

Multilingual Machine Translation with Large Language Models: Empirical Results and Analysis

Anonymous ACL submission

Abstract

Large language models (LLMs) have demonstrated remarkable potential in handling multilingual machine translation (MMT). In this paper, we systematically investigate the advantages and challenges of LLMs for MMT by answering two questions: 1) How well do LLMs perform in translating massive languages? 2) Which factors affect LLMs’ performance in translation? We thoroughly evaluate eight popular LLMs, including ChatGPT and GPT-4. Our empirical results show that translation capabilities of LLMs are continually involving. GPT-4 has beat the strong supervised baseline NLLB in 40.91% of translation directions but still faces a large gap towards the commercial translation system like Google Translate, especially on low-resource languages. Through further analysis, we discover that LLMs exhibit new working patterns when used for MMT. First, LLM can acquire translation ability in a resource-efficient way and generate moderate translation even on zero-resource languages. Second, instruction semantics can surprisingly be ignored when given in-context exemplars. Third, cross-lingual exemplars can provide better task guidance for low-resource translation than exemplars in the same language pairs.

1 Introduction

With the increasing scale of parameters and training corpus, large language models (LLMs) have gained a universal ability to handle a variety of tasks via in-context learning (ICL, Brown et al. 2020), which allows language models to perform tasks with a few given exemplars and human-written instructions as context. One particular area where LLMs have shown outstanding potential is machine translation (MT). Previous studies have shown the surprising performance of LLMs on high-resource bilingual translation, such as English-German translation (Vilar et al., 2022; Zhang et al., 2022), even if these models are not particularly optimized on multilingual data.

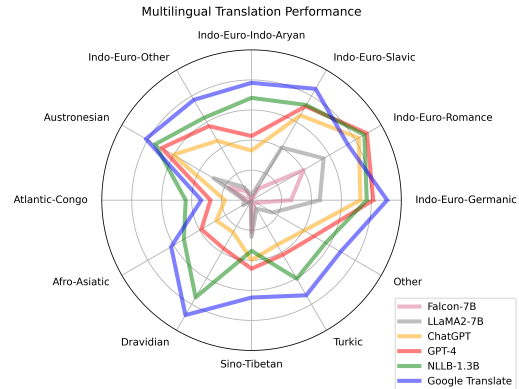


Figure 1: Multilingual translation performance (translating from English to non-English) of some popular LLMs and traditional supervised systems. LLMs have demonstrated great potential in multilingual machine translation.

However, the multilingual translation ability of LLMs remains under-explored. MMT is a challenging task that involves translating text among different languages and requires semantic alignment between languages (Fan et al., 2021; Costa-jussà et al., 2022; Yuan et al., 2023). It is also unclear that how LLM acquires translation ability and which factors affect LLM’s translation ability.

In this paper, we follow ICL paradigm and focus on studying LLMs in multilingual machine translation by answering two questions: 1) *How LLMs perform MMT over massive languages?* 2) *Which factors affect the performance of LLMs?*

For the first question, we evaluate several popular LLMs: English-centric LLMs, including OPT (Zhang et al., 2022), LLaMA2 (Touvron et al., 2023), Falcon (Almazrouei et al., 2023) and multilingual LLMs, including XGLM (Lin et al., 2022), BLOOMZ (Scao et al., 2022), ChatGPT (OpenAI, 2022), GPT-4 (OpenAI, 2023). We consider 102 languages and 606 translation directions (202 English-centric directions, 202 French-centric directions and 202 Chinese-centric direc-

tions). Results show that the multilingual translation capabilities of LLMs are continually improving and GPT-4 reaches new performance height. Compared with the widely-used supervised MMT system NLLB (Costa-jussà et al., 2022), GPT-4 achieves higher performance on 40.91% English-centric translation directions. But compared with the commercial translation system (Google Translate), LLMs still have a long way to go, particularly when it comes to low-resource languages. French-centric and Chinese-centric translation are also more challenging for GPT-4 than English-centric translation, which further indicates its unbalanced capability across languages.

For the second question, we find some new working patterns. First, we discover that LLM can acquire translation ability in a resource-efficient way and generate moderate translation even on zero-resource languages. Second, LLMs are able to perform translation even with unreasonable instructions if in-context learning exemplars are given. However, if given mismatched translation pairs as in-context exemplars, LLMs fail to translate, which is similar to observations from concurrent studies (Wei et al., 2023). This shows the importance of exemplars in ICL for machine translation. Third, we find that cross-lingual translation pairs can be surprisingly good exemplars for low-resource translation, even better than exemplars in the same language.

The main contribution of this paper can be summarized below:

- We benchmark popular LLMs on MMT in 102 languages and 606 translation directions, covering English-centric, French-centric and Chinese-centric translation.
- We systematically compare the results of LLMs and three strong supervised baselines (M2M-100, NLLB, Google Translator) and reveal the gap between two translation paradigms.
- We find some new ICL working patterns of LLMs for MMT and discuss corresponding advantages and challenges.

2 Background

2.1 Large Language Models

Language modeling is a long-standing task in natural language processing (Bengio et al., 2000; Mikolov et al., 2010; Khandelwal et al., 2020),

which is a task to predict the probability of the next token. Transformer (Vaswani et al., 2017) basically is the backbone of existing LLMs.

LLMs show great potential as a universal multi-task learner. Recently, Radford et al. (2019) find that a casual decoder-only language model can be a multi-task learner with merely unsupervised training corpus. Later, Kaplan et al. (2020) reveal the *scaling law* of LLM, indicating that when the scale of neural parameters and training data keeps increasing, LLM can be further strengthened. Wei et al. (2022b) show that scaling the language model also brings astonishing *emergent abilities*, e.g., in-context learning, which is only present in large models. Consequently, more and more efforts have been put into scaling-up language models (Brown et al., 2020; Hoffmann et al., 2022; Scao et al., 2022; Vilar et al., 2022; Ren et al., 2023). Among them, GPT-4 (OpenAI, 2023) and ChatGPT (OpenAI, 2022) are the most representative systems, which shows impressive results in various NLP tasks.

2.2 Emergent Ability: In-context Learning

In-context learning is one of the well-known emergent abilities (Brown et al., 2020; Dong et al., 2022), which enables LLM to learn target tasks according to the prompt without updating any parameters.

Specifically, the prompt is made up of in-context exemplars $\{(\mathcal{X}_i, \mathcal{Y}_i)\}_{i=1}^k$ and in-context template \mathcal{T} . Exemplars are often picked from supervised data, where \mathcal{Y}_i is the ground truth corresponding to the input sentence \mathcal{X}_i . Template \mathcal{T} is usually a human-written instruction related to the target task. Wrapping exemplars with the template and concatenating them together produce the final prompt:

$$\mathcal{P} = \mathcal{T}(\mathcal{X}_1, \mathcal{Y}_1) \oplus \mathcal{T}(\mathcal{X}_2, \mathcal{Y}_2) \oplus \cdots \oplus \mathcal{T}(\mathcal{X}_k, \mathcal{Y}_k)$$

where \oplus denotes the concatenation symbol, e.g., whitespace, line-break. During inference, LLM is able to generate the corresponding output \mathcal{Y} of the test sample \mathcal{X} under the guidance of the prompt:

$$\arg \max_{\mathcal{Y}} p(\mathcal{P} \oplus \mathcal{T}(\mathcal{X}, \mathcal{Y})) \quad (1)$$

For label prediction tasks, the prediction \mathcal{Y} can be obtained in one-step generation. For sequence generation tasks, e.g., machine translation, the prediction \mathcal{Y} can be obtained through sampling strategies like greedy search and beam search.

Language Family	Direction	Translation Performance (BLEU / COMET)									
		XGLM-7.5B	OPT-175B	Falcon-7B	LLaMA2-7B	LLaMA2-7B-Chat	ChatGPT	GPT-4	M2M-12B	NLLB-1.3B	Google
Indo-Euro-Germanic (8)	X⇒Eng	18.54 / 70.09	34.65 / 83.71	27.37 / 67.40	37.28 / 84.73	34.82 / 84.25	45.83 / 89.05	<u>48.51 / 89.48</u>	42.72 / 87.74	46.54 / 88.18	51.16 / 89.36
	Eng⇒X	9.16 / 50.21	18.89 / 71.97	13.19 / 52.93	22.78 / 76.05	19.44 / 73.63	36.34 / 87.83	<u>40.64 / 88.50</u>	37.30 / 86.47	38.47 / 87.31	45.27 / 89.05
Indo-Euro-Romance (8)	X⇒Eng	31.11 / 79.67	38.93 / 87.75	34.06 / 84.40	41.10 / 88.10	37.84 / 87.80	45.68 / 89.61	47.29 / 89.74	42.33 / 88.31	46.33 / 88.99	35.69 / 89.66
	Eng⇒X	21.95 / 69.08	24.30 / 79.07	20.02 / 70.36	27.81 / 82.05	25.50 / 79.67	41.35 / 89.00	44.47 / 88.94	42.98 / 87.56	43.48 / 88.12	37.10 / 88.77
Indo-Euro-Slavic (12)	X⇒Eng	13.20 / 64.24	20.83 / 74.80	13.15 / 57.34	20.24 / 76.30	30.94 / 83.90	39.27 / 87.74	<u>41.19 / 88.15</u>	35.87 / 85.97	39.23 / 87.08	43.61 / 88.18
	Eng⇒X	6.40 / 43.28	8.18 / 54.45	4.34 / 35.73	5.00 / 44.09	16.14 / 69.75	32.61 / 87.90	<u>36.06 / 89.15</u>	35.01 / 86.43	36.56 / 88.74	42.75 / 90.05
Indo-Euro-Indo-Aryan (10)	X⇒Eng	8.68 / 63.93	1.20 / 49.37	1.40 / 45.22	6.68 / 62.63	4.29 / 60.29	25.32 / 84.14	<u>37.30 / 87.79</u>	17.53 / 69.66	40.75 / 88.80	45.66 / 89.43
	Eng⇒X	4.76 / 40.99	0.14 / 31.85	0.13 / 25.84	1.61 / 35.92	1.24 / 34.74	16.50 / 68.43	<u>21.35 / 73.75</u>	14.44 / 65.32	34.04 / 82.55	39.04 / 82.78
Indo-Euro-Other (11)	X⇒Eng	7.32 / 55.29	7.80 / 59.60	7.04 / 51.59	14.27 / 69.87	11.46 / 67.64	29.54 / 84.52	<u>37.29 / 86.76</u>	22.38 / 77.47	36.16 / 86.81	41.68 / 88.29
	Eng⇒X	4.51 / 40.60	3.10 / 40.04	3.38 / 34.64	5.00 / 44.09	4.83 / 43.73	22.81 / 77.33	<u>28.45 / 80.94</u>	19.71 / 74.90	31.65 / 85.82	38.54 / 87.44
Austronesian (6)	X⇒Eng	16.19 / 78.80	25.60 / 78.03	18.62 / 75.36	26.70 / 80.21	24.39 / 80.39	39.95 / 87.29	<u>46.81 / 88.65</u>	31.84 / 84.76	45.41 / 87.85	50.68 / 88.89
	Eng⇒X	10.01 / 73.14	10.68 / 64.97	8.56 / 60.89	14.59 / 74.80	13.29 / 74.88	30.17 / 86.36	<u>34.66 / 87.68</u>	27.03 / 86.83	37.17 / 88.82	40.74 / 89.34
Atlantic-Congo (14)	X⇒Eng	6.67 / 62.00	9.17 / 57.59	6.98 / 0.56	8.76 / 57.72	9.01 / 57.86	19.86 / 79.63	<u>28.27 / 83.42</u>	10.55 / 76.43	32.20 / 84.00	23.55 / 85.44
	Eng⇒X	2.52 / 54.93	1.60 / 34.15	1.89 / 0.34	2.45 / 34.17	3.09 / 38.13	8.91 / 75.26	<u>13.70 / 77.79</u>	6.53 / 75.79	21.99 / 79.95	16.77 / 80.89
Afro-Asiatic (6)	X⇒Eng	6.70 / 54.51	5.93 / 52.90	4.87 / 38.62	10.41 / 57.72	8.65 / 58.27	20.84 / 70.39	<u>30.48 / 78.76</u>	10.00 / 66.98	32.69 / 82.99	36.14 / 84.47
	Eng⇒X	2.07 / 41.48	1.40 / 41.86	1.40 / 27.64	3.22 / 43.04	3.07 / 43.39	13.57 / 67.60	<u>19.36 / 75.56</u>	7.83 / 68.86	26.08 / 82.84	31.00 / 83.78
Turkic (5)	X⇒Eng	7.43 / 61.69	7.89 / 62.47	4.15 / 33.11	9.51 / 65.95	8.88 / 66.15	24.64 / 84.04	<u>31.73 / 86.90</u>	10.25 / 58.52	32.92 / 87.51	37.78 / 88.53
	Eng⇒X	3.48 / 40.32	2.58 / 44.80	1.75 / 20.00	3.28 / 39.65	3.09 / 41.97	17.13 / 74.77	<u>20.96 / 78.50</u>	10.87 / 68.21	30.17 / 88.47	36.54 / 89.38
Dravidian (4)	X⇒Eng	8.04 / 61.95	0.89 / 44.01	1.18 / 24.29	2.65 / 53.17	1.52 / 52.95	20.26 / 82.00	<u>33.10 / 86.91</u>	10.26 / 63.77	39.07 / 88.42	43.17 / 89.10
	Eng⇒X	5.30 / 48.15	0.02 / 32.51	0.03 / 15.31	0.56 / 34.03	0.58 / 35.65	12.34 / 64.74	<u>18.60 / 75.15</u>	6.85 / 62.25	37.33 / 86.32	44.16 / 87.75
Sino-Tibetan (3)	X⇒Eng	9.35 / 58.60	9.32 / 65.32	16.59 / 72.34	18.35 / 74.45	16.88 / 74.20	21.36 / 78.52	<u>27.74 / 84.48</u>	11.09 / 71.35	30.88 / 86.50	35.68 / 87.66
	Eng⇒X	10.14 / 74.16	2.57 / 54.73	10.74 / 66.74	12.24 / 65.99	9.06 / 65.07	19.92 / 76.04	<u>22.81 / 81.11</u>	10.42 / 73.82	16.85 / 80.74	32.40 / 88.52
Other (14)	X⇒Eng	9.71 / 60.43	10.10 / 60.78	5.37 / 47.38	16.00 / 71.15	14.25 / 70.35	25.59 / 82.48	<u>32.62 / 86.21</u>	25.53 / 81.53	35.06 / 86.86	36.95 / 87.93
	Eng⇒X	8.42 / 51.57	3.82 / 46.85	1.73 / 29.73	8.19 / 53.20	7.14 / 52.12	20.26 / 74.31	<u>24.04 / 79.59</u>	23.29 / 77.80	28.54 / 85.84	34.34 / 87.82

