

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

Wenhao Zhu^{1,2}, Hongyi Liu³, Qingxiu Dong⁴, Jingjing Xu²
Shujian Huang¹, Lingpeng Kong⁵, Jiajun Chen¹, Lei Li⁶

¹ National Key Laboratory for Novel Software Technology, Nanjing University

² Shanghai AI Lab ³ Shanghai Jiao Tong University ⁴ Peking University

⁵ The University of Hong Kong ⁶ Language Technologies Institute, Carnegie Mellon University

zhuwh@smail.nju.edu.cn, liu.hong.yi@sjtu.edu.cn, dqx@stu.pku.edu.cn, jingjingxu@pku.edu.cn

huangsj@nju.edu.cn, lpk@cs.hku.hk, chenjj@nju.edu.cn, leili@cs.cmu.edu

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 improving. 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 pairs¹.

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

¹Code will be released at: <https://github.com/NJUNLP/MMT-LLM>.

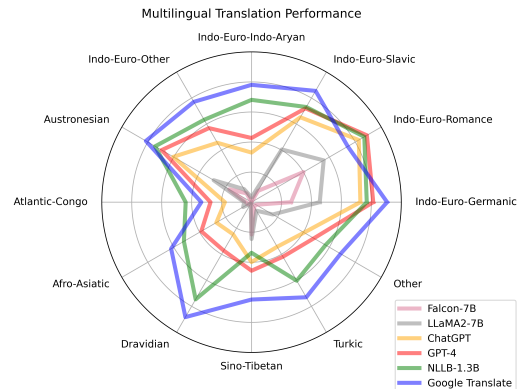


Figure 1: Multilingual translation performance (BLEU) of some popular LLMs and traditional supervised systems in translating from English to non-English. LLMs have demonstrated great potential in multilingual machine translation.

models are not particularly optimized on multilingual data.

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; Team, 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 directions). 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 (Team, 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 show 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:

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

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	<u>44.47 / 88.94</u>	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	34.00 / 84.90	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	20.24 / 76.30	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.

Language Family	Direction	Translation Performance (SEScore)									
		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	-2.08	-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>-9.34</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	-3.76	-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 2: Average SEScore 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 SEScore across models. Underlined text denotes the highest SEScore across LLMs.

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.

3 Experiment Setup

Dataset We benchmark multilingual translation on FLORES-101 (Goyal et al., 2022) dataset², which enables an assessment of model quality on a wide range of languages.

²We 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.

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-4³ (OpenAI, 2023).

ICL strategy For each model, we report its translation performance with eight randomly-picked translation pairs⁴ from the corresponding development set as in-context exemplars and “<X>=<Y>” as in-context template. “<X>” and “<Y>” are the

³We utilized GPT-3.5-TURBO-0301 for ChatGPT (evaluated at April 2023), and GPT-4-0613 for GPT-4 (evaluated at August 2023).

⁴We use the same eight randomly-picked translation pairs as exemplars during evaluation.

