 21 to 30 languages. Notably, many low-resource languages also achieve high degrees of parallelism, providing a solid foundation for multilingual research. (2) TED2025 contains significantly more high-quality translations (with a COMET-QE score greater than 60) compared to previous multi-way corpora. Unlike TED2020, which segments English based on punctuation, or MWccMatrix, which relies on LASER score (Artetxe and Schwenk, 2019), we use human-provided timestamps to generate cleaner and more reliable sentence alignments.

Training. To isolate the effects of multi-way parallel data, we conduct both continued pre-training (Parmar et al., 2024) and instruction tuning (Zhang et al., 2023) on TED2025⁴. We experiment with two model families and sizes: LLaMA-

³For language pairs with more than 10,000 sentence pairs, we randomly sample 10,000.

⁴For RQ1 and RQ2, we focus on continued pretraining to avoid interference from instruction tuning. For RQ3, we instead focus on instruction tuning to assess its specific impact.

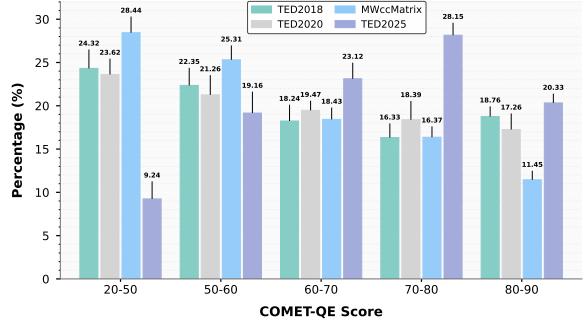


Figure 2: Comparison of translation quality between TED2025 and existing multi-way datasets, including TED2018 (Qi et al., 2018), TED2020 (Reimers and Gurevych, 2020), MWccMatrix (Thompson et al., 2024), using COMET-QE score.

Benchmark	Task	#Langs	Metric
MMMLU	Understanding	14	Acc
XCOPA	Reasoning	11	Acc
FLORES-101	Generation	101	BLEU/COMET
FLORES-200	Generation	204	BLEU/COMET
xIFEval	Instruction Following	17	Acc
SIB	Text Classification	204	Acc

Table 2: Overview of evaluation benchmarks, including task types, the number of languages (#Langs) involved, and the metrics used for assessment.

3.1-8B / LLaMA-3.1-8B-Instruct and Qwen-2.5-14B / Qwen-2.5-14B-Instruct (available on Hugging Face⁵). To make fine-tuning feasible, we employ Low-Rank Adaptation (LoRA) (Hu et al., 2022) instead of performing full parameter updates. Full hyperparameter settings and training configurations are provided in Appendix B.

Evaluation and Metrics. We evaluate our models on five widely adopted multilingual benchmarks in a zero-shot setting, covering a range of tasks: un-

⁵<https://huggingface.co>

derstanding (MMMLU), reasoning (XCOPA; Ponti et al., 2020), generation (FLORES-101; Goyal et al., 2022 and FLORES-200; Costa-Jussà et al., 2022), instruction following (xIFEval; Huang et al., 2025), and text classification (SIB-200; Adelani et al., 2024). Table 2 summarizes these benchmarks along with their associated evaluation metrics. Additionally, we categorize the languages in each benchmark as low-resource or high-resource based on the classification in Costa-Jussà et al. (2022).

4 Effectiveness of Multi-Way Corpora

We investigate the impact of training on multi-way parallel data on the performance of a multilingual LLM across three key dimensions: downstream performance on multilingual benchmarks (Section 4.1), zero-shot cross-lingual transfer to unseen languages (Section 4.2) and cross-lingual alignment of representations within the model’s internal embeddings (Section 4.3).

For fairness, we fix the total continued pre-training data at 5 million tokens (5M) and evaluate two LLM backbones (LLaMA-3-8B and Qwen-2.5-14B) under three conditions: (1) **Multi-Way**: Pretraining on our multi-way parallel corpus (TED2025); (2) **Unaligned**: Pretraining on an unaligned multilingual corpus (DCAD-2000; Shen et al., 2025); (3) **Baseline**: The original pretrained checkpoint without additional data.

4.1 Downstream Performance

We evaluate all variants in a zero-shot setting across four benchmarks—MMMLU, XCOPA, FLORES-101, and xIFEval—that cover understanding, reasoning, generation, and instruction following. The full results are reported in Table 3. Our findings show that, across all tasks, *Multi-Way* consistently outperforms both *Baseline* and *Unaligned* across both low- and high-resource languages. For example, on MMMLU, *Multi-Way* achieves accuracies of 22.48/41.38 in low/high-resource languages, compared to 18.27/33.72 for the *Baseline* and 19.64/36.26 for *Unaligned*. Similar improvements are observed on FLORES-101 and xIFEval. These results demonstrate that multi-way alignment provides stronger cross-lingual supervision, thereby enhancing both discriminative and generative capabilities.

4.2 Zero-Shot Cross-Lingual Transfer

To assess generalization to unseen languages (Lai et al., 2022a; Zhao et al., 2025), we evaluate on

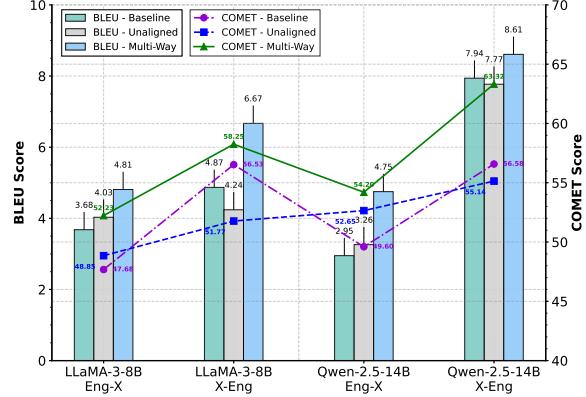


Figure 3: Cross-lingual transfer performance comparison between Baseline, Unaligned and Multi-Way pre-training on the FLORES-200 benchmark with BLEU (bar chart, left y-axis) and COMET (line chart, right y-axis) for LLaMA-3-8B and Qwen-2.5-14B models.

FLORES-200. We exclude all languages in the evaluation subset from training and assess the English↔X translation quality. As shown in Figure 3, the *Multi-Way* model significantly outperforms both *Baseline* and *Unaligned* in both translation directions. This highlights that explicit multi-way supervision promotes language-agnostic representations, enabling robust zero-shot transfer. We further explore this hypothesis in Section 4.3, where we analyze differences in cross-lingual representation alignment across models.

4.3 Cross-Lingual Representation Alignment

We further analyze the alignment of internal representations across models. To be specific, for 32 randomly selected aligned languages (with 100 sentences each) from TED2025, we compute the following four metrics: average *cosine similarity* between parallel sentence embeddings, *Centered Kernel Alignment* (CKA) between representation matrices (Kornblith et al., 2019), Cross-lingual sentence retrieval accuracy at P@1, P@5, and P@10 (Conneau et al., 2017) and SVCCA score (Raghuram et al., 2017).

Table 4 shows that *Multi-Way* outperforms the other models, yielding higher CKA and better retrieval accuracy. Figure 4 further corroborates these results: *Multi-Way* demonstrates denser SVCCA alignments, particularly for linguistically distant language pairs. These metrics confirm that multi-way pretraining promotes a more coherent, language-agnostic embedding space, which drives the observed improvements in downstream performance and cross-lingual transfer.

