

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

Yingli Shen<sup>1\*</sup>, Wen Lai<sup>2\*</sup>, Shuo Wang<sup>1</sup>, Ge Gao<sup>3</sup>  
Kangyang Luo<sup>1</sup>, Alexander Fraser<sup>2</sup>, Maosong Sun<sup>1</sup>

<sup>1</sup>Tsinghua University, <sup>2</sup>Technical University of Munich, <sup>3</sup>Minzu University of China



# 1 Introduction

## 2 TED2025

## 3 Experimental Setup

## 4 RQ1: Multi-Way vs. Unaligned

## 5 RQ2: Impact Factors

## 6 RQ3: Instruction Tuning

## 7 Summary

# Background

- LLMs excel in high-resource languages but lag in low-resource languages.
- Approach to improve performance on low-resource languages:
  - Continue Pretraining / Instruction Tuning using more multilingual data.
- Two types of multilingual data:
  - Unaligned multilingual data
  - Aligned (parallel / multi-way parallel) data

Multi-way parallel data is underexplored for scaling multilingual LLMs.



# Research Goal & Questions

- **Research Goal:**
  - Construct TED2025 (large-scale multi-way corpus) and explore best practices.
- **Research Question:**
  - **RQ1:** Multi-way vs. unaligned → downstream task & cross-lingual transfer & representation alignment.
  - **RQ2:** Best practice on continued pretraining → parallelism & English as Pivot & Language Combinations & Training Data Size.
  - **RQ3:** Best practice on instruction tuning → different objectives.

## ① Introduction

## ② TED2025

## ③ Experimental Setup

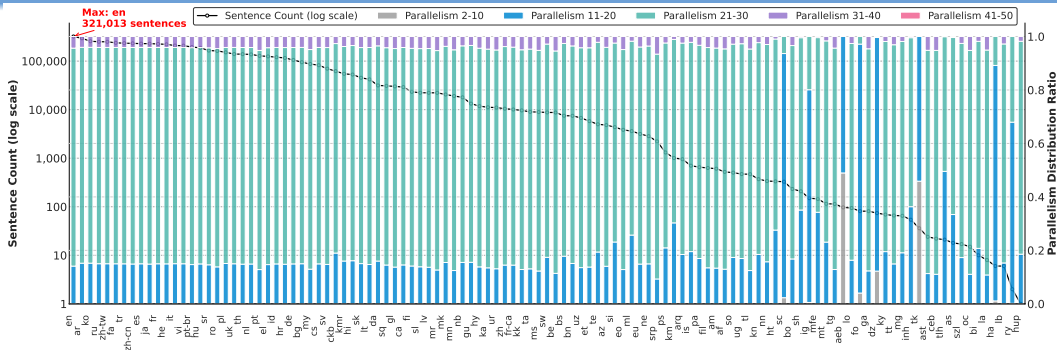
## ④ RQ1: Multi-Way vs. Unaligned

## ⑤ RQ2: Impact Factors

## ⑥ RQ3: Instruction Tuning

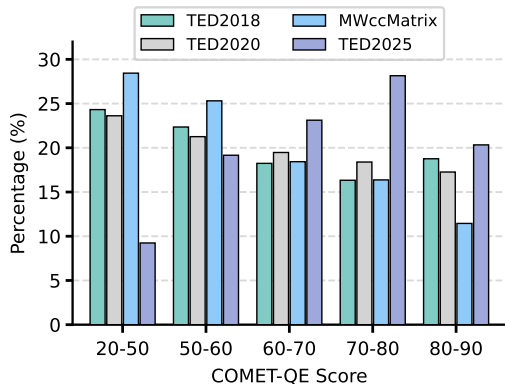
## ⑦ Summary

# Parallelism Distribution



- Contains 113 languages with up to 50-way parallelism.
- Most common parallelism: 21–30 languages.
- Low-Resource languages also achieve high parallelism.