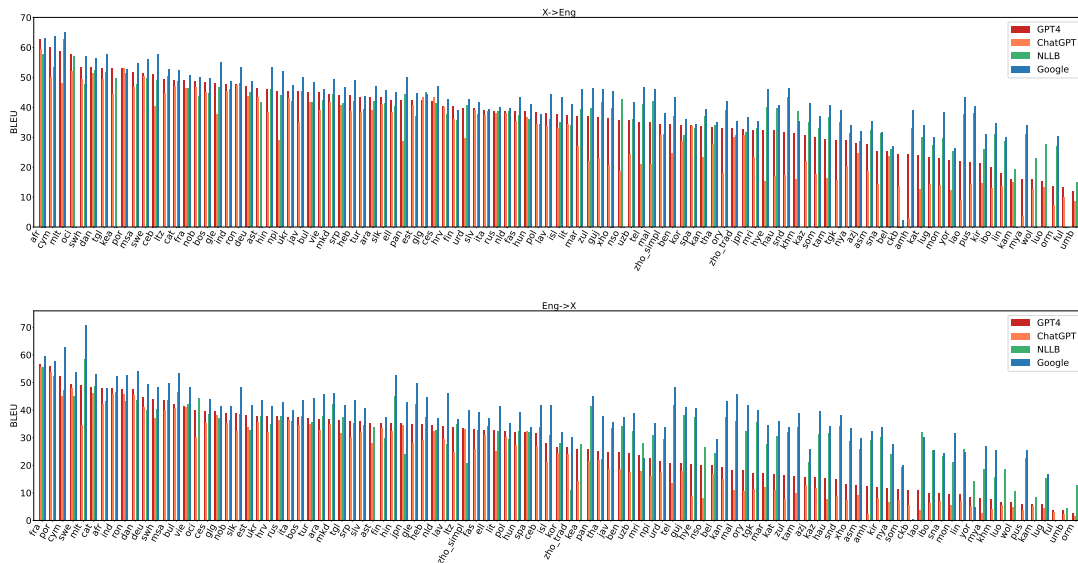


Figure 2: Translation performance (BLEU) of GPT-4, ChatGPT, NLLB and Google Translate on our evaluated languages. “X->Eng” and “Eng->X” denote translating to English and translating from English respectively. In each subfigure, languages are sorted according to BLEU scores of GPT-4.

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*⁵ (Wu et al., 2023).

Supervised baselines We report the performance of the supervised model M2M-100-12B (Fan et al., 2021) and NLLB-1.3B (Team, 2022) (distillation version), which are widely-used many-to-many MMT models. We also report the performance of the powerful commercial translation system, Google Translate⁶.

Metric Following Goyal et al. (2022), we use SentencePiece BLEU⁷ (spBLEU) as evaluation metric, which enables an evaluation of all languages. In addition, we also consider emerging metrics, COMET⁸ (Rei et al., 2020) and SEScore⁹ (Xu et al., 2022b), which have been shown to correlate well with human judgements.

⁵<https://github.com/Shark-NLP/OpenICL>

⁶<https://translate.google.com/>

⁷<https://github.com/mjpost/sacrebleu>

⁸We compute the score with *wmt22-comet-da* model.

⁹We compute the score with *SEScore-2* (Xu et al., 2022a).

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 improving Table 1 and Table 2¹⁰ present evaluation results grouped by language 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.

LLM’s capability is unbalanced across languages In Table 1, we observe a similar trend for all evaluated LLMs: they perform better at

¹⁰Currently, SEScore mainly supports evaluating English translation. Thus we evaluate LLM’s performance on translating other languages to English.

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

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 3 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.

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 2023¹¹, 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 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 example¹². 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 (swh), 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 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

¹¹The news were collected from BBC news, Fox news, ABC news and Yahoo news.

¹²We 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 LLMs or close-source LLMs. (3) The composition of the XGLM’s pre-training corpus is clear, allowing us to analyze the relationship between translation ability and corpus size.

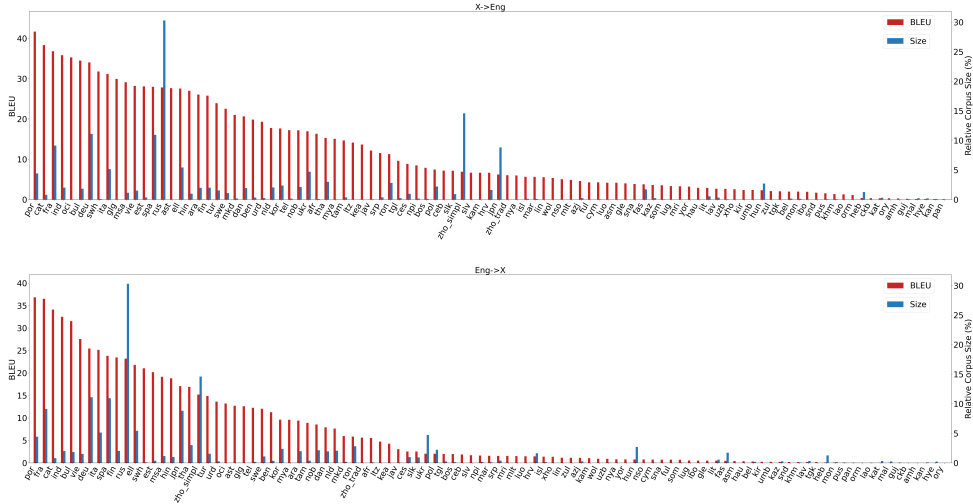


Figure 3: Translation performance (BLEU) of XGLM on evaluated languages and the corpus size of each language relative to English pre-training corpus. In each subfigure, languages are sorted according to BLEU scores of XGLM.