	MMMLU		XCOPA		FLORES-101 (Eng-X)				FLORES-101 (X-Eng)				xIFEval	
					BLEU		COMET		BLEU		COMET			
	low	high	low	high	low	high	low	high	low	high	low	high	low	high
Baseline	18.27	33.72	23.46	34.29	6.03	11.67	57.15	61.03	13.37	22.49	75.24	82.32	17.14	24.43
Unaligned	19.64	36.26	24.62	34.76	6.12	11.78	57.51	62.11	13.84	22.74	75.82	82.58	17.28	24.44
Multi-Way	22.48	41.38	27.58	57.22	6.32	12.08	58.06	67.44	14.45	25.03	76.25	86.43	18.79	27.41

(a) LLaMA-3-8B

	MMMLU		XCOPA		FLORES-101 (Eng-X)				FLORES-101 (X-Eng)				xIFEval	
					BLEU		COMET		BLEU		COMET			
	low	high	low	high	low	high	low	high	low	high	low	high	low	high
Baseline	35.24	49.55	62.25	72.00	7.45	11.05	57.22	67.16	16.54	20.24	67.23	74.29	27.63	32.40
Unaligned	35.61	51.32	62.59	74.06	7.62	11.60	57.85	70.85	16.86	21.02	67.61	75.97	27.92	35.54
Multi-Way	36.64	55.81	63.24	79.52	8.07	13.11	58.94	80.56	17.36	23.26	68.59	81.33	28.64	40.95

(b) Qwen-2.5-14B

Table 3: Performance (%) comparison of different models across multilingual benchmarks. The **Multi-Way** approach with blue background demonstrates consistent superiority in both low-resource (left columns) and high-resource (right columns) scenarios.

	LLaMA-3-8B						Qwen-2.5-14B					
	Baseline	Unaligned	Multi-Way	Baseline	Unaligned	Multi-Way	Baseline	Unaligned	Multi-Way	Baseline	Unaligned	Multi-Way
Cosine (\uparrow)	0.27	0.27	0.30 _{+0.03}	0.29	0.27	0.32 _{+0.03}						
CKA (\uparrow)	0.54	0.54	0.60 _{+0.06}	0.56	0.57	0.63 _{+0.07}						
Retrieval (\uparrow)	P@1 P@5 P@10	0.09 0.24 0.35	0.07 0.23 0.33	-0.02 -0.01 -0.02	0.13 _{+0.04} 0.27 _{+0.03} 0.39 _{+0.04}	0.30 _{+0.03} 0.33 _{+0.07} 0.42 _{+0.06}	0.12	0.14	0.19 _{+0.07}	0.26	0.28	0.33 _{+0.07}
SVCCA (\uparrow)	0.55	0.55	0.61 _{+0.06}	0.57	0.58	0.63 _{+0.06}						

Table 4: Cross-lingual representation alignment results. Improvement margins (colored red/blue) are shown relative to Baseline. **Bolded** Multi-Way results with **red** annotations indicate consistent improvements, particularly in cross-lingual retrieval tasks (e.g., +0.07 P@1 improvement for Qwen-2.5-14B).

5 Impact Factors

In Section 4, we demonstrate that using multi-way parallel data can significantly enhance the multilingual capabilities of LLMs. To better understand the factors driving this improvement, we analyze two key aspects⁶, while keeping the total pretraining fixed at 5 million tokens: (1) **Degree of parallelism**: the number of languages aligned in each training example (Section 5.1). (2) **English as a pivot**: the impact of including versus excluding English in multi-way groups (Section 5.2).

5.1 Degree of Parallelism

We construct datasets with parallelism levels ranging from 2 to 40 languages per example, sampled from the TED2025 corpus, while always keeping 5M tokens. Figure 5 shows model performance

across a range of tasks for each setting. For bidirectional machine translation (FLORES-101, Eng→X and X→Eng), performance steadily improves with higher parallelism, suggesting that broader semantic alignment enhances cross-lingual generation and fluency. In contrast, for non-generative tasks (reasoning and understanding), accuracy tends to peak at small parallelism (around 6–10 languages) before deteriorating. We attribute this decline to two factors: (1) Excessive linguistic diversity can obscure shared semantic patterns. (2) With a fixed token budget, each language receives fewer tokens, limiting the ability to learn language-specific features.

5.2 English as Pivot

English is widely used as a pivot language in multilingual MT and NLP (Kim et al., 2019; Mallinson et al., 2018; Lai et al., 2023b). To explore its im-

⁶Additional analyses on data size and language-family combinations can be found in Appendix C.

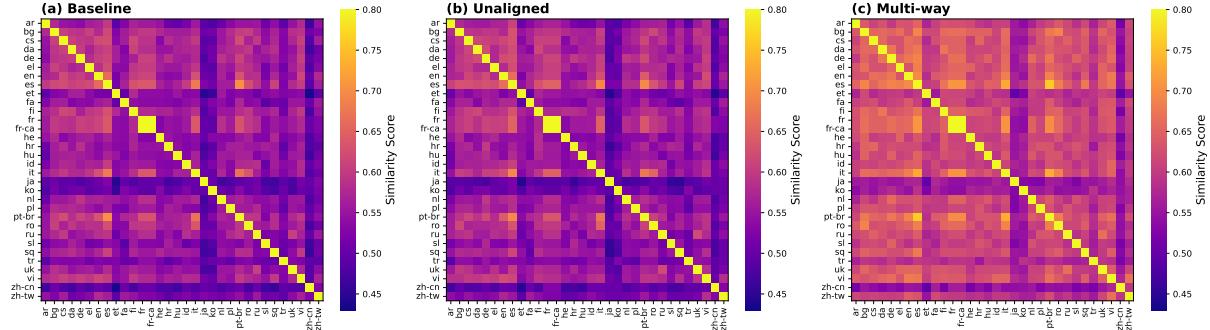


Figure 4: SVCCA alignment comparison between the Multi-Way, Unaligned and Baseline models across 32-way language pairs.

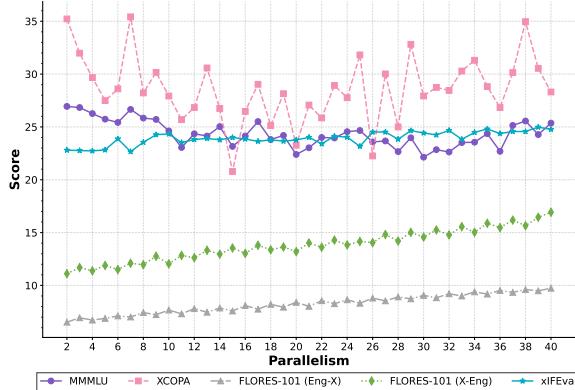


Figure 5: Performance (%) of continued pretraining models on downstream tasks with varying degrees of parallelism.

312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
pact, we form five groups⁷, each with both English-included and English-excluded variants across six languages. Figure 6 compares their performance across various tasks.

In tasks involving understanding and reasoning, groups that include English consistently outperform those in other language groups by an average of 2–4 percentage points. This suggests that English, as a high-resource “semantic anchor”, helps stabilize embedding alignment and facilitates transfer learning, particularly on complex tasks. Interestingly, for machine translation (FLORES-101) and some instruction-following tasks (xIEval), English-inclusive groups show slightly lower performance. We attribute this to two primary factors: (1) English occupies tokens that could otherwise be used to align non-English language pairs directly. (2) The model may overly rely on English as an intermediary, which reduces its ability to directly transfer knowledge between non-English languages. These findings highlight that English’s role as a pivot language is task-dependent. While it

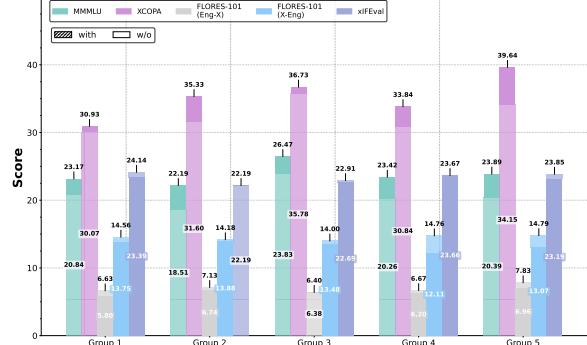


Figure 6: Performance (%) comparison of models with and without (w/o) English across five different language groupings.

can enhance semantic coherence in understanding and reasoning tasks, it may hinder direct multilingual transfer in generative tasks. Consequently, the inclusion of English in multilingual pretraining should be carefully considered based on the target application.