# Data Quality



- Existing Dataset:

- TED2018 (Qi et al., 2018)
- TED2020 (Reimers and Gurevych, 2020)
- MWccMatrix (Thompson et al., 2024)

- Evaluation Metric:

- COMET-QE (Rei et al., 2020)

- TED2025 has significantly more high-quality translations (COMET-QE > 60).
- COMET-QE score for all language pairs can be found in Appendix.



## 1 Introduction

## 2 TED2025

## 3 Experimental Setup

## 4 RQ1: Multi-Way vs. Unaligned

## 5 RQ2: Impact Factors

## 6 RQ3: Instruction Tuning

## 7 Summary

# Setup

- **Datasets:**

- Multi-Way aligned dataset: TED2025
- Unaligned dataset: DCAD-2000 (Shen et al., 2025<sup>a</sup>)

- **Models:**

- LLaMA-3.1-8B and Qwen-2.5-14B

- **Training:**

- Continued pretraining (RQ1 & RQ2) / Instruction tuning (RQ3)
- LoRA instead Full Parameter Training

<sup>a</sup>DCAD-2000: A Multilingual Dataset across 2000+ Languages with Data Cleaning as Anomaly Detection (Shen et al., NeurIPS 2025)

- **Evaluation Benchmarks**

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

- ① Introduction
- ② TED2025
- ③ Experimental Setup
- ④ RQ1: Multi-Way vs. Unaligned
- ⑤ RQ2: Impact Factors
- ⑥ RQ3: Instruction Tuning
- ⑦ Summary

# Downstream Performance

	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

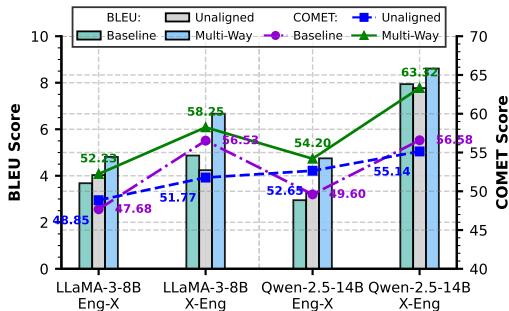
(a) LLaMA-3-8B

	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

- Multi-Way > Unaligned > Baseline (low- & high-resource).
- non-English-centric tasks are provided in Appendix C.

# Zero-Shot Transfer



- Evaluation:

- Machine Translation on FLORES-200.
- We exclude all languages in the evaluation subset from training and assess the English-X translation quality.

- Multi-Way achieves stronger zero-shot transfer on FLORES-200.

# Representation Alignment

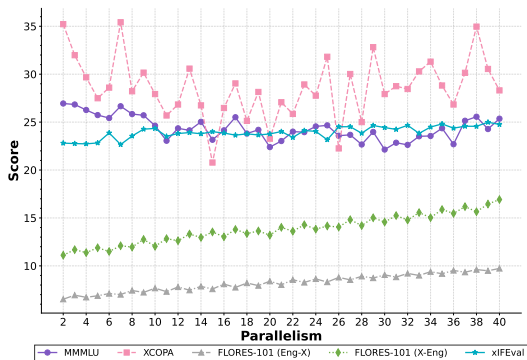
- We randomly select 32 aligned languages (with 100 sentences each) from TED2025
- Four similarity evaluation metric on embedding space:
  - Cosine Similarity
  - Centered Kernel Alignment (CKA; Kornblith et al., 2019)
  - Cross-Lingual Sentence Retrieval (Retrieval; Conneau et al., 2017)
  - Singular Vector Canonical Correlation Analysis (SVCCA; Raghu et al., 2017)