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
$\langle X \rangle \setminus \text{n Translate from [SRC] to [TGT]: } \setminus \text{n } \langle Y \rangle$	37.95	26.29	29.83	20.61	17.56	15.93	24.70
$\langle X \rangle \setminus \text{n Translate to [TGT]: } \setminus \text{n } \langle Y \rangle$	37.69	25.84	29.96	19.61	17.44	16.48	24.50
$\langle X \rangle \setminus \text{n [TGT]: } \langle Y \rangle$	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
$\langle X \rangle \setminus \text{n can be translated to } \setminus \text{n } \langle Y \rangle$	37.55	26.49	29.82	22.14	17.48	16.40	24.98
[SRC]: $\langle X \rangle \setminus \text{n [TGT]: } \langle Y \rangle$	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
$\langle X \rangle \setminus \text{n Translate from [TGT] to [SRC]: } \setminus \text{n } \langle Y \rangle$	38.18	26.21	29.85	20.35	17.75	16.63	24.83
$\langle X \rangle \setminus \text{n Compile to [TGT]: } \setminus \text{n } \langle Y \rangle$	37.39	26.35	29.68	19.91	17.52	16.15	24.50
$\langle X \rangle \setminus \text{n [SRC]: } \langle Y \rangle$	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
$\langle X \rangle \setminus \text{n can be summarized as } \setminus \text{n } \langle Y \rangle$	37.46	26.24	29.42	22.62	17.68	17.15	25.10
[SRC]: $\langle X \rangle \setminus \text{n [SRC]: } \langle Y \rangle$	19.03	8.21	15.96	6.37	7.57	4.40	10.26

Table 4: 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.

applying in-context learning for translation is determining the template. We find that the translation performance varies greatly with different templates (Table 4), 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]: $\langle X \rangle \setminus \text{n [TGT]: } \langle Y \rangle$ ” 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 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

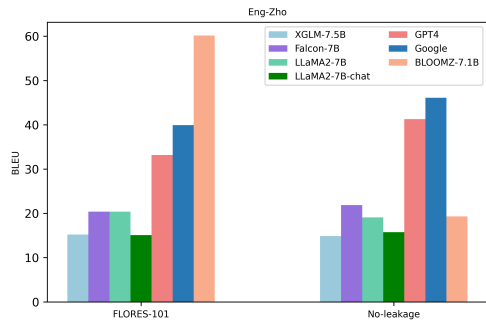


Figure 4: Translation performance of different models on FLORES-101 test set and our annotated no-leakage evaluation set NEWS2023.

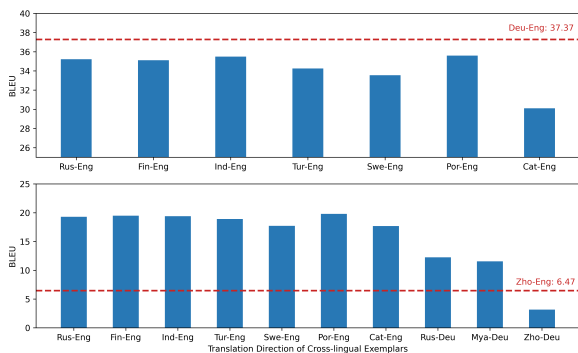


Figure 5: Effects of using cross-lingual exemplars.

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 C), XGLM’s translation performance usually improves 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,

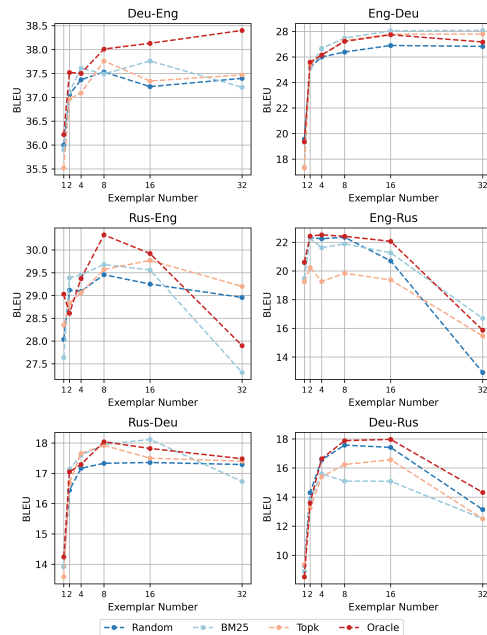


Figure 6: Effects of selecting varying number of in-context exemplars according to different strategies.

namely *Random*¹³, *BM25*¹⁴, *TopK*¹⁵ and *Oracle*¹⁶.