Table 1: Average translation performance of LLMs on different language families. The number in the bracket indicates the number of evaluated languages in the specific language family. Bold text denotes the highest BLEU or COMET score across models. Underlined text denotes the highest BLEU or COMET score across LLMs.

3 Experiment Setup

Dataset We benchmark multilingual translation on FLORES-101 (Goyal et al., 2022) dataset¹, which enables an assessment of model quality on a wide range of languages.

LLMs We evaluate translation performance of eight popular LLMs: XGLM-7.5B (Lin et al., 2022), OPT-175B (Zhang et al., 2022), BLOOMZ-7.1B (Scao et al., 2022), Falcon-7B (Almazrouei et al., 2023), LLaMA2-7B (Touvron et al., 2023), LLaMA2-7B-chat (Touvron et al., 2023), ChatGPT (OpenAI, 2022) and GPT-4 (OpenAI, 2023).

ICL strategy For each model, we report its translation performance with eight randomly-picked translation pairs from the corresponding development set as in-context exemplars and “<X>=<Y>” as in-context template. “<X>” and “<Y>” are the placeholder for the source and target sentence. We use line-break as the concatenation symbol. According to our experiment analysis, this ICL strategy serves as a simple but strong recipe. All implementation is based on *OpenICL*² (Wu et al., 2023).

¹We evaluate LLMs on the first 100 sentences of each direction’s test set in benchmarking experiment, considering the prohibitive API cost of evaluating massive languages. In analysis experiment, we use full test set.

²<https://github.com/Shark-NLP/OpenICL>

Supervised baselines We report the performance of the supervised model M2M-100-12B (Fan et al., 2021) and NLLB-1.3B (Costa-jussà et al., 2022) (distillation version), which are widely-used many-to-many MMT models. We also report the performance of the powerful commercial translation system, Google Translate³.

Metric Following Goyal et al. (2022), we use SentencePiece BLEU⁴ (spBLEU) as evaluation metric, which enables an evaluation of all languages. In addition, we also consider emerging metrics, COMET⁵ (Rei et al., 2020) and SEScore⁶ (Xu et al., 2022b), which have been shown to correlate well with human judgements.

4 Benchmarking LLMs for Massively Multilingual Machine Translation

In this section, we report results on multilingual machine translation and introduce our main findings about LLMs’ translation ability.

The multilingual translation capabilities of LLMs are continually involving Table 1 presents evaluation results⁷ grouped by language

³<https://translate.google.com/>

⁴<https://github.com/mjpost/sacrebleu>

⁵We compute the score with *wmt22-comet-da* model.

⁶We compute the score with *SEScore-2* (Xu et al., 2022a).

⁷Evaluating with SEScore leads to similar findings, thus we report those results in Appendix A. Detailed results for

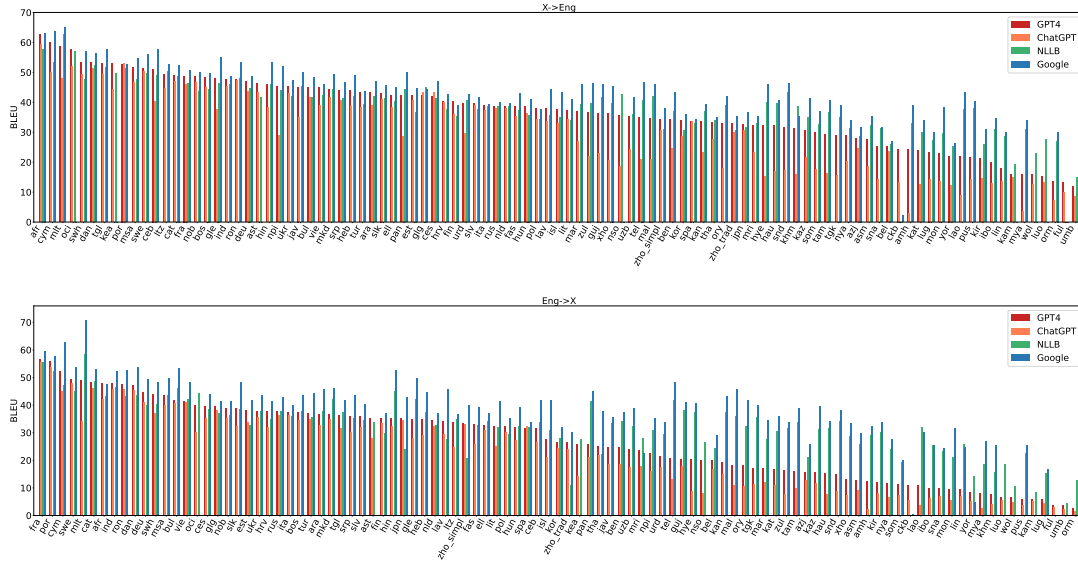


Figure 2: Translation performance (BLEU) of GPT-4, ChatGPT, NLLB and Google Translate on our evaluated languages. “X->Eng” and “Eng->X” denote translating to English and translating from English respectively. In each subfigure, languages are sorted according to BLEU scores of GPT-4.

family. Monolingual pre-trained LLMs present impressive multilingual translation ability, indicating the possibility of aligning multiple languages even with unsupervised data (Garcia et al., 2023). More encouragingly, the multilingual translation capabilities of LLMs are continually improving. The most recent LLMs are reaching new performance heights; for example, LLaMA2-7B outperforms previously released open-source LLMs, and GPT-4 surpasses ChatGPT. Overall, GPT-4 is the best translator among evaluated LLMs and it achieves the highest average BLEU and COMET score on most directions.

Language Family	X→Eng	X→Fra	X→Zho	Eng→X	Fra→X	Zho→X
Indo-Euro-Germanic (8)	48.51	44.23	27.97	40.64	32.34	24.13
Indo-Euro-Romance (8)	47.29	45.16	27.31	44.47	36.05	27.12
Indo-Euro-Slavic (12)	41.19	40.32	25.67	36.06	30.88	23.33
Indo-Euro-Indo-Aryan (10)	37.30	32.81	21.81	21.35	17.26	13.55
Indo-Euro-Other (11)	37.29	35.36	22.70	28.45	22.57	17.50
Austronesian (6)	46.81	39.98	24.40	34.66	25.64	19.52
Atlantic-Congo (14)	28.27	25.02	15.72	13.70	10.42	7.60
Afro-Asiatic (6)	30.48	27.00	17.81	19.36	14.43	10.53
Turkic (5)	31.73	30.90	19.96	20.96	17.80	14.02
Dravidian (4)	33.10	30.61	20.63	18.60	14.47	11.37
Sino-Tibetan (3)	27.74	27.93	20.88	22.81	19.21	16.30
Other (14)	32.62	31.26	21.25	24.04	20.03	16.37

Table 2: Translation performance (BLEU) of GPT-4 on English-centric, French-centric and Chinese-centric translation.

each translation direction are listed in Appendix B.

LLM’s capability is unbalanced across languages In Table 1, we observe a similar trend for all evaluated LLMs: they perform better at translating into English than translating into non-English. LLM’s capability on non-English languages is also unbalanced. For languages that are similar to English, e.g, Indo-European-Germanic languages, LLMs achieve impressive results. For languages that are dissimilar to English, e.g., Sino-Tibetan languages, LLMs often produce less decent results.

Table 2 presents another clue, where we evaluate GPT-4 on French-centric and Chinese-centric translation. Compared to English-centric translation, GPT-4 faces greater challenge when it comes to non-English-centric translation, which again indicates LLM’s unbalanced translation ability across languages.

LLMs still lag behind the strong supervised baseline, especially on low-resource languages

Figure 2 shows the translation performance of the supervised systems and GPT-4 on each language. In 40.91% translation directions, GPT-4 has achieved higher BLEU scores than NLLB, indicating the promising future of this new translation paradigm. But on long-tail low-resource languages, GPT-4 still lags behind NLLB, let alone Google Translate.

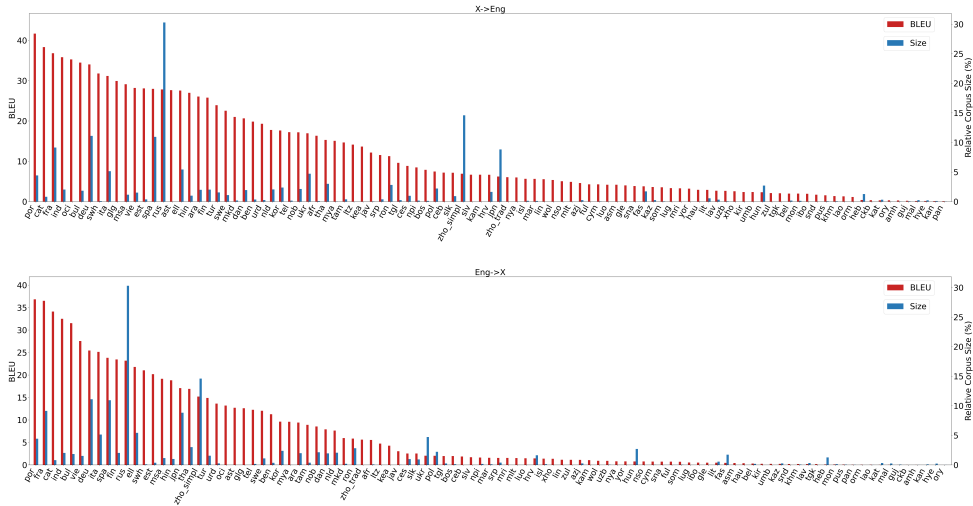


Figure 3: Translation performance (BLEU) of XGLM on evaluated languages and the corpus size of each language relative to English pre-training corpus. In each subfigure, languages are sorted according to BLEU scores of XGLM.

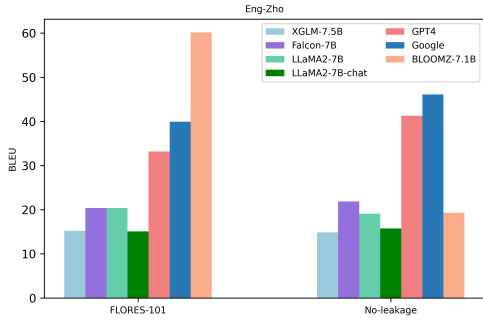


Figure 4: Translation performance of different models on FLORES-101 test set and our annotated no-leakage evaluation set NEWS2023.

Data leakage issue should be considered before evaluating LLMs on public datasets. We do not include BLOOMZ’s performance on FLORES-101 in our report because BLOOMZ is instruction-tuned with XP3 dataset (Scao et al., 2022), which includes FLORES-200 dataset. Thus BLOOMZ may have been exposed to test cases from FLORES-101 during training. If so, the evaluation results can not precisely reflect its translation ability (Elan-govan et al., 2021).