6 Instruction Tuning

In Sections 4 and 5, we demonstrate that using multi-way data can significantly improve multilingual performance. In this section, we further investigate whether instruction tuning can also enhance multilingual performance effectively. Specifically, we address the following key questions⁸: (1) Which of the different instruction fine-tuning objectives, built on multi-way parallel data, is most effective? (2) Do models trained with multi-way parallel data exhibit better generalization across domains?

6.1 Instruction Tuning Objectives

Using our constructed multi-way parallel data (TED2025), we define four instruction tasks: ma-

⁷We provide the configurations of each group in Table 8 of Appendix C.2

⁸Additionally, we explore the cross-lingual alignment of representations in instruction tuned models. Detailed experimental settings and results can be found in the Appendix D.

	MMMLU				XCOPA				FLORES-101 (Eng-X)				FLORES-101 (X-Eng)				xIFEval	
	BLEU		COMET		BLEU		COMET		BLEU		COMET		low		high		low	
	low	high	low	high	low	high	low	high	low	high	low	high	low	high	low	high	low	high
Baseline	41.96	46.68	64.59	66.79	5.51	10.50	56.95	60.02	14.98	18.41	68.29	73.51	35.99	38.56				
MT	45.28	51.06	68.17	69.38	11.26	13.25	63.03	65.61	22.25	21.60	70.76	75.41	41.72	44.52				
CLTS	39.99	45.44	62.87	66.47	3.63	9.90	56.48	59.28	13.25	16.60	66.89	73.15	34.91	38.38				
MTC	40.68	44.49	64.19	65.16	5.02	9.08	56.04	57.84	14.38	18.16	67.86	73.12	34.03	38.01				
CLP	42.49	47.01	65.41	68.42	6.38	11.46	57.23	61.28	15.76	19.27	69.87	73.91	36.17	40.07				
MT + CLTS	42.23	47.94	65.00	68.77	6.19	10.51	58.53	60.33	16.83	18.98	68.70	74.91	36.47	40.30				
MT + MTC	42.69	47.54	65.12	66.83	6.35	10.73	57.53	60.24	15.13	18.45	69.28	74.00	36.10	39.03				
MT + CLP	43.20	49.07	67.29	68.47	7.31	10.51	58.40	62.64	17.75	20.77	69.18	74.47	38.49	39.67				
CLTS + MTC	41.78	45.17	63.21	65.62	5.38	9.77	55.19	58.21	14.67	16.59	67.45	72.35	34.57	36.63				
CLTS + CLP	42.82	46.74	65.30	67.12	5.58	10.88	57.84	60.41	15.41	19.33	68.84	73.58	36.99	38.83				
MTC + CLP	42.62	46.70	65.16	67.56	6.48	10.84	57.74	60.98	15.29	19.22	68.76	73.77	36.87	38.96				
MT + CLTS + MTC	41.03	46.59	64.15	66.05	5.14	10.35	56.14	59.71	14.64	18.06	67.64	73.39	35.24	38.14				
MT + CLTS + CLP	42.72	47.67	64.83	67.57	6.17	10.94	57.78	60.11	15.96	19.16	68.67	74.03	36.54	38.77				
MT + MTC + CLP	42.33	46.72	65.26	67.58	5.92	10.98	57.93	60.76	15.38	18.81	68.98	73.55	36.22	39.15				
CLTS + MTC + CLP	41.56	46.65	64.16	66.45	4.67	10.38	56.20	59.63	14.95	17.62	68.00	73.03	35.27	38.21				
MT + CLTS + MTC + CLP	43.07	47.67	66.64	67.38	7.59	12.42	57.75	60.79	17.62	19.62	71.13	75.24	36.95	40.52				

(a) LLaMA-3-8B-Instruct

	MMMLU				XCOPA				FLORES-101 (Eng-X)				FLORES-101 (X-Eng)				xIFEval	
	BLEU		COMET		BLEU		COMET		BLEU		COMET		low		high		low	
	low	high	low	high	low	high	low	high	low	high	low	high	low	high	low	high	low	high
Baseline	63.61	68.14	78.76	82.22	7.92	12.79	61.49	66.64	16.76	20.48	68.16	70.32	47.68	52.40				
MT	68.52	72.83	82.95	86.95	11.90	16.60	64.96	70.90	20.43	26.22	72.25	73.51	53.24	55.83				
CLTS	61.28	67.61	76.77	79.74	7.20	10.70	60.90	64.50	15.75	18.21	67.01	69.83	45.73	51.83				
MTC	62.59	65.63	76.87	81.46	6.07	12.29	61.48	65.61	13.84	19.71	67.12	68.65	47.53	50.53				
CLP	64.28	70.06	79.69	83.02	9.54	13.10	62.18	68.35	17.54	21.68	69.24	70.68	48.61	52.66				
MT + CLTS	63.93	69.53	78.94	82.75	8.73	13.91	61.87	67.09	18.20	21.50	68.45	71.90	48.16	54.14				
MT + MTC	65.33	69.12	80.27	83.34	8.64	13.94	61.71	67.42	17.87	22.27	68.60	72.01	49.46	53.53				
MT + CLP	64.95	70.51	80.68	83.52	10.30	13.65	62.37	69.32	18.43	22.19	70.01	70.80	48.85	53.15				
CLTS + MTC	62.67	67.19	78.44	81.33	7.61	12.60	61.04	65.90	16.25	19.93	67.47	70.00	46.85	51.63				
CLTS + CLP	64.38	68.65	78.89	82.40	8.65	13.57	61.75	67.57	17.73	20.81	68.65	71.30	47.82	53.04				
MTC + CLP	64.04	68.92	79.16	82.56	8.81	13.05	62.26	67.34	16.97	21.31	68.98	71.20	48.06	53.32				
MT + CLTS + MTC	64.13	68.55	79.51	82.63	8.26	13.30	62.24	67.01	17.69	21.29	68.38	70.74	48.54	52.80				
MT + CLTS + CLP	64.50	68.82	80.43	83.22	9.05	13.93	62.38	67.56	18.66	22.01	68.76	71.62	48.96	53.74				
MT + MTC + CLP	64.71	69.52	80.26	83.32	9.15	13.37	62.82	67.66	17.75	21.70	69.33	71.50	48.93	53.41				
CLTS + MTC + CLP	63.81	69.00	79.70	82.52	8.69	13.48	62.25	67.04	17.55	20.75	68.23	70.96	47.68	52.55				
MT + CLTS + MTC + CLP	64.94	68.19	80.25	82.42	9.87	13.56	62.86	67.04	16.94	21.69	68.17	71.34	49.03	53.72				

(b) Qwen-2.5-14B-Instruct

Table 5: Downstream task performance (%) of LLaMA-3-8B-Instruct and Qwen-2.5-14B-Instruct models trained with different instruction objectives, including MT, CLTS, MTC and CLP.

chine translation (MT), cross-lingual text similarity (CLTS), multilingual text classification (MTC), and cross-language paraphrasing (CLP). Table 5 reports their impact on downstream benchmarks. Table 6 summarizes the prompt templates and output formats for each task.