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 exem-

¹³*Random*: picking exemplars on a random basis.

¹⁴*BM25*: selecting exemplars whose source sentences are similar to the test case’s source sentence according to BM25.

¹⁵*TopK*: selecting exemplars whose source sentences are similar to the test case’s source sentence according to the similarity of sentence embedding.

¹⁶*Oracle*: selecting exemplars whose target sentences are similar to the test case’s according to sentence embedding, which can be seen as the upper bound of selection strategy.

In-context Exemplars	Consistency	Granularity	Diversity	Deu-Eng	Eng-Deu	Zho-Eng	Eng-Zho
Mismatched Translation	✗	✓	✓	0.00	0.00	0.42	1.16
Word-level Translation	✓	✗	✓	25.10	5.84	2.81	2.24
Doc-level Translation	✓	✗	✓	8.01	2.05	4.48	2.20
Duplicated Translation	✓	✓	✗	35.12	19.66	17.87	7.86
Sent-level Translation	✓	✓	✓	37.37	26.49	19.86	11.07

Table 5: 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.

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 6: 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.

plars influence LLM to understand the translation task, we observe LLM’s translation behaviour under abnormal in-context exemplars (Table 5).

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-level¹⁷ and document-level¹⁸ 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 6). We find that reversing exemplars in the

¹⁷We select word pairs from open-source *fasttext* dictionary.

¹⁸We select document translation from Europarl dataset.

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), Hendy et al. (2023) and Peng 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. 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 (Team, 2022; Yuan et al., 2023; Guerreiro et al., 2023; Robinson 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 translation, which could bring more findings concerning with LLM’s translation ability.

Acknowledgement

We would like to thank Fei Yuan, Zhenyu Wu, Yunzhe Lv for their support to this project. Shujian Huang is the corresponding author. This work is partially supported by National Science Foundation of China (No. 62376116, 62176120), the Liaoning Provincial Research Foundation for Basic Research (No. 2022-KF-26-02) and the research project of Nanjing University-China Mobile Joint Institute.

References

- Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2022. In-context examples selection for machine translation. *arXiv preprint arXiv:2212.02437*.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, et al. 2023. Falcon-40b: an open large language model with state-of-the-art performance, 2023. URL <https://huggingface.co/tiiuae/falcon-40b>.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-task, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.
- Rachel Bawden and François Yvon. 2023. Investigating the translation performance of a large multilingual language model: the case of bloom. *arXiv preprint arXiv:2303.01911*.
- Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. 2000. A neural probabilistic language model. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*.
- Aparna Elangovan, Jiayuan He, and Karin Verspoor. 2021. Memorization vs. generalization : Quantifying data leakage in NLP performance evaluation. In *Proceedings of the Conference of the European Chapter of the Association for Computational Linguistics (EACL)*.
- Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. 2021. Beyond english-centric multilingual machine translation. *The Journal of Machine Learning Research (JMLR)*.
- Xavier Garcia, Yamini Bansal, Colin Cherry, George Foster, Maxim Krikun, Fangxiaoyu Feng, Melvin Johnson, and Orhan Firat. 2023. The unreasonable effectiveness of few-shot learning for machine translation. *arXiv preprint arXiv:2302.01398*.