		LLaMA-3-8B			Qwen-2.5-14B		
		Baseline	Unaligned	Multi-Way	Baseline	Unaligned	Multi-Way
Cosine (↑)		0.27	0.27 <sub>-0.00</sub>	<b>0.30</b> <sub>+0.03</sub>	0.29	0.27 <sub>-0.02</sub>	<b>0.32</b> <sub>+0.03</sub>
CKA (↑)		0.54	0.54 <sub>-0.00</sub>	<b>0.60</b> <sub>+0.06</sub>	0.56	0.57 <sub>+0.01</sub>	<b>0.63</b> <sub>+0.07</sub>
Retrieval (↑)	P@1	0.09	0.07 <sub>-0.02</sub>	<b>0.13</b> <sub>+0.04</sub>	0.12	0.14 <sub>+0.02</sub>	<b>0.19</b> <sub>+0.07</sub>
	P@5	0.24	0.23 <sub>-0.01</sub>	<b>0.27</b> <sub>+0.03</sub>	0.26	0.28 <sub>+0.02</sub>	<b>0.33</b> <sub>+0.07</sub>
	P@10	0.35	0.33 <sub>-0.02</sub>	<b>0.39</b> <sub>+0.04</sub>	0.36	0.38 <sub>+0.02</sub>	<b>0.42</b> <sub>+0.06</sub>
SVCCA (↑)		0.55	0.55 <sub>-0.00</sub>	<b>0.61</b> <sub>+0.06</sub>	0.57	0.58 <sub>+0.01</sub>	<b>0.63</b> <sub>+0.06</sub>

- Multi-Way shows denser, language-agnostic embeddings.

- ① Introduction
- ② TED2025
- ③ Experimental Setup
- ④ RQ1: Multi-Way vs. Unaligned
- ⑤ RQ2: Impact Factors
- ⑥ RQ3: Instruction Tuning
- ⑦ Summary

# Degree of Parallelism



- Experiment Setups
  - We fix the tokens to 5M in each parallelism.
- Evaluation

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-200	Text Classification	204	Acc

- MT benefits from higher parallelism.
- Reasoning tasks peak at 6–10 languages.



# English as Pivot

		MMMLU	XCOPA	FLORES (Eng-X)	FLORES (X-Eng)	xIFEval
Group 1	with	<b>23.17</b>	<b>30.93</b>	5.80	13.75	23.39
	w/o	20.84	30.07	<b>6.63</b>	<b>14.56</b>	<b>24.14</b>
Group 2	with	<b>22.19</b>	<b>35.33</b>	6.74	13.88	<b>22.19</b>
	w/o	18.51	31.60	<b>7.13</b>	<b>14.18</b>	<b>22.19</b>
Group 3	with	<b>26.47</b>	<b>36.73</b>	<b>6.40</b>	13.48	22.69
	w/o	23.83	35.78	6.38	<b>14.00</b>	<b>22.91</b>
Group 4	with	<b>23.42</b>	<b>33.84</b>	6.20	12.11	23.66
	w/o	20.26	30.84	<b>6.67</b>	<b>14.76</b>	<b>23.67</b>
Group 5	with	<b>23.89</b>	<b>39.64</b>	6.96	13.07	23.19
	w/o	20.39	34.15	<b>7.83</b>	<b>14.79</b>	<b>23.85</b>

- Five groups, each with both English included and English-excluded variants across ten 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

- English helps reasoning/understanding.
- English slightly hurts direct MT transfer.

## Other Factors

- **Language Combinations**

- Investigating an alternative sampling strategy: whether the selected languages belong to the same language family.

- **Training Data Size**

- 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.

- More details can be found in Section 5.3 and 5.4.

- 1 Introduction
- 2 TED2025
- 3 Experimental Setup
- 4 RQ1: Multi-Way vs. Unaligned
- 5 RQ2: Impact Factors
- 6 RQ3: Instruction Tuning
- 7 Summary

# Objectives

- Four training objectives

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}.