To illustrate this concern, we take 1000 English sentences from the most recent news spanning August 2023 to October 2023⁸, and ask human experts to translate them into Chinese and construct a bilingual no-leakage evaluation set, named NEWS2023. Figure 4 shows that BLOOMZ’s performance significantly deteriorates on this no leakage set, whereas other models maintain a consistent

⁸The news were collected from BBC news, Fox news, ABC news and Yahoo news.

performance across both datasets. Through this example, we wish to draw the community’s attention to the potential data leakage issue when evaluating large language models.

5 Analyzing Factors That Influence LLM’s Translation Performance

To better understand how LLM acquires translation ability and which factors have influence on its performance, we conduct in-depth analysis. For analysis, we choose XGLM-7.5B as an example⁹. Note that, when studying a certain factor, we keep the remaining factors unchanged.

5.1 Findings on Pre-training Corpus Size

LLM can acquire translation ability in a resource-efficient way. As XGLM authors report data distribution of their pre-training corpus, we can investigate the relationship between translation performance and corpus size (Figure 3). We find that for low-resource languages, e.g., Catalan (cat) and Swahili (sw), XGLM can generate moderate translation, showing that LLM can build bilingual mapping between non-English and English with a few non-English monolingual resources (less than 1% of English resources). Even on unseen

⁹We choose XGLM for three reasons: (1) XGLM has a multilingual focus and covers many languages, which can be seen as a representative of multilingual LLM. (2) XGLM-7.5B is an open-source medium-sized LLM. It is more affordable to run experiments with it than large-sized LLM or close-source LLM. (3) The composition of the XGLM’s pre-training corpus is clear, allowing us to analyze the relationship between translation ability and corpus size.

In-context Template	Deu-Eng	Eng-Deu	Rus-Eng	Eng-Rus	Rus-Deu	Deu-Rus	Average
reasonable instructions:							
$\langle X \rangle = \langle Y \rangle$	37.37	26.49	29.66	22.25	17.66	17.31	25.12
$\langle X \rangle \setminus n$ Translate from [SRC] to [TGT]: $\setminus n \langle Y \rangle$	37.95	26.29	29.83	20.61	17.56	15.93	24.70
$\langle X \rangle \setminus n$ Translate to [TGT]: $\setminus n \langle Y \rangle$	37.69	25.84	29.96	19.61	17.44	16.48	24.50
$\langle X \rangle \setminus n$ [TGT]: $\langle Y \rangle$	29.94	17.99	25.22	16.29	12.28	11.71	18.91
$\langle X \rangle$ is equivalent to $\langle Y \rangle$	23.00	4.21	17.76	9.44	8.14	9.84	12.07
$\langle X \rangle \setminus n$ can be translated to $\setminus n \langle Y \rangle$	37.55	26.49	29.82	22.14	17.48	16.40	24.98
[SRC]: $\langle X \rangle \setminus n$ [TGT]: $\langle Y \rangle$	16.95	8.90	14.48	6.88	7.86	4.01	9.85
unreasonable instructions:							
$\langle X \rangle \$ \langle Y \rangle$	37.77	26.43	29.53	20.99	17.72	17.27	24.95
$\langle X \rangle \setminus n$ Translate from [TGT] to [SRC]: $\setminus n \langle Y \rangle$	38.18	26.21	29.85	20.35	17.75	16.63	24.83
$\langle X \rangle \setminus n$ Compile to [TGT]: $\setminus n \langle Y \rangle$	37.39	26.35	29.68	19.91	17.52	16.15	24.50
$\langle X \rangle \setminus n$ [SRC]: $\langle Y \rangle$	27.86	16.69	24.41	18.16	11.98	12.60	18.62
$\langle X \rangle$ is not equivalent to $\langle Y \rangle$	23.50	3.92	16.90	7.80	8.06	9.23	11.57
$\langle X \rangle \setminus n$ can be summarized as $\setminus n \langle Y \rangle$	37.46	26.24	29.42	22.62	17.68	17.15	25.10
[SRC]: $\langle X \rangle \setminus n$ [SRC]: $\langle Y \rangle$	19.03	8.21	15.96	6.37	7.57	4.40	10.26

Table 3: Translation performance (BLEU) of using different templates for in-context learning. The number of in-context exemplars is fixed at eight in this experiment. “ $\langle X \rangle$ ” and “ $\langle Y \rangle$ ” denote the placeholder for source and target sentence respectively. “[SRC]” and “[TGT]” represent the placeholder for source and target language name in English. Bold text denotes the highest score along the column.

languages, e.g., Occitan (oci) and Asturian (ast), XGLM can translate through ICL. These observations indicate a potential advantage of the novel translation paradigm: LLM can learn to translate in a resource-efficient way.

5.2 Findings on In-context Template

The good performance of LLMs relies on carefully-designed template The initial step of applying in-context learning for translation is determining the template. We find that the translation performance varies greatly with different templates (Table 3), where the largest gap in the average performance is up to 16 BLEU. The best template for each direction is also different. Among these templates, “ $\langle X \rangle = \langle Y \rangle$ ” achieves the highest average BLEU score. “[SRC]: $\langle X \rangle \setminus n$ [TGT]: $\langle Y \rangle$ ” achieves the lowest score, although it is a commonly-used template for prompting other LLMs, e.g., PaLM (Vilar et al., 2022), GLM (Zhang et al., 2023). Such phenomena indicate that the template plays a vital role in ICL and it may be challenging to design a universally optimal template for different LLMs and translation directions.

Even unreasonable template can instruct LLM to generate decent translation A common intuition of ICL is that the template instructs LLMs to do the target task (Brown et al., 2020), e.g., the template “ $\langle X \rangle$ can be translated to $\langle Y \rangle$ ” instructs the LLM to perform translation task. However, we find that wrapping translation exemplars with task-unrelated template can also serve as

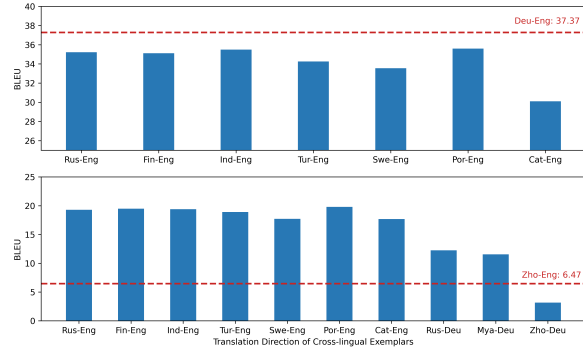


Figure 5: Effects of using cross-lingual exemplars.

an effective prompt. For example, the template like “ $\langle X \rangle$ can be summarized as $\langle Y \rangle$ ” can also instruct LLM to generate translation, rather than guiding it to generate summarization. Given the fact that these unreasonable template are also effective, the community may not fully understand the role of in-context-template.

5.3 Findings on In-context Exemplar

Cross-lingual exemplars help for certain translation directions Translation direction of the exemplar is a unique factor in machine translation. We find that using cross-lingual exemplars does not always causes worse performance and show two cases in Figure 5. When using cross-lingual exemplars for German-English translation, the translation performance degenerates. But when using cross-lingual exemplars for low-resource Chinese-English translation (illustrated in Appendix D), XGLM’s translation performance usually improves

In-context Exemplars	Consistency	Granularity	Diversity	Deu-Eng	Eng-Deu
Mismatched Translation	✗	✓	✓	0.00	0.00
Word-level Translation	✓	✗	✓	25.10	5.84
Doc-level Translation	✓	✗	✓	8.01	2.05
Duplicated Translation	✓	✓	✗	35.12	19.66
Sent-level Translation	✓	✓	✓	37.37	26.49

Table 4: Translation performance (BLEU) of XGLM when using different contents as in-context exemplars. “Consistency” column denotes whether source and target sentence are semantically consistent. “Granularity” column denotes whether the exemplar is a sentence-level pair. “Diversity” column denotes whether exemplars in the context are different from each other.

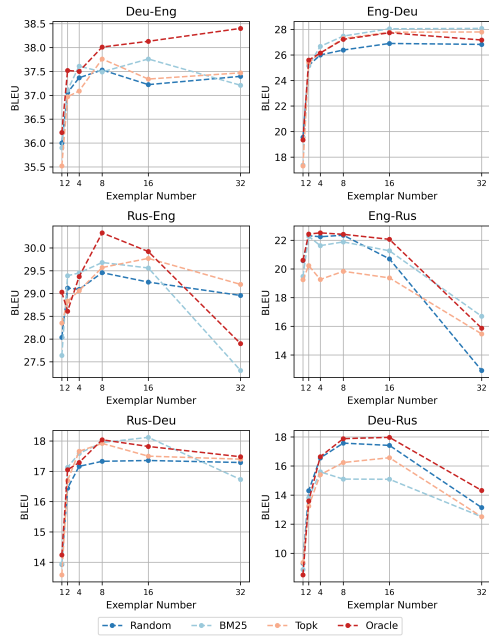


Figure 6: Effects of selecting varying number of in-context exemplars according to different strategies.

significantly, even when both source and target language is changed. This phenomenon indicates the potential usage of cross-lingual exemplars in a broader range of tasks (Lin et al., 2022), and we will explore more about this in the future.

Semantically-related exemplars does not brings more benefits than randomly-picked exemplars

In this paper, we use development set for exemplar selection, which has been found to be a high-quality candidate pool (Vilar et al., 2022), and we compare four ways of selecting in-context exemplars, namely *Random*¹⁰, *BM25*¹¹, *TopK*¹² and *Oracle*¹³.

¹⁰*Random*: picking exemplars on a random basis.

¹¹*BM25*: selecting exemplars whose source sentences are similar to the test case’s source sentence according to BM25.

¹²*TopK*: selecting exemplars whose source sentences are similar to the test case’s source sentence according to the similarity of sentence embedding.

¹³*Oracle*: selecting exemplars whose target sentences are similar to the test case’s according to sentence embedding,

Rev ratio	Deu-Eng		Eng-Deu	
	Head	Tail	Head	Tail
0 / 8	37.37	37.37	26.49	26.49
1 / 8	37.74	36.05	26.75	23.96
2 / 8	37.29	36.79	26.89	24.66
3 / 8	36.82	35.67	26.44	24.34
4 / 8	36.60	35.18	26.23	22.17
5 / 8	35.61	31.93	25.58	17.47
6 / 8	30.49	20.71	22.42	8.73
7 / 8	14.60	5.36	12.51	3.19
8 / 8	3.42	3.42	3.10	3.10

Table 5: Effects of reversing in-context examples’ translation direction. “Rev ratio” means the number of exemplars that are reversed. “Head” and “Tail” represents reversing the exemplars in the head and tail of the prompt respectively.

Effects of selecting varying number of in-context exemplars with different approaches are shown in Figure 6. The general trend in all dataset is similar. As the number of examples grows from 1 to 8, the BLEU score increases rapidly. Afterwards, the translation performance plateaus regardless of selection strategy. When more exemplars are added, e.g., 32 exemplars, the BLEU score usually starts to decline, shows an opposite phenomenon against the observation in natural language understanding tasks (Li et al., 2023).

Compared to semantically-related exemplars, randomly-picked exemplars gives comparable translation performance. Even the performance of oracle selection is on par with random selection. Based on these observations, we suggest that translation exemplars can teach LLM to translate but LLM may struggle to acquire helpful translation knowledge from semantically-related exemplars.

Exemplars teach LLM the core feature of translation task To better understand how ICL exemplars influence LLM to understand the translation task, we observe LLM’s translation behaviour under abnormal in-context exemplars (Table 4).

which can be seen as the upper bound of selection strategy.

We can see that LLM completely fails when mismatched translation is used as exemplars, indicating that LLM needs to learn from the context to keep source and target sentence semantically consistent. Word-level¹⁴ and document-level¹⁵ translation exemplar degenerates LLM’s translation performance, which demonstrates that the translation granularity of exemplar matters as well. Another interesting phenomenon is that LLM performs worse when duplicated translation is used as the exemplar, indicating that keeping in-context exemplars diverse is also important. In general, these comparison results show that LLM learns the core feature of translation task through in-context learning.

The exemplar in the tail of the prompt has more impact on the LLM’s behaviour During our analysis, we find that reversing the translation direction of exemplars will cause LLM to fail. Based on this observation, we conduct experiments to investigate the importance of different parts of the prompt (Table 5). We find that reversing exemplars in the tail of the prompt consistently produced worse results compared to reversing exemplars in the head, which suggests that exemplars in the tail of the prompt have larger influence on LLM’s behavior.