- 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. *Transactions of the Association for Computational Linguistics (ACL)*.
- Nuno M Guerreiro, Duarte Alves, Jonas Waldendorf, Barry Haddow, Alexandra Birch, Pierre Colombo, and André FT Martins. 2023. Hallucinations in large multilingual translation models. *arXiv preprint arXiv:2303.16104*.
- Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. An empirical analysis of compute-optimal large language model training. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang Huang, Xing Wang, and Zhaopeng Tu. 2023. Is chatgpt a good translator? yes with gpt-4 as the engine. *arXiv preprint arXiv:2301.08745*.
- Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics (ACL)*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through memorization: Nearest neighbor language models. In *International Conference on Learning Representations (ICLR)*.
- Mukai Li, Shansan Gong, Jiangtao Feng, Yiheng Xu, Jun Zhang, Zhiyong Wu, and Lingpeng Kong. 2023. In-context learning with many demonstration examples. *arXiv preprint arXiv:2302.04931*.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona Diab, Veselin Stoyanov, and Xian hLi. 2022. Few-shot learning with multilingual generative language models. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan Černocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. *Interspeech*.
- Yasmin Moslem, Rejwanul Haque, and Andy Way. 2023. Adaptive machine translation with large language models. *arXiv preprint arXiv:2301.13294*.
- OpenAI. 2022. <https://openai.com/blog/chatgpt>.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2023. Towards making the most of ChatGPT for machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Xiaozhe Ren, Pingyi Zhou, Xinfan Meng, Xinjing Huang, Yadao Wang, Weichao Wang, Pengfei Li, Xiaoda Zhang, Alexander Podolskiy, Grigory Arshinov, Andrey Bout, Irina Piontkovskaya, Jiansheng Wei, Xin Jiang, Teng Su, Qun Liu, and Jun Yao. 2023. Pangu- σ : Towards trillion parameter language model with sparse heterogeneous computing. *arXiv preprint arXiv:2303.10845*.
- Nathaniel Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. ChatGPT MT: Competitive for high- (but not low-) resource languages. In *Proceedings of the Eighth Conference on Machine Translation (WMT)*.
- Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- NLLB Team. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2022. Prompting palm for translation: Assessing strategies and performance. *arXiv preprint arXiv:2211.09102*.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2022a. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations (ICLR)*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022b. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.
- Jerry W. Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. 2023. Larger language models do in-context learning differently. *CoRR*, abs/2303.03846.
- Zhenyu Wu, YaoXiang Wang, Jiacheng Ye, Jiangtao Feng, Jingjing Xu, Yu Qiao, and Zhiyong Wu. 2023. Openicl: An open-source framework for in-context learning. *arXiv preprint arXiv:2303.02913*.
- Wenda Xu, Xian Qian, Mingxuan Wang, Lei Li, and William Yang Wang. 2022a. Sescore2: Retrieval augmented pretraining for text generation evaluation. *arXiv preprint arXiv:2212.09305*.
- Wenda Xu, Yi-Lin Tuan, Yujie Lu, Michael Saxon, Lei Li, and William Yang Wang. 2022b. Not all errors are equal: Learning text generation metrics using stratified error synthesis. In *Findings of the Association for Computational Linguistics: EMNLP 2022*.
- Fei Yuan, Yinqun Lu, Wenhao Zhu, Lingpeng Kong, Lei Li, Yu Qiao, and Jingjing Xu. 2023. Lego-mt: Towards detachable models in massively multilingual machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023*.
- Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study. *arXiv preprint arXiv:2301.07069*.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.

A Detailed Results on Each Language

We report detailed results of our evaluated models in Table 7 (BLEU), Table 8 (COMET), Table 9 (SEScore) and Figure 8. One thing that needs to be mentioned is that BLEU supports all translation directions, whereas COMET and SCScore only support a subset of these translation directions.

B Lists of Language

We evaluate 102 languages in this paper. Table 10 lists the name, ISO code and language family of these languages.

C Cross-lingual Exemplars

In Figure 7, we show an example of using cross-lingual in-context exemplars (Russian-English exemplars for Chinese-English translation).

D Used Scientific Artifacts

Below lists scientific artifacts that are used in our work. For the sake of ethic, our use of these artifacts is consistent with their intended use.

- *OpenICL (Apache-2.0 license)*, a framework that provides an easy interface for in-context learning.

- *Transformers (Apache-2.0 license)*, a framework that provides thousands of pretrained models to perform tasks on different modalities such as text, vision, and audio.

[Input]

Этот фильм с участием Райана Гослинга и Эммы Стоун получил номинации во всех главных категориях.=The movie, featuring Ryan Gosling and Emma Stone, received nominations in all major categories.

"Теперь у нас есть четырёхмесячные мыши, у которых больше нет диабета", — добавил он.=“We now have 4-month-old mice that are non-diabetic that used to be diabetic,” he added.

Гослинг и Стоун получили номинации на лучшего актёра и актрису соответственно.=Gosling and Stone received nominations for Best Actor and Actress respectively.

Находка также позволяет ознакомиться с эволюцией перьев у птиц.=The find also grants insight into the evolution of feathers in birds.

Канцелярия губернатора сообщила, что 19 из раненных были офицерами полиции.=The governor’s office said nineteen of the injured were police officers.

Стандарт 802.11n работает на обеих частотах – 2.4 ГГц и 5.0 ГГц.=The 802.11n standard operates on both the 2.4Ghz and 5.0Ghz frequencies.

Он сказал, что создал дверной звонок, работающий от WiFi.=He built a WiFi door bell, he said.