# Instruction Tuning Objectives

	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	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	<b>45.28</b>	<b>51.06</b>	<b>68.17</b>	<b>69.38</b>	<b>11.26</b>	<b>13.25</b>	<b>63.03</b>	<b>65.61</b>	<b>22.25</b>	<b>21.60</b>	<b>70.76</b>	<b>75.41</b>	<b>41.72</b>	<b>44.52</b>
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	<b>68.77</b>	6.19	10.51	<b>58.53</b>	60.33	16.83	18.98	68.70	<b>74.91</b>	36.47	<b>40.30</b>
MT + MTC	42.69	47.54	65.12	66.83	6.35	10.73	57.53	60.24	15.13	18.45	<b>69.28</b>	74.00	36.10	39.03
MT + CLP	43.20	<b>49.07</b>	<b>67.29</b>	68.47	<b>7.31</b>	10.51	58.40	<b>62.64</b>	<b>17.75</b>	<b>20.77</b>	69.18	74.47	<b>38.49</b>	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	<b>42.82</b>	46.74	65.30	67.12	5.58	<b>10.88</b>	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	<b>42.72</b>	<b>47.67</b>	64.83	67.57	<b>6.17</b>	10.94	57.78	60.11	<b>15.96</b>	<b>19.16</b>	68.67	<b>74.03</b>	<b>36.54</b>	38.77
MT + MTC + CLP	42.33	46.72	<b>65.26</b>	<b>67.58</b>	5.92	<b>10.98</b>	<b>57.93</b>	<b>60.76</b>	15.38	18.81	<b>68.98</b>	73.55	36.22	<b>39.15</b>
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

- The improvements in MT are the largest and most stable in both high-resource and low-resource languages.

# Instruction Tuning Objectives

	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	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	<b>45.28</b>	<b>51.06</b>	<b>68.17</b>	<b>69.38</b>	<b>11.26</b>	<b>13.25</b>	<b>63.03</b>	<b>65.61</b>	<b>22.25</b>	<b>21.60</b>	<b>70.76</b>	<b>75.41</b>	<b>41.72</b>	<b>44.52</b>
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	<b>68.77</b>	6.19	10.51	<b>58.53</b>	60.33	16.83	18.98	68.70	<b>74.91</b>	36.47	<b>40.30</b>
MT + MTC	42.69	47.54	65.12	66.83	6.35	10.73	57.53	60.24	15.13	18.45	<b>69.28</b>	74.00	36.10	39.03
MT + CLP	43.20	<b>49.07</b>	<b>67.29</b>	68.47	<b>7.31</b>	10.51	58.40	<b>62.64</b>	<b>17.75</b>	<b>20.77</b>	69.18	74.47	<b>38.49</b>	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	<b>42.82</b>	46.74	65.30	67.12	5.58	<b>10.88</b>	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	<b>42.72</b>	<b>47.67</b>	64.83	67.57	<b>6.17</b>	10.94	57.78	60.11	<b>15.96</b>	<b>19.16</b>	68.67	<b>74.03</b>	<b>36.54</b>	38.77
MT + MTC + CLP	42.33	46.72	<b>65.26</b>	<b>67.58</b>	5.92	<b>10.98</b>	<b>57.93</b>	<b>60.76</b>	15.38	18.81	<b>68.98</b>	73.55	36.22	<b>39.15</b>
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

- Interestingly, the combination of tasks did not significantly improve performance.

## 1 Introduction

## 2 TED2025

## 3 Experimental Setup

## 4 RQ1: Multi-Way vs. Unaligned

## 5 RQ2: Impact Factors

## 6 RQ3: Instruction Tuning

## 7 Summary

# Conclusion

- We introduce TED2025, a large-scale, high-quality multi-way parallel dataset covering 113 languages, with a maximum parallel degree of 50.
- Using TED2025, we investigate:
  - The effectiveness of using multi-way parallel data to improve the multilingual capabilities.
  - The best practice (impact factors) on using multi-way parallel data.
  - The best practice on using multi-way parallel data on instruction tuning.



Thank you!

Questions & Comments?



Paper



Code

**Contact:** wen.lai@tum.de    syl@mail.tsinghua.edu.cn