6 Related Work

In-context learning for machine translation

Using LLMs for multilingual machine translation is attracting more and more attention. Lin et al. (2022) evaluate GPT-3 and XGLM-7.5B on 182 directions. Bawden and Yvon (2023) evaluates BLOOM on 30 directions. Bang et al. (2023), Jiao et al. (2023) and Hendy et al. (2023) evaluate ChatGPT on 6 to 18 directions. In this paper, we thoroughly evaluate multilingual translation performance of popular LLMs on 102 languages and 606 directions and compare them with state-of-the-art translation engines, such as NLLB and Google Translate, which provides a more comprehensive benchmark result and highlights the challenges involved in optimizing this emerging translation paradigm.

To find better ICL recipe for machine translation, many efforts have been put into designing exemplars selection strategy (Agrawal et al., 2022; Zhang et al., 2023; Moslem et al., 2023). Similar to the findings of Zhang et al. (2023), we find that random selection is a simple but effective strategy.

¹⁴We select word pairs from open-source *fasttext* dictionary.

¹⁵We select document translation from Europarl dataset.

We also find that even oracle selection can not result in consistently better performance. Wei et al. (2022a) shows few-shot exemplars improve translation performance. And we further demonstrate the dynamic variations of translation performance with the number of in-context exemplars and the usage of cross-lingual exemplars. Besides, Vilar et al. (2022) find that using a high-quality pool, e.g., development set, for ICL example selection is better and Zhang et al. (2023) analyze why the quality of translation exemplars matters. In this paper, we reveal how in-context exemplars teach LLM to translate by analyzing LLM’s behaviour under different kinds of exemplars.

Multilingual machine translation Developing a bilingual translation system for each direction becomes impossible when the number of supporting languages increases. Therefore, multilingual machine translation is proposed (Johnson et al., 2017). But how to build a high-quality yet efficient MMT system remains an on-going challenge (Costa-jussà et al., 2022; Yuan et al., 2023; Guerreiro et al., 2023). In this paper, we focus on LLM and reveal its potential in MMT.

7 Conclusion

In this paper, we evaluate the multilingual translation ability of popular LLMs, including ChatGPT and GPT-4, on 102 languages and 606 directions, which presents the advantages and challenges of LLMs for MMT. We find that translation capabilities of LLMs are continually involving and GPT-4 reaches new performance height. However, even for GPT-4, it still face challenge on low-resource languages. In our analysis, we find that LLMs exhibit new working patterns when used for MMT. For example, instruction semantics can be ignored during in-context learning and cross-lingual exemplars can provide better task instruction for low-resource translation. More importantly, we find that LLM can acquire translation ability in a resource-efficient way, which indicates the promising future of LLM in multilingual machine translation.

Limitations

In this paper, we mainly evaluate LLM’s English-centric, French-centric and Chinese-centric translation ability. In the future, we would like to investigate more translation directions, e.g., Russian-centric translation, Arabic-centric trans-

471	lation, which could bring more findings concerning	Chaudhary, et al. 2021. Beyond english-centric multi-	524
472	with LLM’s translation ability.	lingual machine translation. <i>The Journal of Machine</i>	525
		<i>Learning Research (JMLR)</i> .	526
473	References		
474	Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke	Xavier Garcia, Yamini Bansal, Colin Cherry, George	527
475	Zettlemoyer, and Marjan Ghazvininejad. 2022. In-	Foster, Maxim Krikun, Fangxiaoyu Feng, Melvin	528
476	context examples selection for machine translation.	Johnson, and Orhan Firat. 2023. The unreasonable	529
477	<i>arXiv preprint arXiv:2212.02437</i> .	effectiveness of few-shot learning for machine trans-	530
		lation. <i>arXiv preprint arXiv:2302.01398</i> .	531
478	Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Al-	Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-	532
479	shamsi, Alessandro Cappelli, Ruxandra Cojocaru,	Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Kr-	533
480	Merouane Debbah, Etienne Goffinet, Daniel Heslow,	ishnan, Marc’ Aurelio Ranzato, Francisco Guzmán,	534
481	Julien Launay, Quentin Malartic, et al. 2023. Falcon-	and Angela Fan. 2022. The Flores-101 evaluation	535
482	40b: an open large language model with state-of-	benchmark for low-resource and multilingual ma-	536
483	the-art performance, 2023. URL https://huggingface.	chine translation. <i>Transactions of the Association for</i>	537
484	co/tiiuae/falcon-40b .	<i>Computational Linguistics (TACL)</i> .	538
485	Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-	Nuno M Guerreiro, Duarte Alves, Jonas Waldendorf,	539
486	liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei	Barry Haddow, Alexandra Birch, Pierre Colombo,	540
487	Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-	and André FT Martins. 2023. Hallucinations in	541
488	task, multilingual, multimodal evaluation of chatgpt	large multilingual translation models. <i>arXiv preprint</i>	542
489	on reasoning, hallucination, and interactivity. <i>arXiv</i>	<i>arXiv:2303.16104</i> .	543
490	<i>preprint arXiv:2302.04023</i> .		
491	Rachel Bawden and François Yvon. 2023. Investigating	Amr Hendy, Mohamed Abdelrehim, Amr Sharaf,	544
492	the translation performance of a large multilingual	Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita,	545
493	language model: the case of bloom. <i>arXiv preprint</i>	Young Jin Kim, Mohamed Afify, and Hany Hassan	546
494	<i>arXiv:2303.01911</i> .	Awadalla. 2023. How good are gpt models at ma-	547
		chine translation? a comprehensive evaluation. <i>arXiv</i>	548
		<i>preprint arXiv:2302.09210</i> .	549
495	Yoshua Bengio, Réjean Ducharme, and Pascal Vincent.	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch,	550
496	2000. A neural probabilistic language model. <i>Ad-</i>	Elena Buchatskaya, Trevor Cai, Eliza Rutherford,	551
497	<i>vances in Neural Information Processing Systems</i>	Diego de Las Casas, Lisa Anne Hendricks, Johannes	552
498	<i>(NeurIPS)</i> .	Welbl, Aidan Clark, et al. 2022. An empirical analy-	553
		sis of compute-optimal large language model training.	554
		<i>Advances in Neural Information Processing Systems</i>	555
		<i>(NeurIPS)</i> .	556
499	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	Wenxiang Jiao, Wenxuan Wang, JT Huang, Xing	557
500	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind	Wang, and ZP Tu. 2023. Is chatgpt a good trans-	558
501	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	lator? yes with gpt-4 as the engine. <i>arXiv preprint</i>	559
502	Askeel, et al. 2020. Language models are few-shot	<i>arXiv:2301.08745</i> .	560
503	learners. <i>Advances in Neural Information Processing</i>		
504	<i>Systems (NeurIPS)</i> .	Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim	561
		Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat,	562
		Fernanda Viégas, Martin Wattenberg, Greg Corrado,	563
		Macduff Hughes, and Jeffrey Dean. 2017. Google’s	564
		multilingual neural machine translation system: En-	565
		abling zero-shot translation. <i>Transactions of the As-</i>	566
		<i>sociation for Computational Linguistics (TACL)</i> .	567
505	Marta R Costa-jussà, James Cross, Onur Çelebi, Maha	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B	568
506	Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe	Brown, Benjamin Chess, Rewon Child, Scott Gray,	569
507	Kalbassi, Janice Lam, Daniel Licht, Jean Maillard,	Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.	570
508	et al. 2022. No language left behind: Scaling	Scaling laws for neural language models. <i>arXiv</i>	571
509	human-centered machine translation. <i>arXiv preprint</i>	<i>preprint arXiv:2001.08361</i> .	572
510	<i>arXiv:2207.04672</i> .		
511	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke	573
512	Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and	Zettlemoyer, and Mike Lewis. 2020. Generalization	574
513	Zhifang Sui. 2022. A survey for in-context learning.	through memorization: Nearest neighbor language	575
514	<i>arXiv preprint arXiv:2301.00234</i> .	models. In <i>International Conference on Learning</i>	576
		<i>Representations (ICLR)</i> .	577
515	Aparna Elangovan, Jiayuan He, and Karin Verspoor.		
516	2021. Memorization vs. generalization : Quantify-		
517	ing data leakage in NLP performance evaluation. In		
518	<i>Proceedings of the Conference of the European Chap-</i>		
519	<i>ter of the Association for Computational Linguistics</i>		
520	<i>(EACL)</i> .		
521	Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi		
522	Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep		
523	Baines, Onur Celebi, Guillaume Wenzek, Vishrav		

578	Mukai Li, Shansan Gong, Jiangtao Feng, Yiheng Xu,	David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo,	633
579	Jun Zhang, Zhiyong Wu, and Lingpeng Kong. 2023.	Viresh Ratnakar, and George Foster. 2022. Prompt-	634
580	In-context learning with many demonstration exam-	ing palm for translation: Assessing strategies and	635
581	ples. <i>arXiv preprint arXiv:2302.04931</i> .	performance. <i>arXiv preprint arXiv:2211.09102</i> .	636
582	Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,	637
583	Wang, Shuohui Chen, Daniel Simig, Myle Ott, Nam-	Adams Wei Yu, Brian Lester, Nan Du, Andrew M	638
584	anman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth	Dai, and Quoc V Le. 2022a. Finetuned language	639
585	Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav	models are zero-shot learners. In <i>International Con-</i>	640
586	Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettle-	<i>ference on Learning Representations (ICLR)</i> .	641
587	moyer, Zornitsa Kozareva, Mona Diab, Veselin Stoy-	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	642
588	anov, and Xian hLi. 2022. Few-shot learning with	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	643
589	multilingual generative language models. In <i>Pro-</i>	Maarten Bosma, Denny Zhou, Donald Metzler, et al.	644
590	<i>ceedings of the Conference on Empirical Methods in</i>	2022b. Emergent abilities of large language models.	645
591	<i>Natural Language Processing (EMNLP)</i> .	<i>arXiv preprint arXiv:2206.07682</i> .	646
592	Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan	Jerry W. Wei, Jason Wei, Yi Tay, Dustin Tran, Al-	647
593	Černocký, and Sanjeev Khudanpur. 2010. Recurrent	bert Webson, Yifeng Lu, Xinyun Chen, Hanxiao	648
594	neural network based language model. <i>Interspeech</i> .	Liu, Da Huang, Denny Zhou, and Tengyu Ma. 2023.	649
595	Yasmin Moslem, Rejwanul Haque, and Andy Way. 2023.	Larger language models do in-context learning dif-	650
596	Adaptive machine translation with large language	ferently. <i>CoRR</i> , abs/2303.03846.	651
597	models. <i>arXiv preprint arXiv:2301.13294</i> .	Zhenyu Wu, YaoXiang Wang, Jiacheng Ye, Jiangtao	652
598	OpenAI. 2022. https://openai.com/blog/chatgpt .	Feng, Jingjing Xu, Yu Qiao, and Zhiyong Wu. 2023.	653
599	OpenAI. 2023. Gpt-4 technical report .	Openicl: An open-source framework for in-context	654
600	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	learning. <i>arXiv preprint arXiv:2303.02913</i> .	655
601	Dario Amodei, and Ilya Sutskever. 2019. Language	Wenda Xu, Xian Qian, Mingxuan Wang, Lei Li, and	656
602	models are unsupervised multitask learners.	William Yang Wang. 2022a. Sescore2: Retrieval	657
603	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	augmented pretraining for text generation evaluation.	658
604	Lavie. 2020. COMET: A neural framework for MT	<i>arXiv preprint arXiv:2212.09305</i> .	659
605	evaluation. In <i>Proceedings of Conference on Em-</i>	Wenda Xu, Yi-Lin Tuan, Yujie Lu, Michael Saxon, Lei	660
606	<i>pirical Methods in Natural Language Processing</i>	Li, and William Yang Wang. 2022b. Not all errors are	661
607	<i>(EMNLP)</i> .	equal: Learning text generation metrics using strati-	662
608	Xiaozhe Ren, Pingyi Zhou, Xinfan Meng, Xinjing	fied error synthesis. In <i>Findings of the Association</i>	663
609	Huang, Yadao Wang, Weichao Wang, Pengfei Li,	<i>for Computational Linguistics: EMNLP 2022</i> .	664
610	Xiaoda Zhang, Alexander Podolskiy, Grigory Arshi-	Fei Yuan, Yinquan Lu, Wenhao Zhu, Lingpeng Kong,	665
611	nov, Andrey Bout, Irina Piontkovskaya, Jiansheng	Lei Li, Yu Qiao, and Jingjing Xu. 2023. Lego-mt:	666
612	Wei, Xin Jiang, Teng Su, Qun Liu, and Jun Yao. 2023.	Towards detachable models in massively multilingual	667
613	Pangu- σ : Towards trillion parameter language	machine translation. In <i>Findings of the Association</i>	668
614	model with sparse heterogeneous computing. <i>arXiv</i>	<i>for Computational Linguistics: ACL 2023</i> .	669
615	<i>preprint arXiv:2303.10845</i> .	Biao Zhang, Barry Haddow, and Alexandra Birch. 2023.	670
616	Teven Le Scao, Angela Fan, Christopher Akiki, El-	Prompting large language model for machine transla-	671
617	lie Pavlick, Suzana Ilić, Daniel Hesslow, Roman	tion: A case study. <i>arXiv preprint arXiv:2301.07069</i> .	672
618	Castagné, Alexandra Sasha Luccioni, François Yvon,	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	673
619	Matthias Gallé, et al. 2022. Bloom: A 176b-	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	674
620	parameter open-access multilingual language model.	wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022.	675
621	<i>arXiv preprint arXiv:2211.05100</i> .	Opt: Open pre-trained transformer language models.	676
622	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	<i>arXiv preprint arXiv:2205.01068</i> .	677
623	bert, Amjad Almahairi, Yasmine Babaei, Nikolay		
624	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti		
625	Bhosale, et al. 2023. Llama 2: Open founda-		
626	tion and fine-tuned chat models. <i>arXiv preprint</i>		
627	<i>arXiv:2307.09288</i> .		
628	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob		
629	Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz		
630	Kaiser, and Illia Polosukhin. 2017. Attention is all		
631	you need. In <i>Advances in Neural Information Pro-</i>		
632	<i>cessing Systems (NeurIPS)</i> .		