We observe that the improvements in MT are the largest and most stable in both high-resource and low-resource languages. This can primarily be attributed to the fact that, as a token-level supervised generation task, translation strengthens cross-lingual syntactic and semantic consistency, making it a particularly robust task for broad multilingual generalization. In contrast, CLTS and MTC

show smaller drops across different tasks, which we attribute to the coarser granularity of similarity judgments that may not provide fine-grained alignment signals. Moreover, discrete class labels lack the expressive power to capture subtle semantic distinctions. Although CLP shows similar advantages to MT in generation tasks, its advantages are narrower in scope and cannot be widely generalized.

Interestingly, the combination of tasks did not significantly improve performance, which may be due to several reasons: First, the objectives of MT and CLP differ; while MT emphasizes accuracy, CLP focuses more on the diversity of expression, which may make it difficult for the model to bal-

Task	Prompt
Machine Translation (MT)	Translate the following {src_lang_1}, {src_lang_2}, ..., {src_lang_m} sentence to {tgt_lang_1}, {tgt_lang_2}, ..., {tgt_lang_n}.\\n {src_lang_1} Sentence: {src_txt_1}.\\n {src_lang_2} Sentence: {src_txt_2}.\\n ... {src_lang_m} Sentence: {src_txt_m}.\\n Translation:\\n {tgt_lang_1} Sentence: {tgt_txt_1}.\\n {tgt_lang_2} Sentence: {tgt_txt_2}.\\n ... {tgt_lang_n} Sentence: {tgt_txt_n}.\\n
Cross-Lingual Text Similarity (CLTS)	Given the sentences below in different languages, rate how similar their meanings are on a scale of 0 to 1, where 0 means completely dissimilar and 1 means identical meanings.\\n {lang_1} Sentence: {txt_1}.\\n {lang_2} Sentence: {txt_2}.\\n ... {lang_m} Sentence: {txt_m}.\\n Similarity: {sim_score}.
Multilingual Text Classification (MTC)	Classify the following sentence in {lang_1}, {lang_2}, ..., {lang_m} into one of the following categories: {domain_list}.\\n {lang_1} Sentence: {txt_1}.\\n {lang_2} Sentence: {txt_2}.\\n ... {lang_m} Sentence: {txt_m}.\\n Categories: {target_domain}.
Cross-Lingual Paraphrasing (CLP)	Paraphrase the following {src_lang} sentence in {tgt_lang}.\\n {src_lang} Sentence: {src_txt}.\\n Paraphrasing:\\n {tgt_lang} Sentence: {tgt_txt}.

Table 6: Instruction prompts used for four multilingual tasks: machine translation (MT), cross-lingual text similarity (CLTS), multilingual text classification (MTC), and cross-lingual paraphrasing (CLP). Each prompt is designed to reflect the task’s specific objective and structure.

382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
ance these two goals during training, thereby affecting performance. Second, interference between multiple tasks may arise, making it challenging for the model to focus on optimizing the optimal goal of each task, particularly when the task objectives are too similar, amplifying the interference effect. Furthermore, there is overlap in task design between MT and CLP, which may cause the model to encounter redundant training signals when processing both tasks, preventing it from fully utilizing the unique value of each task. Finally, joint training of multiple tasks increases the complexity of model training, potentially leading to unstable gradient updates and affecting the model’s convergence.

6.2 Cross-Domain Generalization

To evaluate cross-domain transfer (Lai et al., 2022b; Liu et al., 2023a), we extract domain-labeled subsets from TED2025 according to the taxonomy of the SIB-200 benchmark. Instruction tuning is then performed on each domain using the MTC objective, with the resulting models evaluated across all other domains within SIB-200. Figure 7 illustrates the transfer performance.

Overall, instruction tuning with multi-way parallel data significantly improves domain transfer. The rich cross-lingual and cross-domain signals allow the model to learn domain-invariant features, enhancing robustness when confronted with novel topics and linguistic contexts. However, transfer performance remains limited in domains such as politics, sports, travel, and geography. We hypothesize that the high topical diversity, coupled with the relative sparsity of domain-specific examples in the training data, hinders the model’s ability to capture specialized patterns. Overcoming these limitations may require more balanced domain coverage or

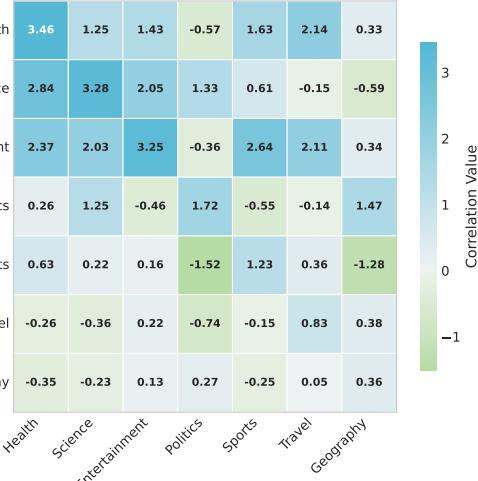


Figure 7: Cross-domain generalization performance of instruction-tuned models using multi-way parallel data. Models trained on one domain are evaluated on all other domains from the SIB-200 benchmark.

targeted data augmentation strategies.

7 Conclusion

In this paper, we construct a large-scale, high-quality multi-way parallel dataset covering 113 languages, with a maximum parallel degree of 50. This dataset provides a strong foundation for investigating the multilingual adaptation of LLMs. Using this dataset, we systematically explore best practices for adapting LLMs to multilingual tasks via multi-way parallel data. Our experiments reveal that multi-way data offers substantial advantages for both continued pretraining and instruction tuning, resulting in improved cross-lingual and cross-domain generalization. We further analyze key factors influencing model performance, including the degree of parallelism, the language combination strategies, and instruction training objectives.

435 Limitations

436 This work has the following limitations: **(i)** Due
437 to limited computational resources, we employed
438 parameter-efficient fine-tuning (PEFT) methods
439 (LoRA) instead of full-parameter fine-tuning.
440 While recent studies have demonstrated that LoRA
441 achieves performance comparable to full fine-
442 tuning across various tasks, our conclusions may
443 still benefit from validation under full fine-tuning or
444 alternative PEFT methods such as adapters or prefix
445 tuning. **(ii)** Although our constructed dataset sur-
446 passes existing multi-way parallel corpora in both
447 language coverage and maximum parallel degree,
448 its overall size remains modest compared to large-
449 scale unaligned multilingual datasets. To fully un-
450 lock the potential of multi-way parallel data for
451 LLM adaptation, future work will focus on scal-
452 ing up the dataset to further enhance multilingual
453 performance.

