

Evaluation of Pre-Trained Models for Many-to-English Translation

Stefan Liemawan Adjii

August 14, 2024

Abstract

Machine translation (MT) has become an essential tool for overcoming language barriers in an increasingly globalised world. Today, there are pre-trained models (PTMs) specifically designed for automatic machine translation. This paper evaluates the performance of various pre-trained models (PTMs) for many-to-English translation. The models assessed include OPUS-MT [44], mBART-50 [25], NLLB-200 [40], and M2M-100 [14]. Using a dataset curated from the Tatoeba repository, comprising 1,323 parallel sentence pairs across 14 source languages, translation quality is measured using BLEU, SacreBLEU, and METEOR scores. The results indicate that OPUS-MT, a one-to-one PTM, outperforms multilingual models in one-to-English translation tasks, with mBART-50 leading among the multilingual PTMs. These findings highlight the strengths and limitations of both specific language and multilingual PTMs in machine translation. (Findings)

1 Introduction

According to Ethnologue [12], 7,164 languages currently exist and in use today, with 40% of them considered endangered. As of July 2024, 243 languages are supported by Google Translate (according to Wikipedia [49]). In modern times, the need for translation services has surged due to the growing exchange of information across different regions that speak various languages [29].

Machine translation (MT) is the task of automatically translating from one language to another. This can be done through text or audio. It can be traced back to 1949 [47], with the first public demonstration of an MT system on January 7, 1954, in collaboration with IBM, where 49 Russian sentences were translated into English using a limited vocabulary of 250 words and 6 grammar rules [19]. However, over the next several decades, growth were limited for machine translation, with 1956-1966 considered the decade of high expectation and disillusion, and 1967-1976 dubbed 'the quiet decade' [20]. Then in 1989, the dominance of the rule-based approach has been challenged by the rise of new methods and strategies, collectively referred to as 'corpus-based' methods (data-driven) [17, 18]. Subsequently, statistics-based approaches for MT re-emerged, bolstered by the recent success of probabilistic techniques in speech recognition. Statistical machine translation [26] dominated the domain between late 1990s through the early 2010s, before largely being surpassed by neural machine translation (NMT) [9, 38].

Since the introduction of Transformers in 2017 [46], Natural Language Processing (NLP) and machine translation in particular reached a giant milestone. The following years saw the birth of Large Language Models (LLMs) such as BERT [11], GPT [30], and T5 [35], which revolutionised both MT and the whole field of NLP. Then in early 2020s, several pre-trained models (PTM) that are specifically designed for machine translation emerged, OPUS-MT [44] produced numerous PTMs designed for one-to-one translations. Additionally, multilingual models: mBART-50 [25], NLLB-200 [40], and M2M-100 [14] have also been published, growing popularity as multilinguality has been shown to allow for sharing information among languages [15]. Most of these models are trained on diverse set of languages, allowing for many-to-many translation: able to translate between any of the supported pair of languages. This allows the models to generalise over shared lexical and linguistic among languages, and have been shown to increase performance compared to one-to-one translation models [25].

Despite these advancements, pre-trained models are often evaluated using different set of languages and benchmarks [25, 40, 14, 48], making it difficult to gauge their relative effectiveness across various languages. Nevertheless, there does not seem to be much work on comparing or benchmarking different pre-trained models in machine translation. Intento published 'The State of Machine Translation 2024' [42] providing an in-depth evaluation of popular MT engines and LLMs. However, the biggest drawback

in this report is that the selected LLMs are general LLMs such as GPT, LLaMa, Mistral, instead of MT-specific LLMs.

Through simple experimentations, this paper aims to evaluate the performance of existing pre-trained models (PTMs) on many-to-English translation across 14 source languages. Although fine-tuning multilingual PTMs has been proven to increase model performance [10], no pre-training or fine-tuning is performed in this study for simplicity reasons. A dataset is curated from the Tatoeba repository [39], containing 1,323 parallel sentence pairs across source and target languages. The models include a one-to-one PTMs: OPUS-MT [45, 44], and multilingual PTMs: such as mBART-50 [25], NLLB-200 [40], and M2M-100 [14]. The performance of these models is evaluated using the BLEU [33], SacreBLEU [34], and METEOR [23].

The result indicates that OPUS-MT outperforms all other multilingual models in one-to-English translation, highlighting the superiority of specific languages PTMs compared to multilingual PTMs. mBART-50 shows the best performance among other multilingual PTMs: M2M-100 and NLLB-200. (More findings here)

2 Literature Review

2.1 Pre-Trained Models for Machine Translation

The encoder-decoder approach [9] remains as the foundation architecture for many sequence-to-sequence models in machine translation. Attention mechanism has been used for Neural Machine Translation (NMT) [5, 27] before the release of the transformer [46]

Currently, Pre-Trained-Models (PTMs) have become a cornerstone in machine translation, leveraging large datasets and powerful architectures to achieve state-of-the-art performance [16]. Most PTMs support two kinds of translations: one-to-one translation and many-to-many translation (multilingual models).