678	A Evaluating LLM’s translation performance with SEScore	E Used Scientific Artifacts	700
679			
680	Table 6 presents average SEScore of LLMs on	Below lists scientific artifacts that are used in our	701
681	different language families. Currently, SEScore	work. For the sake of ethic, our use of these arti-	702
682	mainly supports evaluating English translation.	facts is consistent with their intended use.	703
683	Thus we evaluate LLM’s performance on trans-		
684	lating other languages to English.	• <i>OpenICL (Apache-2.0 license)</i> , a framework	704
685		that provides an easy interface for in-context	705
686	B Detailed Results on Each Language	learning.	706
687	We report detailed results of our evaluated mod-	• <i>Transformers (Apache-2.0 license)</i> , a frame-	707
688	els in Table 7 (BLEU), Table 8 (COMET), Table	work that provides thousands of pretrained	708
689	9 (SEScore) and Figure 8. One thing that needs	models to perform tasks on different modali-	709
690	to be mentioned is that BLEU supports all transla-	ties such as text, vision, and audio.	710
691	tion directions, whereas COMET and SEScore only		
692	support a subset of these translation directions.		
693			
694	C Lists of Language		
695	We evaluate 102 languages in this paper. Table 10		
696	lists the name, ISO code and language family of		
697	these languages.		
698			
699	D Cross-lingual Exemplars		
700	In Figure 7, we show an example of using cross-		
701	lingual in-context exemplars (Russian-English ex-		
702	emplars for Chinese-English translation).		
703			
704	[Input]		
705	Этот фильм с участием Райана Гослинга и Эммы Стоун		
706	получил номинации во всех главных категориях.=The		
707	movie, featuring Ryan Gosling and Emma Stone, received		
708	nominations in all major categories.		
709	"Теперь у нас есть четырёхмесячные мыши, у которых		
710	больше нет диабета", — добавил он.=“We now have 4-		
711	month-old mice that are non-diabetic that used to be diabetic,"		
712	he added.		
713	Гослинг и Стоун получили номинации на лучшего актера и		
714	актрису соответственно.=Gosling and Stone received		
715	nominations for Best Actor and Actress respectively.		
716	Находка также позволяет ознакомиться с эволюцией перьев		
717	у птиц.=The find also grants insight into the evolution of		
718	feathers in birds.		
719	Канцелярия губернатора сообщила, что 19 из раненных		
720	были офицерами полиции.=The governor’s office said		
721	nineteen of the injured were police officers.		
722	Стандарт 802.11n работает на обеих частотах – 2.4 ГГц и		
723	5.0 ГГц.=The 802.11n standard operates on both the 2.4Ghz		
724	and 5.0Ghz frequencies.		
725	Он сказал, что создал дверной звонок, работающий от		
726	WiFi.=He built a WiFi door bell, he said.		
727	В конце 2017 года Симинофф появился на торговом		
728	телеканале QVC.=In late 2017, Siminoff appeared on		
729	shopping television channel QVC.		
730	伊拉克研究小组于格林尼治时间 (GMT) 今天 12 点提交了		
731	报告。 =		
732	[Output]		
733	The Iraqi research team submitted a report at Greenwich time		
734	(GMT) today at 12 noon.		

Figure 7: An example of using cross-lingual in-context exemplars

Language Family	Direction	Translation Performance (SESore)									
		XGLM-7.5B	OPT-175B	Falcon-7B	LLaMA-7B	LLaMA-7B-Chat	ChatGPT	GPT4	M2M-12B	NLLB-1.3B	Google
Indo-Euro-Germanic (8)	X⇒Eng	-11.78	-6.00	-8.34	-5.41	-5.90	-2.52	<u>-2.16</u>	-3.15	-2.78	-1.85
Indo-Euro-Romance (8)	X⇒Eng	-6.54	-4.01	-5.57	-3.72	-4.14	-2.30	<u>-2.08</u>	-3.08	-2.54	-2.12
Indo-Euro-Slavic (12)	X⇒Eng	-14.29	-10.31	-13.46	-5.11	-5.75	-3.55	<u>-3.17</u>	-4.21	-3.70	-2.80
Indo-Euro-Indo-Aryan (10)	X⇒Eng	-16.45	-22.15	-21.65	-17.15	-19.46	-7.64	<u>-4.69</u>	-11.77	-3.53	-2.80
Indo-Euro-Other (11)	X⇒Eng	-18.36	-17.81	-18.09	-13.61	-15.42	-6.74	<u>-4.62</u>	-7.57	-3.75	-4.40
Austronesian (6)	X⇒Eng	-14.06	-10.08	-12.30	-9.61	-10.48	-4.48	<u>-3.03</u>	-5.37	-3.47	-2.56
Atlantic-Congo (14)	X⇒Eng	-19.42	-17.61	-18.44	-17.59	-18.48	-12.38	<u>-3.34</u>	-14.16	-6.88	-5.75
Afro-Asiatic (6)	X⇒Eng	-18.85	-18.91	-19.17	-16.61	-17.66	-12.16	<u>-8.28</u>	-14.41	-4.46	-3.49
Turkic (5)	X⇒Eng	-17.15	-16.99	-18.66	-15.50	-16.47	-7.63	<u>-5.50</u>	-15.29	-4.89	-3.93
Dravidian (4)	X⇒Eng	-16.52	-22.58	-21.91	-20.18	-21.96	-9.26	<u>-5.35</u>	-13.69	-6.77	-3.07
Sino-Tibetan (3)	X⇒Eng	-19.41	-15.20	-12.37	-11.33	-12.01	-10.43	<u>-6.79</u>	-11.93	-5.50	-4.30
Other (14)	X⇒Eng	-16.74	-16.56	-18.70	-13.05	-14.17	-8.51	<u>-6.07</u>	-6.91	-4.94	-3.80

Table 6: Average SEScore of LLMs in the specific language family. The number in the bracket indicates the number of evaluated languages in the specific language family. Bold text denotes the highest SEScore across models. Underlined text denotes the highest SEScore across LLMs.

Language Family	Language	X⇒Eng (BLEU)										Eng⇒X (BLEU)										
		XGLM-7.5B	OPT-175B	Falcon-7B	LLaMA2-7B	LLaMA-7B-Chat	ChatGPT	GPT4	M2M-12B	NLLB-1.3B	Google	XGLM-7.5B	OPT-175B	Falcon-7B	LLaMA2-7B	LLaMA-7B-Chat	ChatGPT	GPT4	M2M-12B	NLLB-1.3B	Google	
Indo-European-Germanic (8)	atn	16.34	48.49	34.73	47.89	42.89	59.28	62.65	52.86	57.76	63.15	5.96	20.75	14.45	22.98	28.42	42.18	48.02	41.41	43.39	47.89	
	dan	20.65	43.54	35.31	48.33	45.83	51.23	53.18	48.32	52.35	56.44	7.91	26.81	14.80	32.79	20.19	45.49	47.46	41.12	43.81	53.23	
	fin	17.78	43.02	41.62	44.11	41.45	46.13	48.81	38.60	34.52	38.68	7.64	21.38	16.09	32.86	21.57	36.66	31.79	31.79	32.73	39.71	
	deu	34.03	39.15	34.60	41.94	39.44	43.46	47.04	42.79	44.79	48.52	25.44	23.38	20.65	30.46	26.01	41.02	44.69	40.18	40.20	49.32	
	nl	5.65	12.68	18.18	15.41	12.28	32.98	37.58	29.47	30.07	43.19	1.40	3.10	2.77	5.53	21.26	27.69	27.80	38.40	44.65	41.00	
	hi	14.13	17.96	13.60	21.87	18.36	44.57	49.20	40.04	50.37	52.52	4.74	5.54	5.10	6.32	5.72	24.65	33.89	28.04	35.08	36.80	
	nob	17.19	39.45	28.38	41.91	42.08	46.62	48.51	45.38	43.76	49.94	8.55	23.18	12.90	26.01	20.35	35.44	39.10	37.09	36.33	41.40	
	sv	22.54	44.67	37.30	46.47	44.62	50.32	51.34	48.37	49.29	55.86	12.04	27.00	18.12	33.69	28.49	46.09	49.29	47.02	45.00	53.96	
	Average		18.54	34.65	27.37	37.28	34.82	45.83	45.83	42.72	46.54	51.16	9.16	18.89	13.19	22.78	19.44	36.34	40.64	37.09	38.47	42.27
	Indo-European-Romance (8)	cat	27.65	32.20	28.84	33.88	30.90	43.18	46.41	39.06	41.65	41.00	12.70	13.11	10.96	12.89	11.24	28.24	35.45	33.43	34.01	41.00
bos		38.33	41.45	27.52	44.48	40.97	47.04	49.10	44.21	48.72	52.46	34.30	20.10	23.49	13.95	36.18	35.31	46.33	48.34	48.49	48.79	53.23
fra		36.81	43.02	41.62	44.11	41.45	46.13	48.81	38.60	34.52	38.68	7.64	21.38	16.09	32.86	21.57	36.66	31.79	31.79	32.73	39.71	
glg		29.93	36.57	29.30	37.98	35.43	43.33	42.18	38.13	43.12	44.18	12.60	18.53	12.30	16.07	14.38	38.07	39.54	38.29	37.11	41.49	
ces		8.84	22.06	23.03	39.44	35.74	43.25	42.08	41.87	44.12	47.00	2.54	15.47	8.09	21.73	21.73	35.22	39.72	37.21	38.62	41.41	
mkd		21.00	8.32	5.63	27.81	27.81	41.76	34.36	39.99	44.34	49.21	1.97	1.52	2.06	12.80	8.58	34.94	36.69	42.58	42.31	46.31	
sl		7.46	28.63	23.95	33.02	31.44	34.31	38.12	32.65	34.27	37.74	2.62	14.15	7.96	26.79	17.93	30.16	32.27	29.26	29.67	35.29	
pol		11.56	6.57	4.70	26.97	33.54	40.71	44.09	37.56	41.40	46.75	1.55	0.86	1.30	24.58	19.85	30.39	36.18	30.00	35.35	43.50	
sik		7.15	31.21	16.86	31.80	29.03	40.92	43.13	38.57	41.97	45.11	7.54	10.24	5.80	13.66	10.38	32.48	37.71	37.84	38.73	48.36	
skv		6.67	25.64	13.08	33.26	29.52	39.94	39.70	35.88	37.73	41.69	1.71	9.10	4.78	17.98	16.37	32.04	36.03	36.99	34.77	40.58	
Average		19.84	29.95	23.19	36.97	34.02	42.97	45.02	39.67	43.43	49.51	11.63	15.85	11.35	23.13	19.76	36.17	39.77	37.94	39.09	41.16	
Indo-European-Slavic (12)	bel	1.98	4.48	1.88	12.85	9.48	23.71	25.12	15.62	26.00	27.03	0.31	0.35	0.39	3.39	1.89	16.95	20.13	13.59	24.55	29.64	
	bul	7.88	34.37	21.26	39.24	37.13	44.86	48.34	41.24	44.47	49.75	1.97	10.05	7.41	23.37	18.71	34.44	37.52	33.78	37.77	43.67	
	hrv	6.66	33.37	19.48	36.35	34.68	40.02	40.22	36.28	37.62	42.60	1.44	15.71	6.19	21.96	17.66	31.90	37.84	32.54	34.94	41.41	
	ces	8.84	22.06	23.03	39.44	35.74	43.25	42.08	41.87	44.12	47.00	2.54	15.47	8.09	21.73	21.73	35.22	39.72	37.21	38.62	41.41	
	mkd	21.00	8.32	5.63	27.81	27.81	41.76	34.36	39.99	44.34	49.21	1.97	1.52	2.06	12.80	8.58	34.94	36.69	42.58	42.31	46.31	
	sl	7.46	28.63	23.95	33.02	31.44	34.31	38.12	32.65	34.27	37.74	2.62	14.15	7.96	26.79	17.93	30.16	32.27	29.26	29.67	35.29	
	pol	11.56	6.57	4.70	26.97	33.54	40.71	44.09	37.56	41.40	46.75	1.55	0.86	1.30	24.58	19.85	30.39	36.18	30.00	35.35	43.50	
	srp	7.15	31.21	16.86	31.80	29.03	40.92	43.13	38.57	41.97	45.11	7.54	10.24	5.80	13.66	10.38	32.48	37.71	37.84	38.73	48.36	
	ukr	6.67	25.64	13.08	33.26	29.52	39.94	39.70	35.88	37.73	41.69	1.71	9.10	4.78	17.98	16.37	32.04	36.03	36.99	34.77	40.58	
	Average		19.84	29.95	23.19	36.97	34.02	42.97	45.02	39.67	43.43	49.51	11.63	15.85	11.35	23.13	19.76	36.17	39.77	37.94	39.09	41.16
Indo-European-Indo-Aryan (10)	asm	4.18	1.11	1.17	3.82	1.27	18.58	27.47	-1.00	32.32	35.35	0.42	0.05	0.05	0.21	0.07	9.08	12.74	-1.00	26.02	29.77	
	ben	19.84	1.12	1.66	6.72	2.71	24.63	34.23	30.60	36.97	43.37	11.27	0.03	0.11	2.09	0.78	18.65	24.74	28.39	34.31	37.66	
	guj	0.21	1.49	1.65	4.46	1.61	18.28	27.78	31.90	34.76	45.97	10.45	0.05	0.13	1.38	0.11	10.05	14.60	7.33	27.37	39.99	
	hin	26.99	1.17	1.26	21.04	14.89	38.15	48.88	40.72	45.83	51.37	18.81	0.42	0.27	5.84	1.18	33.44	38.30	40.54	44.97	52.86	
	mar	5.63	0.97	1.00	7.37	4.78	26.94	37.08	27.29	39.25	46.02	15.88	0.06	0.07	2.17	1.83	18.22	13.12	18.27	27.66	34.71	
	nep	8.47	2.31	3.17	9.88	6.62	16.83	26.55	19.00	41.00	51.91	1.63	0.16	0.16	1.42	1.65	16.16	22.77	14.08	36.00	52.34	
	ori	0.31	0.82	1.14	1.55	1.33	17.83	33.07	0.64	39.02	42.00	0.01	0.06	0.02	0.05	0.02	10.70	18.12	0.60	32.57	41.11	
	san	0.13	1.00	1.16	2.56	1.40	24.65	42.65	42.65	42.65	42.65	0.06	0.06	0.06	0.10	0.17	14.92	10.37	14.92	17.16	21.66	
	snd	1.70	1.72	0.65	4.27	3.25	17.29	31.53	31.81	33.42	46.23	0.20	0.39	0.31	0.82	0.60	8.75	14.97	13.15	34.34	38.15	
	urd	19.31	0.79	1.09	4.95	2.49	29.53	40.27	23.94	40.07	45.69	0.33	0.20	0.20	0.43	0.27	17.58	10.24	11.37	25.54	30.64	
Average		16.91	23.38	17.45	28.00	26.20	38.33	42.99	33.85	42.46	44.07	9.62	11.71	8.40	12.47	14.89	30.99	34.92	31.76	37.76	41.12	
Indo-European-Other (11)	hye	0.15	0.32	0.74	3.87	2.05	15.30	32.20	20.70	39.99	45.84	0.02	0.05	0.01	1.19	1.51	9.02	20.47	9.89	37.54	40.91	
	ell	27.54	9.42	5.70	24.18	17.56	38.39	42.36	35.74	44.04	44.84	21.79	1.07	0.51	2.88	2.37	31.12	32.90	36.02	34.35	37.27	
	grk	4.02	17.81	16.83	13.																	