454 References

- 455 David Ifeoluwa Adelani, Hannah Liu, Xiaoyu Shen,
456 Nikita Vassilyev, Jesujoba O. Alabi, Yanke Mao, Hao-
457 nan Gao, and En-Shiun Annie Lee. 2024. **SIB-200:**
458 **A simple, inclusive, and big evaluation dataset for**
459 **topic classification in 200+ languages and dialects.**
460 In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 226–245, St. Julian’s, Malta. Association for Computational Linguistics.
- 465 Mikel Artetxe and Holger Schwenk. 2019. **Mas-**
466 **sively multilingual sentence embeddings for zero-**
467 **shot cross-lingual transfer and beyond.** *Transactions*
468 *of the Association for Computational Linguistics*,
469 7:597–610.
- 470 Christos Christodoulopoulos and Mark Steedman.
471 2015. A massively parallel corpus: the bible in
472 100 languages. *Language resources and evaluation*,
473 49:375–395.
- 474 Alexis Conneau, Guillaume Lample, Marc’Aurelio Ran-
475 zato, Ludovic Denoyer, and Hervé Jégou. 2017.
476 Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.
- 477 Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha
478 Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe
479 Kalbassi, Janice Lam, Daniel Licht, Jean Maillard,
480 and 1 others. 2022. No language left behind: Scaling
481 human-centered machine translation. *arXiv preprint arXiv:2207.04672*.
- 482 Raj Dabre and Sadao Kurohashi. 2017. Mmcr4nlp: mul-
483 tilingual multiway corpora repository for natural lan-
484 guage processing. *arXiv preprint arXiv:1710.01025*.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. A survey on in-context learning. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1107–1128, Miami, Florida, USA. Association for Computational Linguistics.	487
Markus Freitag and Orhan Firat. 2020. Complete multi-lingual neural machine translation. In <i>Proceedings of the Fifth Conference on Machine Translation</i> , pages 550–560, Online. Association for Computational Linguistics.	488
Jay P Gala, Pranjal A Chitale, AK Raghavan, Varun Gumma, Sumanth Doddapaneni, Kumar M Aswanth, Janki Atul Nawale, Anupama Sujatha, Ratish Pudupully, Vivek Raghavan, and 1 others. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. <i>Transactions on Machine Learning Research</i> , 2023.	489
Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. <i>Transactions of the Association for Computational Linguistics</i> , 10:522–538.	490
Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Author, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, and 24 others. 2024. OLMo: Accelerating the science of language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15789–15809, Bangkok, Thailand. Association for Computational Linguistics.	491
Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. <i>arXiv preprint arXiv:2009.03300</i> .	492
Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In <i>International Conference on Learning Representations</i> .	493
Haoyang Huang, Tianyi Tang, Dongdong Zhang, Xin Zhao, Ting Song, Yan Xia, and Furu Wei. 2023. Not all languages are created equal in LLMs: Improving multilingual capability by cross-lingual-thought prompting. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 12365–12394, Singapore. Association for Computational Linguistics.	494
Kaiyu Huang, Fengran Mo, Xinyu Zhang, Hongliang Li, You Li, Yuanchi Zhang, Weijian Yi, Yulong Mao,	495

545	Jinchen Liu, Yuzhuang Xu, and 1 others. 2024. A survey on large language models with multilingualism: Recent advances and new frontiers. <i>arXiv preprint arXiv:2405.10936</i> .	600
546		601
547		602
548		603
549	Xu Huang, Wenhao Zhu, Hanxu Hu, Conghui He, Lei Li, Shujian Huang, and Fei Yuan. 2025. Benchmax: A comprehensive multilingual evaluation suite for large language models. <i>arXiv preprint arXiv:2502.07346</i> .	604
550		605
551		606
552		607
553		608
554	Kenji Imamura and Eiichiro Sumita. 2018. Multilingual parallel corpus for global communication plan. In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	609
555		610
556		611
557		612
558		613
559		614
560	Shaoxiong Ji, Zihao Li, Indraneil Paul, Jaakko Paavola, Peiqin Lin, Pinzhen Chen, Dayyán O’Brien, Hengyu Luo, Hinrich Schütze, Jörg Tiedemann, and 1 others. 2024. Emma-500: Enhancing massively multilingual adaptation of large language models. <i>arXiv preprint arXiv:2409.17892</i> .	615
561		616
562		617
563		618
564		619
565		620
566	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM computing surveys</i> , 55(12):1–38.	621
567		622
568		623
569		624
570		625
571	Yunsu Kim, Petre Petrov, Pavel Petrushkov, Shahram Khadivi, and Hermann Ney. 2019. Pivot-based transfer learning for neural machine translation between non-English languages. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 866–876, Hong Kong, China. Association for Computational Linguistics.	626
572		627
573		628
574		629
575		630
576		631
577		632
578		633
579		634
580	Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. 2019. Similarity of neural network representations revisited. In <i>International conference on machine learning</i> , pages 3519–3529. PMLR.	635
581		636
582		637
583		638
584		639
585	Viet Dac Lai, Nghia Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Huu Nguyen. 2023a. ChatGPT beyond English: Towards a comprehensive evaluation of large language models in multilingual learning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 13171–13189, Singapore. Association for Computational Linguistics.	640
586		641
587		642
588		643
589		644
590		645
591		646
592		647
593	Wen Lai, Alexandra Chronopoulou, and Alexander Fraser. 2022a. m^4 adapter: Multilingual multi-domain adaptation for machine translation with a meta-adapter. In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 4282–4296, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	648
594		649
595		650
596		651
597		652
598		653
599		654
	Wen Lai, Alexandra Chronopoulou, and Alexander Fraser. 2023b. Mitigating data imbalance and representation degeneration in multilingual machine translation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 14279–14294, Singapore. Association for Computational Linguistics.	655
		656
	Wen Lai, Jindřich Libovický, and Alexander Fraser. 2022b. Improving both domain robustness and domain adaptability in machine translation. In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 5191–5204, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	657
		658
	Wen Lai, Mohsen Mesgar, and Alexander Fraser. 2024. LLMs beyond English: Scaling the multilingual capability of LLMs with cross-lingual feedback. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 8186–8213, Bangkok, Thailand. Association for Computational Linguistics.	659
		660
	Fangyu Liu, Qianchu Liu, Shruthi Bannur, Fernando Pérez-García, Naoto Usuyama, Sheng Zhang, Tristan Naumann, Aditya Nori, Hoifung Poon, Javier Alvarez-Valle, Ozan Oktay, and Stephanie L. Hyland. 2023a. Compositional zero-shot domain transfer with text-to-text models. <i>Transactions of the Association for Computational Linguistics</i> , 11:1097–1113.	661
		662
	Zeming Liu, Ping Nie, Jie Cai, Haifeng Wang, Zheng-Yu Niu, Peng Zhang, Mrinmaya Sachan, and Kaiping Peng. 2023b. XDailyDialog: A multilingual parallel dialogue corpus. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12240–12253, Toronto, Canada. Association for Computational Linguistics.	663
		664
	Jonathan Mallinson, Rico Sennrich, and Mirella Lapata. 2018. Sentence compression for arbitrary languages via multilingual pivoting. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2453–2464, Brussels, Belgium. Association for Computational Linguistics.	665
		666
	Yongyu Mu, Peinan Feng, Zhiqian Cao, Yuzhang Wu, Bei Li, Chenglong Wang, Tong Xiao, Kai Song, Tongran Liu, Chunliang Zhang, and JingBo Zhu. 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, Miami, Florida, USA. Association for Computational Linguistics.	667
		668
	Jupinder Parmar, Sanjeev Satheesh, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. 2024. Reuse, don’t retrain: A recipe for continued pre-training of language models. <i>arXiv preprint arXiv:2407.07263</i> .	669
		670
	Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020.	671
		672