В конце 2017 года Симинофф появился на торговом телеканале QVC.=In late 2017, Siminoff appeared on shopping television channel QVC.

Иракский исследовательский小组于格林尼治时间 (GMT) 今天 12 点提交了报告。 =

[Output]

The Iraqi research team submitted a report at Greenwich time (GMT) today at 12 noon.

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

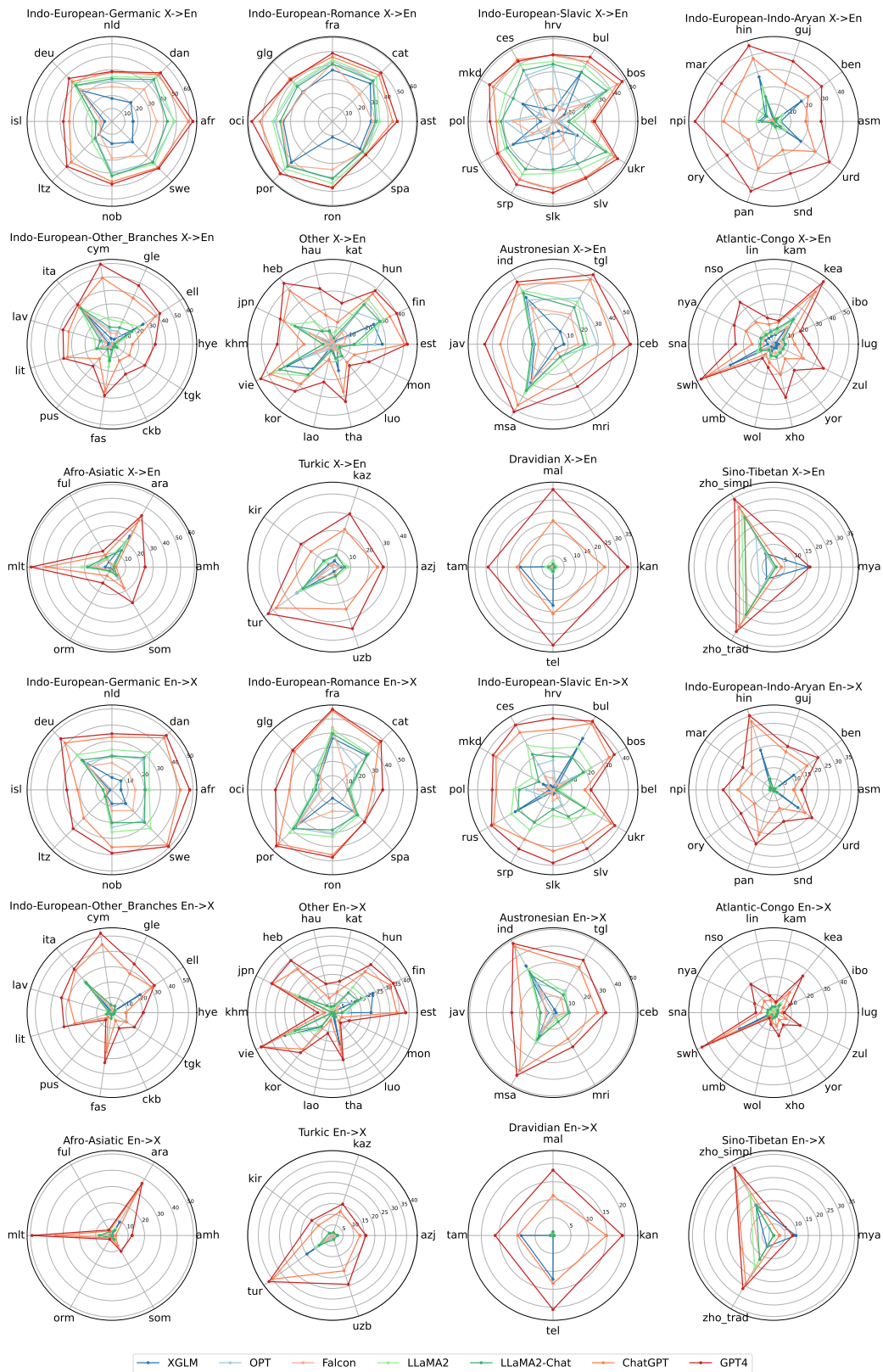


Figure 8: Comparison results (BLEU) between our evaluated LLMs on different language families.

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.