Language Family	Language	X→Eng (COMET)										Eng→X (COMET)										
		XGLM-7.5B	OPT-175B	Falcon-7B	LLaMA2-7B	LLaMA2-7B-Chat	ChatGPT	GPT4	M2M-12B	NLLB-1.3B	Google	XGLM-7.5B	OPT-175B	Falcon-7B	LLaMA2-7B	LLaMA2-7B-Chat	ChatGPT	GPT4	M2M-12B	NLLB-1.3B	Google	
Indo-European-Germanic (7)	afz	62.96	86.21	80.54	84.93	83.45	89.87	90.33	87.24	88.36	89.73	44.67	69.54	61.85	72.62	68.30	87.00	87.88	85.64	86.56	86.94	
	din	72.74	88.46	84.15	89.10	89.04	90.54	90.67	89.74	90.03	90.76	45.95	80.09	62.31	84.14	79.85	90.79	90.45	88.89	88.92	91.02	
	nd	73.32	86.41	83.69	87.42	87.27	88.62	88.64	87.42	88.04	88.39	47.82	81.52	72.40	82.56	82.11	88.87	89.03	87.43	88.39	88.97	
	deu	86.13	87.75	86.71	88.08	88.30	89.39	89.61	88.48	88.98	89.50	80.23	78.06	76.54	82.88	79.99	87.92	87.58	88.51	86.26	89.13	
	itl	50.24	62.09	54.25	66.08	64.91	85.35	87.13	83.43	85.25	87.50	29.53	32.46	32.01	37.53	42.71	79.22	82.28	81.51	82.92	87.16	
	nob	69.85	86.82	81.47	87.55	87.55	89.25	89.47	88.08	86.94	89.07	48.01	81.45	64.72	83.68	80.21	89.86	89.86	87.58	88.37	89.67	
	swe	75.42	88.26	8.87	89.32	89.22	90.32	90.50	89.78	89.67	90.54	54.84	80.67	69.69	86.26	82.20	91.09	90.90	88.61	89.74	90.46	
Average		70.99	83.71	67.40	84.73	84.25	89.05	89.48	87.74	88.18	89.26	50.21	71.97	52.93	70.05	73.63	87.83	88.30	86.47	87.31	89.05	
Indo-European-Romance (4)	cat	86.13	88.73	81.25	88.23	87.86	89.53	89.77	87.86	88.91	89.73	83.64	72.95	61.49	83.44	81.46	88.46	88.23	86.96	87.71	88.31	
	fra	86.60	88.73	88.09	88.89	88.67	89.68	89.92	88.61	89.16	89.69	81.81	82.83	84.83	83.97	81.81	88.61	88.39	86.71	87.66	88.39	
	gle	82.41	87.06	83.89	86.19	85.90	89.02	88.92	87.42	88.80	88.70	71.62	76.01	71.58	76.02	75.54	87.99	87.92	86.82	87.31	87.52	
	ron	59.62	84.95	79.26	87.53	86.66	88.97	89.17	87.79	88.03	89.63	29.94	65.99	50.67	82.80	55.16	90.57	91.27	88.71	90.18	90.57	
Average		73.58	85.18	73.58	85.95	85.54	89.25	89.57	87.95	88.47	57.07	74.55	59.27	78.23	75.82	88.25	88.66	86.87	87.60	88.95		
Indo-European-Slavic (12)	bel	48.56	60.07	50.69	70.78	66.92	83.19	84.09	75.99	83.73	84.42	31.23	33.97	31.81	42.44	44.57	79.07	82.53	71.62	85.61	87.53	
	bul	57.24	85.93	74.98	87.71	87.66	89.48	89.69	87.83	88.76	89.77	30.09	77.04	45.68	83.56	75.99	89.32	90.69	88.62	90.09	90.53	
	hrv	55.20	85.84	74.71	86.48	85.38	88.03	88.27	86.93	86.76	88.55	28.41	78.03	43.08	83.55	73.81	89.41	90.54	87.98	89.08	90.70	
	mkd	74.16	64.07	59.09	83.56	81.27	88.06	88.73	87.17	87.59	89.08	59.72	35.44	35.24	68.26	57.43	86.85	87.12	88.96	88.86	89.64	
	pol	56.44	84.42	80.90	86.10	85.97	87.26	87.25	86.13	86.63	87.42	31.42	74.55	62.60	83.39	77.50	89.16	89.85	86.65	88.26	90.13	
	rus	84.02	73.03	72.87	86.11	86.06	87.36	87.49	86.12	86.88	87.53	81.30	58.23	46.24	83.74	79.29	89.64	89.76	88.67	87.75	89.38	
	srp	65.14	59.20	57.21	86.46	86.03	87.91	88.51	86.59	87.24	88.15	43.18	32.69	32.46	82.08	73.98	86.57	89.03	84.72	87.75	89.38	
	slk	56.65	82.74	71.95	84.08	83.62	89.40	88.65	86.94	87.56	88.42	30.19	54.34	42.07	62.71	56.39	89.26	89.78	88.06	89.71	90.85	
	slv	56.46	80.35	6.70	85.41	83.30	88.11	88.45	86.25	87.07	88.53	28.92	52.41	0.39	77.55	73.24	87.58	89.48	86.88	88.39	90.02	
	ukr	71.47	69.11	65.61	86.68	85.81	87.75	88.22	86.00	87.13	87.84	40.25	48.07	0.37	82.14	76.28	88.41	89.78	87.71	88.57	90.26	
	Average		68.70	79.77	65.11	85.40	84.68	88.46	88.83	86.91	87.75	78.79	49.87	64.06	46.99	77.22	72.66	88.07	88.64	88.18	89.19	90.52
	Indo-European-Indo-Aryan (10)	asm	64.09	48.51	48.71	57.33	55.02	79.55	84.57	86.19	86.94	80.33	33.47	31.49	34.05	28.06	66.34	72.13	68.02	82.67	82.90	
ben		83.47	48.51	48.99	66.59	60.13	86.55	88.38	86.22	88.88	89.71	73.13	32.46	29.99	36.90	31.03	75.85	81.77	84.03	86.36	86.30	
gju		47.25	49.48	51.09	52.61	53.72	85.19	89.27	78.57	89.05	90.83	23.10	36.95	34.06	38.91	48.29	73.62	78.90	62.98	87.86	88.33	
hin		85.88	51.18	50.17	80.11	77.50	89.20	90.64	88.76	90.37	90.79	70.40	29.92	26.39	44.57	41.04	76.72	78.61	79.05	81.60	82.40	
mar		62.11	49.64	48.68	68.73	62.18	83.99	87.06	82.09	86.24	88.84	34.20	26.60	22.67	34.53	33.37	90.24	66.22	69.08	74.35	75.59	
ori		73.64	83.66	85.26	74.31	70.13	79.31	90.88	79.24	91.22	91.84	39.40	30.51	24.37	37.98	36.34	70.87	77.36	53.47	81.33	83.59	
ory		47.95	45.52	50.42	52.09	52.18	81.20	87.24	44.71	88.79	89.09	19.94	35.21	32.98	32.16	34.95	60.85	70.37	60.70	83.72	80.50	
pan		47.09	60.42	49.41	52.32	52.45	78.27	89.25	78.38	89.35	89.84	29.70	31.70	29.26	33.88	39.62	77.34	59.40	84.32	84.69		
und		46.69	50.29	48.52	57.19	55.39	76.01	81.96	51.82	87.13	87.91	33.10	35.43	26.87	29.17	33.65	53.11	56.16	60.01	80.44	80.30	
urd		81.11	46.53	6.48	68.14	63.22	88.08	88.54	81.15	87.88	88.53	74.95	30.25	0.28	37.90	37.12	77.07	78.66	74.31	82.83	80.20	
Average		67.26	78.56	59.08	78.50	77.29	87.15	88.51	82.86	88.07	88.99	47.18	54.30	40.58	46.71	61.77	82.12	84.32	80.64	86.48	87.48	
Indo-European-Other (9)	huc	39.23	45.76	45.65	55.59	54.10	76.40	85.17	75.72	88.41	89.26	24.05	35.48	32.49	23.90	34.59	52.00	69.40	66.32	89.80	89.72	
	ghe	84.35	65.83	60.60	79.99	76.96	87.76	88.33	86.92	87.48	88.26	85.01	46.81	36.63	46.75	43.10	88.13	88.73	88.48	88.31	89.00	
	cyj	45.61	56.53	55.41	67.80	64.88	84.67	87.30	37.70	84.79	87.84	33.91	34.20	34.79	39.83	42.27	74.21	77.68	33.84	80.53	81.99	
	eym	47.05	59.34	63.57	67.87	63.40	87.82	88.68	79.02	80.67	82.98	32.98	34.94	36.86	39.27	38.69	84.54	86.26	79.27	85.56	88.78	
	ita	85.44	86.66	87.07	87.65	87.48	88.52	89.02	87.27	88.31	88.82	82.89	82.61	84.83	84.62	82.76	88.56	89.91	87.24	88.05	88.67	
	lav	51.23	60.28	56.54	68.28	65.59	86.28	87.24	84.44	86.98	88.50	28.88	34.48	31.68	37.24	47.86	87.26	87.87	85.79	87.79		
	lit	50.67	61.28	59.07	66.68	66.03	87.26	87.54	86.52	85.84	87.43	26.84	33.25	34.31	40.70	40.16	88.09	89.82	87.86	88.62	90.59	
	pus	38.23	49.79	49.46	57.78	58.45	73.69	77.52	79.36	85.55	85.99	22.74	32.01	28.86	30.25	33.93	48.69	53.48	68.49	79.48	80.31	
	fas	55.82	49.28	49.94	77.43	71.82	87.28	88.21	86.21	87.16	88.50	30.37	28.53	27.85	43.62	39.45	84.42	86.33	84.45	86.24	87.13	
Average		64.69	68.21	57.48	76.65	75.22	86.59	88.14	81.05	87.80	45.77	51.24	39.31	60.29	47.83	81.09	83.60	79.38	86.34	87.47		
Austronesian (3)	ind	86.76	84.82	82.85	87.85	87.72	89.74	90.14	88.37	89.21	89.62	85.86	77.45	71.69	85.54	82.96	91.42	91.58	89.34	90.47	91.93	
	jav	67.78	68.26	61.53	66.81	67.84	82.77	85.33	77.56	85.71	87.10	52.22	45.55	43.42	57.62	63.86	78.23	81.82	83.30	86.65	86.24	
	ma	81.85	83.00	81.69	85.87	85.61	89.36	90.27	88.36	88.62	89.96	81.34	71.91	67.57	81.25	77.82	89.43	89.64	89.84	89.34	89.84	
Average		65.63	68.86	58.67	76.89	75.57	86.63	88.17	81.31	87.80	88.84	47.60	52.16	40.75	61.26	58.59	81.44	83.87	79.89	86.51	87.59	
Atlantic-Congo (2)	svh	81.12	61.01	6.58	61.34	60.51	87.71	88.25	83.77	86.24	88.16	77.72	33.58	34.28	34.28	39.07	85.51	86.04	83.43	85.65	85.74	
	xho	42.89	54.17	6.53	54.10	55.21	71.55	78.59	69.09	81.76	82.72	32.13	34.72	0.37	34.05	37.29	65.00	69.53	68.15	74.26	76.04	
Average																						