656	XCOPA: A multilingual dataset for causal common-	713
657	sense reasoning.	714
658	In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2362–2376, Online. Association for Computational Linguistics.	
661	Ye Qi, Devendra Sachan, Matthieu Felix, Sarguna Pad-	
662	manabhan, and Graham Neubig. 2018. When and	
663	why are pre-trained word embeddings useful for neu-	
664	ral machine translation?	
665	In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 529–535, New Orleans, Louisiana. Association for Computational Linguistics.	
666		
667		
668		
669		
670	Libo Qin, Qiguang Chen, Yuhang Zhou, Zhi Chen,	
671	Yinghui Li, Lizi Liao, Min Li, Wanxiang Che, and	
672	Philip S Yu. 2024. Multilingual large language	
673	model: A survey of resources, taxonomy and fron-	
674	tiers.	
675	<i>arXiv preprint arXiv:2404.04925</i> .	
676		
677		
678		
679		
680	Maithra Raghu, Justin Gilmer, Jason Yosinski, and	
681	Jascha Sohl-Dickstein. 2017. Svcca: Singular vec-	
682	tor canonical correlation analysis for deep learning	
683	dynamics and interpretability.	
684	<i>Advances in neural information processing systems</i> , 30.	
685		
686		
687		
688		
689		
690		
691	Gowtham Ramesh, Sumanth Doddapaneni, Aravindh	
692	Bheemraj, Mayank Jobanputra, Raghavan AK,	
693	Ajitesh Sharma, Sujit Sahoo, Harshita Diddee, Ma-	
694	halakshmi J, Divyanshu Kakwani, Navneet Kumar,	
695	Aswin Pradeep, Srihari Nagaraj, Kumar Deepak,	
696	Vivek Raghavan, Anoop Kunchukuttan, Pratyush Ku-	
697	mar, and Mitesh Shantadevi Khapra. 2022. Saman-	
698	tar: The largest publicly available parallel corpora	
699	collection for 11 Indic languages.	
700	<i>Transactions of the Association for Computational Linguistics</i> , 10:145–	
701	162.	
702		
703	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	
704	Lavie. 2020. COMET: A neural framework for MT	
705	evaluation.	
706	In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2685–2702, Online. Association for Computational Linguistics.	
707	Nils Reimers and Iryna Gurevych. 2020. Making	
708	monolingual sentence embeddings multilingual us-	
709	ing knowledge distillation.	
710	In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4512–4525,	
711	Online. Association for Computational Linguistics.	
712		
713	Philip Resnik, Mari Broman Olsen, and Mona Diab.	
714	1999. The bible as a parallel corpus: Annotating the	
715	‘book of 2000 tongues’.	
716	<i>Computers and the Humanities</i> , 33:129–153.	
717		
718		
719		
720		
721		
722	Philip Resnik and Noah A. Smith. 2003. The web as a	
723	parallel corpus.	
724	<i>American Journal of Computational Linguistics</i> , 29(3):349–380.	
725		
726	Yingli Shen, Wen Lai, Shuo Wang, Xueren Zhang,	
727	Kangyang Luo, Alexander Fraser, and Maosong Sun.	
728	2025. Dcad-2000: A multilingual dataset across	
729		
730		
731		
732		
733		
734		
735		
736		
737	Di Wu, Shaomu Tan, Yan Meng, David Stap, and	
738	Christof Monz. 2024. How far can 100 samples	
739	go? unlocking zero-shot translation with tiny multi-	
740	parallel data.	
741	In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 15092–	
742	15108, Bangkok, Thailand. Association for Compu-	
743	tational Linguistics.	
744	Yulin Xu, Zhen Yang, Fandong Meng, and Jie Zhou.	
745	2022. EAG: Extract and generate multi-way aligned	
746	corpus for complete multi-lingual neural machine	
747	translation.	
748	In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8141–8153, Dublin,	
749	Ireland. Association for Computational Linguistics.	
750		
751	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang,	
752	Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tian-	
753	wei Zhang, Fei Wu, and 1 others. 2023. Instruction	
754	tuning for large language models: A survey.	
755	<i>arXiv preprint arXiv:2308.10792</i> .	
756	Yiran Zhao, Wenxuan Zhang, Huiming Wang, Kenji	
757	Kawaguchi, and Lidong Bing. 2025. AdaMergeX:	
758	Cross-lingual transfer with large language models via	
759	adaptive adapter merging.	
760	In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 9785–9800, Albuquerque, New Mexico.	
761	Association for Computational Linguistics.	
762		
763		
764		
765	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan	
766	Ye, and Zheyuan Luo. 2024. LlamaFactory: Unified	
767	efficient fine-tuning of 100+ language models.	
768	In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3:</i>	
769		

770 *System Demonstrations*), pages 400–410, Bangkok,
771 Thailand. Association for Computational Linguistics.

772 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Sid-
773 dhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,
774 and Le Hou. 2023. Instruction-following eval-
775 uation for large language models. *arXiv preprint*
776 *arXiv:2311.07911*.

777 Michał Ziemska, Marcin Junczys-Dowmunt, and Bruno
778 Pouliquen. 2016. **The United Nations parallel cor-**
779 **pus v1.0.** In *Proceedings of the Tenth International*
780 *Conference on Language Resources and Evaluation*
781 (*LREC’16*), pages 3530–3534, Portorož, Slovenia.
782 European Language Resources Association (ELRA).

A Statistics of the Constructed TED2025

In Section 3, we introduce TED2025 dataset and present the distribution of sentence counts across languages, along with the overall degree of parallelism (Figure 1) and translation quality compared with other multi-way corpora (Figure 2). To provide a more comprehensive overview, we further analyze domain coverage, fine-grained variations in parallelism, and the distribution of bitexts with respect to both quantity and quality.

Domain Coverage. The TED2025 dataset is a multi-domain, multi-way parallel corpus encompassing 352 domains, with domain labels derived from TED Talks. Table 9 presents statistics on the number of talks per domain. We observe that the domains of global issues, education, technology, art, and business constitute the top five in terms of talk count. This statistical overview offers valuable insights into the dataset’s structure and facilitates a deeper understanding of its composition. Moreover, the domain labels enhance the dataset’s suitability for cross-lingual text classification tasks, positioning TED2025 as a robust benchmark for multilingual text classification.

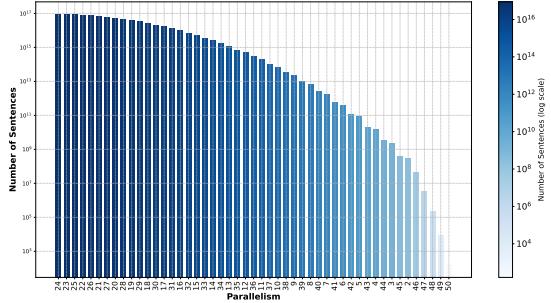


Figure 8: Fine-grained parallelism in TED2025 dataset by tuple size.

Fine-grained Parallelism. In Figure 1, we present an approximate count (ratio) of parallelism across different languages. For a more intuitive understanding of fine-grained parallelism in the TED2025 dataset, Figure 8 shows the tuple size corresponding to parallelism. Note that we count tuples rather than individual sentences. For instance, in a combination of six languages—English, French, Spanish, Russian, Arabic, and Chinese—a large tuple (en, fr, es, ru, ar, zh) encompasses all possible language combinations. The total number of such combinations is given by $C_6^1 + C_6^2 + C_6^3 + C_6^4 + C_6^5 + C_6^6 = 6 + 15 + 20 + 15 + 6 + 1 = 63$

In contrast, the corresponding sentence count for this tuple is $6 \times 1 + 15 \times 2 + 20 \times 3 + 15 \times 4 + 6 \times 5 + 1 \times 6 = 192$ sentences. This tuple-based counting method offers a more precise analysis of parallelism in the dataset.

Quantity and Quality of Bitext. In Figure 2, we compare the translation quality of TED2025 with existing multi-way parallel datasets, demonstrating the overall effectiveness of TED2025 translations. To offer a more comprehensive and intuitive view of translation quality, Table 10, 11, 12, 13, 14, and 15 report COMET-QE score for all 4,765 language pairs included in the dataset. This analysis highlights that, despite being a multi-way parallel corpus, TED2025 can be readily decomposed into bilingual sentence pairs, making it suitable for training machine translation models or fine-tuning LLMs on specific language pairs. By providing both the number of bilingual pairs and their associated translation quality, we aim to support researchers in selecting suitable data for their specific translation tasks and in optimizing performance on targeted language pairs.

B Details of Experimental Setup.

LoRA		Training	
rank	8	batch size	8
alpha	32	learning rate	1e-04
dropout	0.1	lr schedule	cosine
target	all	warmup ratio	0.1

Table 7: Hyper-parameters for continued pretraining and instruction tuning using LLaMA-Factory.