Language Family	Language	X→Eng (SEScore)									
		XGLM-7.5B	OPT-175B	Falcon-7B	LLaMA2-7B	LLaMA2-7B-Chat	ChatGPT	GPT4	M2M-12B	NLLB-1.3B	Google
Indo-European-Germanic (8)	afr	-13.42	-3.14	-6.72	-3.59	-4.30	-0.48	-0.19	-1.77	-1.67	-0.24
	dan	-10.98	-3.18	-5.76	-2.66	-2.68	-1.35	-1.15	-1.98	-1.65	-0.93
	nld	-10.66	-4.76	-6.31	-4.25	-4.53	-3.66	-3.59	-4.17	-3.89	-3.41
	deu	-4.44	-3.41	-4.56	-3.17	-3.23	-2.21	-2.04	-2.74	-2.37	-1.93
	isl	-20.12	-15.12	-18.05	-13.19	-14.49	-4.99	-4.09	-5.59	-5.04	-3.16
	ltz	-12.83	-11.46	-13.68	-10.39	-11.89	-3.14	-2.12	-4.06	-2.01	-1.52
	nob	-12.37	-3.96	-7.10	-3.60	-3.54	-2.52	-2.33	-2.85	-2.88	-2.24
swe	-9.42	-2.96	-4.52	-2.40	-2.57	-1.78	-1.80	-2.02	-2.33	-1.34	
Average		-11.78	-6.00	-8.34	-5.41	-5.90	-2.52	-2.16	-3.15	-2.78	-1.85
Indo-European-Romance (8)	ast	-6.71	-5.74	-6.92	-5.61	-5.95	-2.82	-2.42	-4.02	-3.47	-
	cat	-4.07	-3.69	-7.42	-2.93	-3.18	-2.00	-1.77	-2.66	-2.14	-1.34
	fra	-4.30	-2.97	-3.39	-2.87	-3.04	-2.08	-1.82	-2.62	-2.61	-1.84
	glg	-6.40	-4.17	-6.25	-4.36	-4.80	-2.64	-2.68	-3.32	-2.48	-2.67
	por	-3.53	-2.66	-3.39	-2.35	-2.76	-1.17	-1.39	-2.18	-1.82	-1.40
	ron	-15.54	-3.15	-5.78	-2.76	-3.06	-2.10	-1.92	-2.35	-2.31	-1.30
	spa	-5.88	-5.02	-5.41	-4.78	-5.13	-4.14	-4.12	-4.92	-4.57	-4.18
Average		-9.16	-5.01	-6.95	-4.56	-5.02	-2.41	-2.12	-3.11	-2.66	-1.96
Indo-European-Slavic (12)	bel	-20.43	-17.35	-19.92	-12.05	-13.99	-6.79	-6.12	-9.50	-6.25	-5.86
	bos	-17.40	-4.83	-10.39	-3.74	-4.31	-2.46	-2.28	-3.29	-2.98	-1.86
	bul	-4.33	-13.68	-15.29	-3.58	-4.49	-2.80	-2.25	-4.78	-2.93	-1.67
	hrv	-18.09	-5.29	-10.96	-4.60	-4.81	-3.73	-3.60	-3.96	-4.18	-3.24
	ces	-17.19	-5.39	-9.01	-3.88	-4.55	-2.80	-2.77	-3.54	-3.42	-2.18
	mkl	-9.69	-16.10	-17.54	-5.61	-6.86	-3.02	-2.34	-3.28	-3.04	-1.83
	pol	-17.76	-5.89	-7.88	-4.94	-4.32	-4.14	-3.61	-4.62	-4.24	-2.59
	rus	-5.97	-11.16	-11.51	-4.64	-4.76	-3.83	-3.61	-4.47	-3.68	-3.17
	srp	-13.95	-17.28	-18.10	-3.88	-4.43	-3.05	-2.56	-3.20	-3.24	-2.36
	slk	-17.84	-6.57	-11.48	-5.78	-6.38	-3.32	-2.93	-3.67	-3.54	-2.75
	slv	-17.77	-7.74	-12.32	-5.21	-5.31	-3.45	-3.31	-4.17	-3.90	-3.10
	ukr	-11.07	-12.44	-16.56	-3.39	-3.94	-3.01	-2.40	-3.57	-2.97	-1.98
Average		-11.36	-7.28	-9.74	-4.80	-5.34	-2.90	-2.57	-3.59	-3.10	-2.35
Indo-European-Indo-Aryan (10)	asm	-17.62	-22.44	-21.75	-19.23	-21.61	-10.07	-7.23	-5.46	-5.02	-5.02
	ben	-8.64	-22.29	-21.85	-16.07	-19.59	-6.90	-4.90	-5.50	-4.06	-3.18
	guj	-22.60	-22.48	-21.66	-21.04	-22.00	-7.84	-4.68	-23.21	-3.36	-2.52
	hin	-6.78	-21.75	-21.65	-9.46	-11.55	-4.05	-2.65	-3.56	-2.54	-1.69
	mar	-17.74	-22.28	-21.98	-16.22	-19.49	-7.44	-4.84	-6.94	-3.60	-2.75
	mri	-15.26	-21.15	-20.97	-14.08	-17.05	-6.54	-2.94	-10.52	-2.88	-1.41
	ory	-22.89	-22.85	-21.75	-21.22	-22.60	-10.30	-5.71	-22.74	-3.91	-3.52
	pan	-22.96	-22.04	-21.63	-20.72	-22.01	-6.20	-3.30	-8.54	-3.01	-2.27
	smd	-21.45	-21.71	-21.57	-18.93	-21.03	-11.57	-7.19	-17.98	-3.25	-2.66
	urd	-8.52	-22.49	-21.67	-14.55	-17.63	-5.52	-4.43	-6.98	-3.86	-2.99
Average		-12.70	-11.19	-12.88	-8.05	-9.05	-4.15	-3.13	-5.58	-3.22	-2.47
Indo-European-Other (11)	hye	-23.28	-22.38	-22.07	-18.87	-21.02	-11.09	-5.55	-8.90	-3.47	-2.53
	eil	-5.76	-14.88	-16.79	-7.70	-9.70	-3.66	-7.29	-3.77	-2.77	-
	gle	-20.94	-17.48	-17.96	-12.56	-14.43	-4.07	-2.30	-21.85	-3.03	-1.36
	cym	-20.63	-17.06	-17.53	-12.22	-15.38	-1.95	-0.45	-8.77	-1.78	0.24
	ita	-5.34	-4.82	-4.90	-4.12	-3.92	-3.60	-3.34	-4.02	-3.53	-3.07
	lav	-16.43	-17.88	-17.88	-14.59	-15.73	-4.44	-3.61	-4.12	-4.38	-2.79
	lit	-20.00	-16.62	-17.12	-13.55	-14.81	-4.49	-3.96	-4.27	-4.67	-3.27
	fas	-23.23	-21.32	-21.25	-18.69	-20.25	-12.83	-10.00	-7.44	-4.25	-4.08
	pus	-19.32	-21.16	-20.69	-10.03	-12.97	-4.29	-3.54	-4.59	-4.17	-2.77
	ckb	-22.58	-22.60	-21.94	-20.98	-21.42	-13.33	-7.76	-	-2.41	-2.36
	tgk	-21.03	-21.16	-20.90	-18.74	-20.04	-10.41	-6.04	-	-4.64	-3.64
Average		-13.97	-12.68	-14.05	-9.30	-10.48	-4.73	-3.46	-5.97	-3.33	-2.93
Austronesian (6)	ceb	-18.67	-9.20	-13.22	-11.12	-12.24	-4.20	-2.08	-8.10	-2.67	-1.09
	tgl	-18.04	-5.94	-10.11	-7.18	-8.10	-2.25	-1.53	-5.82	-2.43	-1.53
	ind	-4.84	-6.02	-8.15	-4.01	-4.23	-2.56	-2.33	-3.13	-2.90	-2.28
	jav	-14.95	-14.93	-16.14	-14.04	-14.98	-6.25	-3.94	-6.97	-3.56	-2.83
	msa	-6.84	-6.73	-7.99	-4.95	-5.20	-2.75	-1.68	-2.82	-2.74	-1.53
	mri	-21.00	-17.56	-18.19	-16.36	-18.08	-8.88	-6.64	-	-6.51	-6.09
Average		-13.98	-12.39	-13.86	-9.33	-10.48	-4.70	-3.42	-5.91	-3.34	-2.88
Atlantic-Congo (14)	hlg	-20.90	-18.15	-18.81	-18.11	-18.81	-14.28	-10.65	-19.72	-8.06	-7.53
	ibo	-21.95	-19.14	-19.17	-18.70	-19.60	-15.44	-11.98	-12.60	-6.29	-6.56
	kea	-14.56	-9.88	-13.84	-10.94	-11.97	-3.44	-1.64	-	-2.92	-
	kam	-19.35	-17.92	-19.02	-18.11	-18.57	-15.87	-14.95	-	-10.85	-
	lin	-19.99	-17.51	-18.27	-17.44	-18.70	-14.65	-11.86	-17.78	-6.37	-6.55
	nso	-20.19	-17.77	-18.22	-17.61	-19.05	-13.09	-7.08	-17.27	-4.41	-
	nya	-20.06	-17.96	-18.32	-18.05	-18.51	-11.50	-8.07	-	-6.96	-6.49
	sna	-21.21	-18.09	-18.80	-18.03	-19.02	-12.83	-8.95	-	-7.00	-7.08
	swb	-6.38	-16.24	-17.33	-15.75	-17.28	-2.42	-1.59	-4.08	-2.85	-3.16
	umb	-21.73	-19.15	-19.69	-19.03	-20.22	-17.99	-17.05	-	-17.75	-
	wol	-20.21	-17.34	-18.81	-17.96	-18.85	-16.39	-13.95	-18.38	-10.07	-
	xho	-22.28	-18.37	-19.12	-18.18	-18.82	-10.39	-6.29	-8.70	-4.49	-3.58
	yor	-20.94	-19.45	-19.33	-19.56	-20.03	-14.89	-11.11	-19.97	-8.74	-8.86
	zul	-22.18	-19.63	-19.41	-18.74	-19.36	-10.17	-5.55	-8.98	-4.50	-3.72
Average		-15.08	-13.45	-14.79	-11.01	-12.11	-6.26	-4.62	-7.15	-4.07	-3.30
Afro-Asiatic (6)	amh	-22.90	-22.15	-21.98	-21.75	-21.78	-19.81	-8.12	-12.19	-5.34	-3.84
	ara	-7.46	-20.31	-18.95	-8.72	-11.46	-3.72	-3.02	-4.33	-3.21	-2.36
	ful	-19.87	-18.57	-18.88	-18.57	-19.21	-17.33	-16.43	-21.80	-	-
	mlt	-19.57	-14.71	-15.96	-11.51	-12.84	-2.64	-0.69	-	-0.36	0.17
	orm	-22.10	-20.04	-19.91	-19.80	-20.86	-17.83	-14.09	-	-7.58	-6.36
	som	-21.21	-17.68	-19.32	-19.31	-19.28	-11.63	-7.32	-19.33	-3.79	-5.08
Average		-15.38	-13.89	-15.14	-11.45	-12.55	-6.73	-4.91	-7.60	-4.10	-3.31
Turkic (5)	azj	-18.08	-16.01	-18.91	-15.32	-17.02	-7.13	-6.08	-15.70	-5.94	-5.26
	kaz	-19.54	-20.68	-20.30	-17.39	-17.92	-8.76	-5.99	-20.87	-4.75	-3.48
	kir	-20.36	-21.28	-20.15	-17.95	-18.49	-11.35	-8.38	-	-5.96	-5.36
	tur	-7.45	-8.08	-15.06	-9.06	-10.21	-3.59	-2.74	-4.05	-3.84	-2.72
	uzb	-20.32	-18.90	-18.89	-17.79	-18.70	-7.34	-4.32	-20.52	-3.94	-2.85
Average		-15.50	-14.08	-15.36	-11.71	-12.79	-6.79	-4.95	-8.05	-4.15	-3.36
Dravidian (4)	kan	-22.74	-22.73	-22.14	-21.22	-22.52	-8.51	-5.29	-22.71	-4.08	-3.73
	mal	-22.96	-22.72	-21.88	-19.81	-22.07	-9.12	-4.77	-6.48	-3.15	-2.41
	tam	-10.36	-22.89	-22.09	-18.83	-21.50	-10.17	-6.08	-11.87	-4.31	-3.49
	tel	-10.00	-21.96	-21.55	-20.87	-21.75	-9.24	-5.25	-	-3.49	-2.65
Average		-15.54	-14.49	-15.67	-12.11	-13.23	-6.91	-4.97	-8.29	-4.13	-3.34
Sino-Tibetan (3)	mya	-11.16	-22.86	-22.11	-21.43	-22.57	-21.11	-11.22	-17.98	-5.48	-4.77
	zho_simpl	-22.28	-10.92	-7.05	-6.14	-6.41	-5.14	-4.53	-5.88	-3.80	-3.80
	zho_trad	-23.78	-11.81	-7.94	-6.41	-7.04	-5.03	-4.62	-	-5.44	-4.34
Average		-15.68	-14.51	-15.56	-12.08	-13.19	-7.03	-5.03	-8.39	-4.18	-3.38
Other (14)	est	-5.69	-8.04	-17.97	-13.04	-14.73	-3.16	-3.14	-3.99	-4.34	-2.54
	fin	-6.27	-5.99	-16.63	-4.99	-5.39	-3.81	-3.30	-4.53	-4.55	-3.41
	hun	-21.72	-8.73	-17.42	-5.54	-5.94	-4.10	-3.74	-4.60	-4.86	-3.72
	kat	-22.75	-22.83	-22.02	-17.02	-19.63	-11.63	-7.29	-11.45	-5.48	-4.62
	hau	-20.67	-18.03	-18.60	-18.46	-18.6					

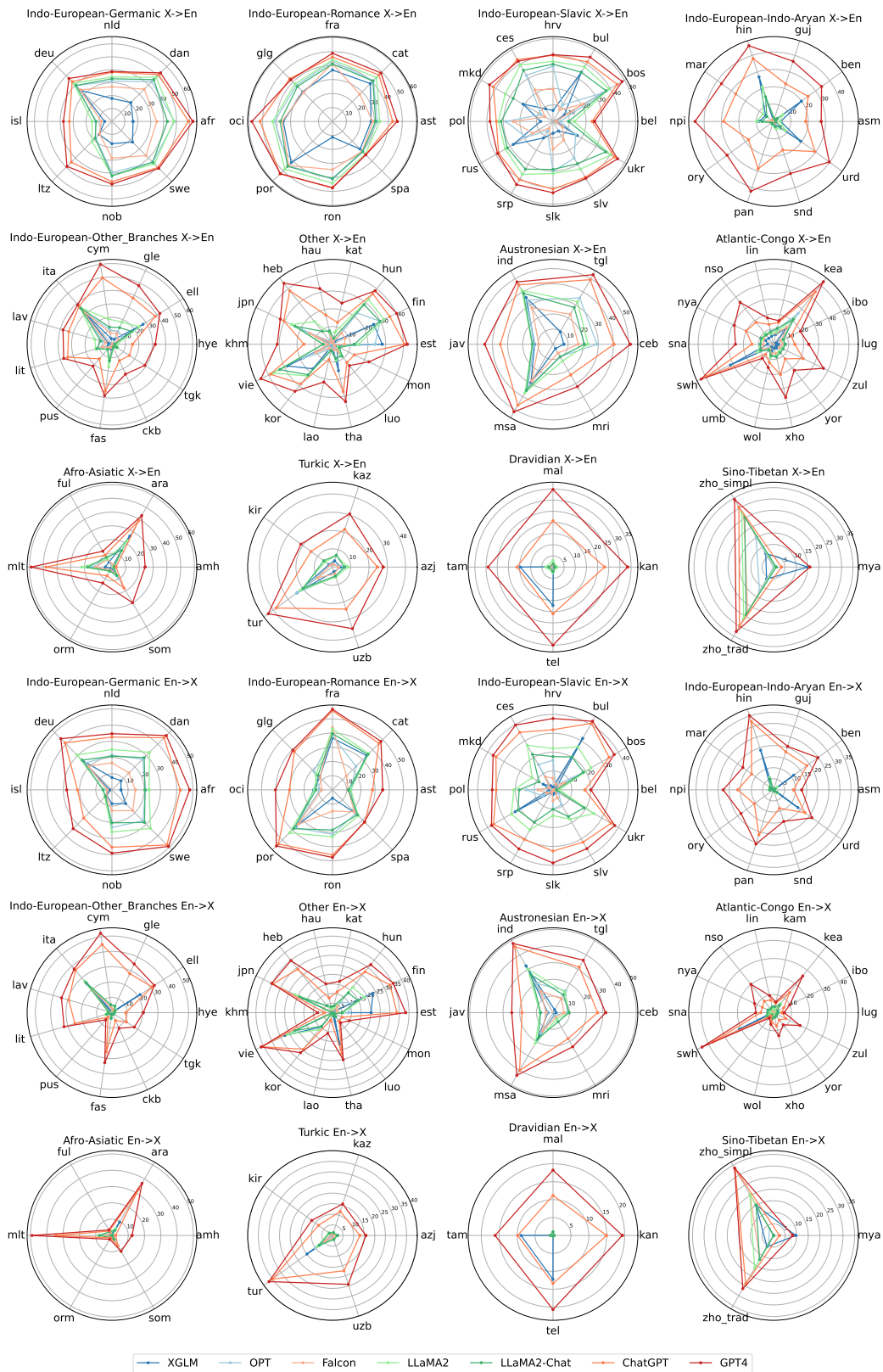


Figure 8: Comparison results between our evaluated LLMs on different language families.

Language	ISO 639-1	ISO 639-2/T	Language family	Language	ISO 639-1	ISO 639-2/T	Language family
Afrikaans	af	afr	Indo-European-Germanic	Latvian	lv	lav	Indo-European-Other
Amharic	am	amh	Afro-Asiatic	Lingala	ln	lin	Atlantic-Congo
Arabic	ar	ara	Afro-Asiatic	Lithuanian	lt	lit	Indo-European-Other
Armenian	hy	hye	Indo-European-Other	Luo	luo	luo	Other
Assamese	as	asm	Indo-European-Indo-Aryan	Luxembourgish	lb	ltz	Indo-European-Germanic
Asturian	ast	ast	Indo-European-Romance	Macedonian	mk	mkd	Indo-European-Slavic
Azerbaijani	az	azj	Turkic	Malay	ms	msa	Austronesian
Belarusian	be	bel	Indo-European-Slavic	Malayalam	ml	mal	Dravidian
Bengali	bn	ben	Indo-European-Indo-Aryan	Maltese	mt	mlt	Afro-Asiatic
Bosnian	bs	bos	Indo-European-Slavic	Maori	mi	mri	Austronesian
Bulgarian	bg	bul	Indo-European-Slavic	Marathi	mr	mar	Indo-European-Indo-Aryan
Burmese	my	mya	Sino-Tibetan	Mongolian	mn	mon	Other
Catalan	ca	cat	Indo-European-Romance	Nepali	ne	npi	Indo-European-Indo-Aryan
Cebuano	ceb	ceb	Austronesian	Northern Sotho	ns	nso	Atlantic-Congo
Chinese (Simpl)	zh	zho_simpl	Sino-Tibetan	Norwegian	no	nob	Indo-European-Germanic
Chinese (Trad)	zhttrad	zho_trad	Sino-Tibetan	Nyanja	ny	nya	Atlantic-Congo
Croatian	hr	hrv	Indo-European-Slavic	Occitan	oc	oci	Indo-European-Romance
Czech	cs	ces	Indo-European-Slavic	Oriya	or	ory	Indo-European-Indo-Aryan
Danish	da	dan	Indo-European-Germanic	Oromo	om	orm	Afro-Asiatic
Dutch	nl	nld	Indo-European-Germanic	Pashto	ps	pus	Indo-European-Other
English	en	eng	Indo-European-Germanic	Persian	fa	fas	Indo-European-Other
Estonian	et	est	Other	Polish	pl	pol	Indo-European-Slavic
Tagalog	tl	tgl	Austronesian	Portuguese	pt	por	Indo-European-Romance
Finnish	fi	fin	Other	Punjabi	pa	pan	Indo-European-Indo-Aryan
French	fr	fra	Indo-European-Romance	Romanian	ro	ron	Indo-European-Romance
Fulah	ff	ful	Afro-Asiatic	Russian	ru	rus	Indo-European-Slavic
Galician	gl	glg	Indo-European-Romance	Serbian	sr	srp	Indo-European-Slavic
Luganda	lg	lug	Atlantic-Congo	Shona	sn	sna	Atlantic-Congo
Georgian	ka	kat	Other	Sindhi	sd	snd	Indo-European-Indo-Aryan
German	de	deu	Indo-European-Germanic	Slovak	sk	slk	Indo-European-Slavic
Greek	el	ell	Indo-European-Other	Slovenian	sl	slv	Indo-European-Slavic
Gujarati	gu	guj	Indo-European-Indo-Aryan	Somali	so	som	Afro-Asiatic
Hausa	ha	hau	Other	Kurdish	ku	ckb	Indo-European-Other
Hebrew	he	heb	Other	Spanish	es	spa	Indo-European-Romance
Hindi	hi	hin	Indo-European-Indo-Aryan	Swahili	sw	swh	Atlantic-Congo
Hungarian	hu	hun	Other	Swedish	sv	swe	Indo-European-Germanic
Icelandic	is	isl	Indo-European-Germanic	Tajik	tg	tgk	Indo-European-Other
Igbo	ig	ibo	Atlantic-Congo	Tamil	ta	tam	Dravidian
Indonesian	id	ind	Austronesian	Telugu	te	tel	Dravidian
Irish	ga	gle	Indo-European-Other	Thai	th	tha	Other
Italian	it	ita	Indo-European-Other	Turkish	tr	tur	Turkic
Japanese	ja	jpn	Other	Ukrainian	uk	ukr	Indo-European-Slavic
Javanese	jav	jav	Austronesian	Umbundu	umb	umb	Atlantic-Congo
Kabuverdianu	kea	kea	Atlantic-Congo	Urdu	ur	urd	Indo-European-Indo-Aryan
Kamba	kam	kam	Atlantic-Congo	Uzbek	uz	uzb	Turkic
Kannada	kn	kan	Dravidian	Vietnamese	vi	vie	Other
Kazakh	kk	kaz	Turkic	Welsh	cy	cym	Indo-European-Other
Khmer	km	khm	Other	Wolof	wo	wol	Atlantic-Congo
Korean	ko	kor	Other	Xhosa	xh	xho	Atlantic-Congo
Kyrgyz	ky	kir	Turkic	Yoruba	yo	yor	Atlantic-Congo
Lao	lo	lao	Other	Zulu	zu	zul	Atlantic-Congo

Table 10: For each language, we list its language name, ISO code and language family.