Training and Inference Setup. Due to computational resource constraints, we adopt LoRA (Hu et al., 2022) for continued pretraining and instruction fine-tuning of LLMs. All experiments are conducted using the LLaMA-Factory platform⁹ (Zheng et al., 2024), with training hyperparameters detailed in Table 7. Each experiment is run on 8 NVIDIA A100 (80GB) GPUs. For inference, we utilize the vLLM toolkit¹⁰, and the prompt templates used for each benchmark are partially sourced from PromptSource¹¹ as well as the respective original papers.

⁹<https://github.com/hiyouga/LLaMA-Factory>

¹⁰<https://github.com/vllm-project/vllm>

¹¹<https://github.com/bigscience-workshop/promptsource>

856
857
858
859
860
861
862
863
864

Evaluation Benchmarks. We evaluate the trained model across a diverse set of tasks to comprehensively assess its capabilities in natural language understanding, commonsense reasoning, text generation, instruction following, and text classification. These tasks span multiple languages and domains, enabling us to measure both general and multilingual performance. Below, we outline the benchmarks used for each evaluation category:

- 865
866
867
868
869
870
871
872
873
874
- **Natural Language Understanding:** We use the MMMLU benchmark, a multilingual extension of the widely adopted MMLU dataset (Hendrycks et al., 2020), designed for evaluating the multitask language understanding abilities of large language models. MMMLU covers 14 languages and includes questions from a wide range of domains. We report accuracy across all tasks to measure performance.
 - **Commonsense Reasoning:** We evaluate the model using the XCOPA dataset (Ponti et al., 2020). XCOPA tests a model’s ability to perform causal commonsense reasoning in multiple languages. The task involves selecting the most plausible cause or effect of a given premise from two alternatives, thus requiring both language understanding and reasoning skills.
 - **Text Generation:** We assess multilingual text generation performance using two benchmarks. First, FLORES-101 (Goyal et al., 2022) is used to evaluate the model’s general generation quality across a broad set of high- and low-resource languages. Second, FLORES-200 (Costa-Jussà et al., 2022) is employed to test zero-shot cross-lingual transfer capabilities—specifically, the model’s ability to generate high-quality outputs in languages it was not directly trained on.
 - **Instruction Following:** We use the multilingual variant of the IFEval benchmark (Zhou et al., 2023), implemented by the Benchmax framework (Huang et al., 2025), to evaluate the model’s ability to follow human instructions across diverse languages and task types. This benchmark focuses on the alignment between user instructions and model responses, which is critical for real-world applications of instruction-tuned models.

905
906
907
908
909
910
911
912

- **Text Classification:** For evaluating domain-robust classification performance, we adopt the SIB-200 benchmark (Adelani et al., 2024), which contains text classification tasks across 200 languages and multiple domains. This benchmark is particularly suited for testing the generalization and robustness of instruction-tuned models in a multilingual setting.

C Impact Factors: Additional Analysis

913
914
915
916
917
918
919
920
921
922

In Section 5, we investigated two key factors affecting the effectiveness of leveraging multi-way parallel data to enhance the multilingual capabilities of LLMs: the degree of parallelism and the presence of English in the language combinations. In this section, we extend our analysis to two additional factors: the size of the training data and the linguistic characteristics of the languages involved in the combinations.

C.1 Training Data Size

923
924
925
926
927
928
929
930
931
932
933
934

Figure 9 presents the impact of training data size on model performance. We conduct experiments by randomly sampling varying amounts of tokens—10K, 50K, 100K, 500K, 1M, 5M, 10M, 50M, 100M, 500M, and 1B—from the constructed TED2025 dataset. The results demonstrate that model performance is notably constrained when trained on smaller datasets (typically under 100K tokens). However, as the size of training data increases, performance consistently improves across all evaluated tasks.

935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951

For MMMLU and XCOPA, performance exhibits early gains with additional data but plateaus beyond a certain threshold. This trend likely reflects the nature of these tasks, which emphasize general language understanding and reasoning. Once the model acquires the necessary core linguistic and world knowledge, the marginal gains from further data diminish. Interestingly, performance on FLORES and xIFEval continues to improve with increasing data volume. These tasks, which involve cross-lingual understanding and translation—particularly for low-resource languages or semantically nuanced alignments—appear to benefit more substantially from large-scale data. This suggests that extensive training data is crucial for enhancing translation quality and evaluation accuracy in such settings.

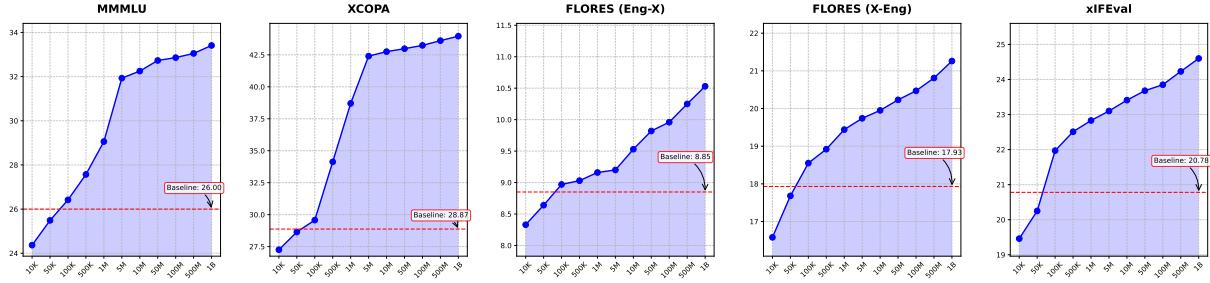


Figure 9: Impact of training data size on model performance across different token amounts (ranging from 10K to 1B) sampled from the TED2025 dataset.

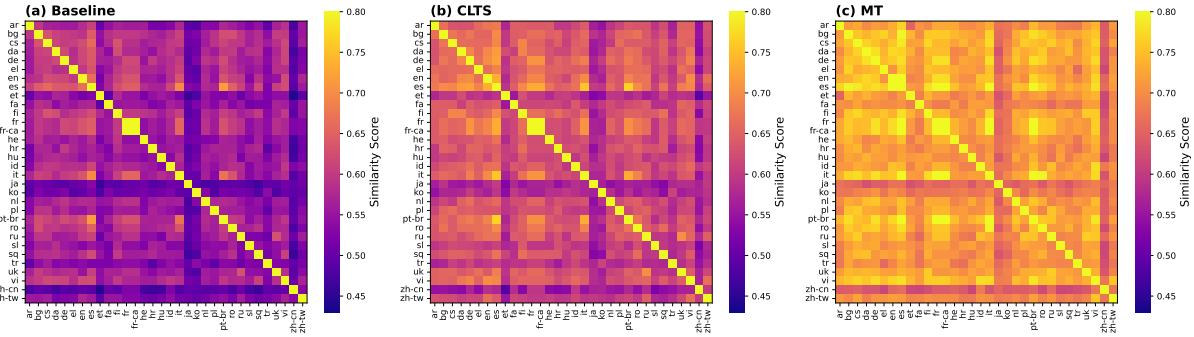


Figure 10: Representation alignment with instruction tuning across different training objectives.

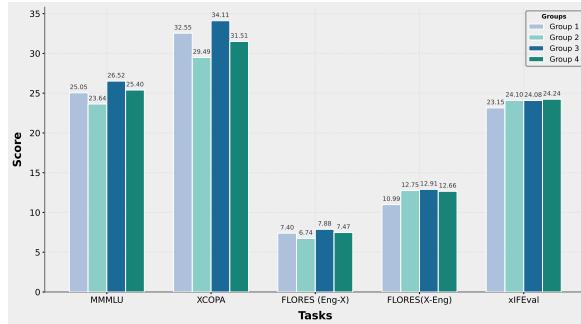


Figure 11: Impact of language family composition on model performance.

C.2 Language Combinations

In Section 5.2, we investigated the influence of including English in sampling combinations on model performance. In this section, we extend that analysis by investigating an alternative sampling strategy: whether the selected languages belong to the same language family. To this end, we construct four language groups. Groups 1 and 2 consist of languages from the same language family, whereas Groups 3 and 4 include languages from diverse families. The specific language configurations for each group are detailed in Table 8.

Figure 11 illustrates the impact of language family composition on model performance. The re-

sults indicate that sampling from cross-family language combinations more effectively leverages the benefits of multi-way parallel corpora, resulting in greater improvements in multilingual task performance. In contrast, sampling languages from the same family yields only marginal gains. This may be attributed to the structural and lexical similarities among related languages, which can introduce redundancy during training and limit the model’s ability to generalize across typologically diverse languages.

English Pivoting Combinations	
Group	Language List
group 1	en,vi,ar,bg,de,es,fr,he,it,ja,ko
group 2	en,nl,pl,pt-br,zh-cn,ar,bg,de,es,fr,he
group 3	en,el,vi,ar,bg,de,es,fr,he,it,ja
group 4	en,el,ar,bg,de,es,fr,he,it,ja,ko
group 5	en,fa,hu,ar,bg,de,es,fr,he,it,ja

Language Family Combinations		
Group	Language List	Language Family
group 1	bg,pl,ru,sr,uk	Slavic
group 2	el,it,kmr,qr,tr	Indo-European
group 3	en,ko,my,sq,zh-tw	Indo-European (en, sq), Sino-Tibetan (zh-tw), Koreanic (ko)
group 4	fr,hu,hy,lv,vi	Indo-European (fr, hy), Uralic (hu, lv), Austroasiatic (vi)

Table 8: Language configurations for different sampling groups.

977

978 **D Representation Alignment with**
Instruction Tuning

979 In this section, we explore how each tuning ob-
980 jective reshapes the model’s internal multilingual
981 embeddings, using the same alignment metrics as
982 in Section 4.3. Figure 10 shows that MT-tuned
983 models achieve the highest SVCCA score, indicat-
984 ing tighter alignment among semantically equiva-
985 lent sentences across languages. In contrast, CLTS
986 yields minimal alignment gains despite its similar-
987 ity focus, likely because binary similarity labels
988 lack the contextual richness of translation pairs.

Topic	#Talks	Topic	#Talks	Topic	#Talks	Topic	#Talks	Topic	#Talks
global issues	10334	race	154	coronavirus	70	inclusion	36	resources	18
education	9126	writing	154	sex	69	solar system	35	paleontology	17
technology	6323	physics	153	pandemic	69	sight	35	dinosaurs	17
art	5858	gender	144	product design	68	Audacious Project	35	television	17
business	5694	exploration	143	algorithm	68	industrial design	34	exercise	17
Life	5623	neuroscience	142	aging	67	sound	34	blockchain	17
health	5561	family	140	empathy	67	behavioral economics	34	fungi	16
science	4819	architecture	138	justice system	67	online privacy	34	public speaking	16
entertainment	4717	emotions	137	goals	66	travel	33	nuclear energy	16
design	2313	humor	136	urban planning	65	magic	33	PTSD	16
Humanities	1048	happiness	132	crime	65	trees	33	driverless cars	16
TED-Ed	941	universe	131	machine learning	65	weather	33	bees	15
animation	879	AI	131	compassion	64	public space	33	geology	15
culture	847	illness	131	investing	62	TED Prize	32	smell	15
social change	783	entrepreneur	130	fish	62	bioethics	31	Alzheimer's	15
TEDx	734	decision-making	130	demo	62	military	30	shopping	15
society	712	language	130	disability	62	illusion	30	Autism spectrum disorder	15
history	617	self	130	conservation	62	human rights	29	vulnerability	15
innovation	518	photography	127	feminism	61	theater	28	TED Connects	15
humanity	500	policy	126	Planets	60	solar energy	28	toys	14
biology	482	religion	124	renewable energy	57	biosphere	28	Antarctica	14
communication	453	media	121	code	57	maps	27	science fiction	14
future	434	democracy	120	women in business	57	spoken word	27	Islam	14
creativity	413	india	120	chemistry	57	heart	27	neurology	14
climate change	412	energy	118	Middle East	57	Slavery	27	Moon	13
community	407	motivation	115	immigration	57	sexual violence	27	3D printing	13
personal growth	395	philosophy	114	Best of the Web	56	surveillance	27	body language	13
environment	392	film	113	natural disaster	55	manufacturing	27	hearing	13
activism	371	violence	113	dance	55	trust	27	meditation	13
sustainability	341	literature	112	consciousness	54	archaeology	26	Brazil	12
performance	325	parenting	111	marine biology	54	AIDS	25	graphic design	12
psychology	325	journalism	111	virus	53	virtual reality	25	ebola	12
medicine	323	money	111	statistics	52	gardening	25	suicide	12
brain	323	potential	110	depression	52	nanotechnology	25	wind energy	12
music	316	ancient world	109	microbiology	52	Europe	25	coral reefs	12
work	316	social media	108	china	51	biomimicry	24	rivers	12
economics	308	love	107	electricity	51	drones	24	international relations	12
health care	307	poetry	107	plants	51	mindfulness	24	glaciers	12
nature	300	law	107	Vaccines	51	quantum	24	augmented reality	12
collaboration	288	pollution	106	Asia	49	aliens	24	worklife	12
politics	279	biodiversity	105	ethics	49	friendship	24	rocket science	11
animals	271	software	103	bacteria	48	encryption	24	Christianity	11
women	269	visualizations	103	farming	48	medical imaging	24	Sun	11
TED Fellows	267	teaching	100	prison	47	South America	23	bullying	11
human body	267	international development	99	refugees	47	telescopes	23	CRISPR	11
storytelling	259	finance	99	TED en Español	47	birds	23	String theory	10
invention	250	genetics	96	gaming	46	Big Bang	22	homelessness	10
identity	245	death	96	fear	46	Mission Blue	22	grammar	9
kids	243	books	94	natural resources	46	Surgery	22	typography	8
engineering	241	TED Residency	94	LGBTQIA+	46	protest	22	Buddhism	8
leadership	234	biotech	92	terrorism	44	painting	22	asteroid	8
government	233	beauty	92	philanthropy	44	interview	21	deextinction	8
equality	223	work-life balance	89	personality	44	library	21	cryptocurrency	8
public health	219	robots	88	marketing	43	plastic	21	metaverse	8
medical research	218	poverty	87	drugs	43	Mars	21	conducting	7
Internet	214	water	84	sociology	43	Egypt	21	bionics	7
computers	201	transportation	83	curiosity	43	pregnancy	21	NASA	7
Africa	197	cancer	83	TEDMED	43	synthetic biology	21	atheism	5
data	194	astronomy	82	indigenous peoples	42	Transgender	21	pain	5
cities	193	agriculture	82	sleep	41	prosthetics	20	forensics	4
disease	192	DNA	79	diversity	40	museums	20	microbes	4
space	187	success	79	fashion	40	addiction	20	NFTs	4
war	183	ecology	77	discovery	40	TED Membership	19	Judaism	3
United States	181	infrastructure	77	corruption	39	primates	18	street art	3
food	174	cognitive science	76	time	39	cyber security	18	Hinduism	1
math	173	sports	74	consumerism	38	astrobiology	18	reproductive health	1
ocean	168	youth	73	productivity	38	dark matter	18	crowdsourcing	1
Countdown	168	comedy	73	Anthropocene	38	botany	18	veganism	1
evolution	162	insects	71	capitalism	38	fossil fuels	18		
mental health	159	memory	70	flight	36	blindness	18		
relationships	154	anthropology	70	UX design	36	TED Books	18		

Table 9: Statistics on the number of talks per domain in the TED2025 dataset.