One-to-One Translation refers to a translation approach where a model is specifically trained to translate between one source language and one target language. This setup is characterised by having a dedicated model for each unique language pair. OPUS-MT [44] and Facebook Fairseq [32] are examples of PTMs that support one-to-one translation. Multilingual models have also gained significant prominence, enabling many-to-many translation, where a single model can translate between multiple source and target languages [3].

OPUS-MT [44] is a project dedicated to creating and providing free resources for machine translation, providing over a thousand of PTMs of direct one-to-one translations. In contrast, mBART [25], M2M [14], NLLB-200 [40], are massively multilingual models that support many-to-many translation language.

2.2 Parallel Corpora

Corpora are large and structured sets of texts used for linguistic research and analysis. Thus, parallel corpora are defined as sets of texts in a given source language along with their translations in another target language [24]. It can be bilingual or multilingual, and are crucial for training and evaluating machine translation. Languages without or with limited parallel corpora are referred as **low-resources language**.

Since the early 2000s, English has often been used as an intermediary language, meaning texts originally written in languages are first translated into English, and then from English into other languages, resulting in most target texts being translations of translations [24]. This is called indirect translation, and posed for many problems in translations as terms can be ambiguous with varying or similar meanings [4].

Neural Machine Translation (NMT) systems require vast amounts of training data, and thus the availability of parallel corpora is crucial for building effective models [22]. The lack of extensive parallel corpora, especially for low-resource languages, leads to suboptimal performance in NMT techniques compared to their high-resource counterparts [36]. OPUS [31] is a comprehensive collection of open-source parallel corpora used extensively in the field of machine translation (MT). It includes corpora for 744 languages and contains over 1,210 different datasets, amassing a total of 45,945,946,108 sentence pairs. Tatoeba [39] is another prominent resource in the field of MT and NLP, known for its extensive collection of translated sentences. As of July 2024, it contains 12,186,207 sentences over 423 supported languages, growing daily through volunteer contributions.

Several datasets are often used for evaluation in machine translation. The Tatoeba Challenge [43] covers 487 languages in 4,024 language pairs, including 657 test sets sourced from Tatoeba website, covering 138 languages. The TED Talks dataset [52] contains transcripts from TED talks for more than 50 languages. IWSLT [2] also contains TED talks data, but paired with English translations. Finally, WMT [6] is an annual event that organises tasks for machine translation and provides a collection of datasets for benchmarking and evaluating translation systems.

While these datasets can be used to evaluate translation performance between any language pairs, they do not contain parallel corpora between different languages. Corpus between English and French for example, contains different text to the corpus between English and Spanish. Therefore, they do not support easy evaluation for many-to-English translations. Furthermore, papers often use different datasets and versions to evaluate their model performance, making direct comparison challenging.

mBart [25] uses WMT19 for English-German translation and TED15 for Chinese-English translation. M2M-100 [14] uses 7 different datasets shared across languages. NLLB-200 [40] uses FLORES dataset, designed for low-resources languages. OPUS-MT [45, 44] models use datasets from Tatoeba challenges.

3 Methodology

3.1 Dataset











Tatoeba is a vast, continuously expanding database consisting sentences and their translations, built through the contributions of thousands of volunteers, offering a tool that allows users to see examples of how words are used in sentences [39]. They currently have 12,132,349 sentences and 423 supported languages, with around one to two thousand new sentences added daily, on average. The English sentence dataset contains 1,905,089 sentences, the largest one in their repository, with Russian in the second place with 1,066,633 sentences. Some languages supported on the website is shown in Figure 1 and Figure 2, sorted from the biggest corpus.



Figure 1: Tatoeba’s languages repository with 10,000+ sentences and 100,000+ sentences [39]

Table 1 show the 14 languages selected for this project. Languages are chosen based on its resources’ availability in Tatoeba, as well as considering supported languages in most PTMs models. Thus, the languages chosen here can be considered as moderate to high resources languages.

To build the dataset, sentences in English are first downloaded, containing 1,898,494 sentences (it is unclear why it is less than the number stated in the Tatoeba website). Then for each language, sentence pairs between English and source languages are downloaded individually and compiled. The result is a single Dataframe containing 1,323 parallel sentences in all 14 languages, this will be treated as a test set to evaluate the models performance on each language. Note that in the case of multiple reference

		Language	Sentences
1		eng English	1,906,613
2		rus Russian	1,067,167
3		ita Italian	881,287
4		epo Esperanto	760,064
5		tur Turkish	734,083
6		kab Kabyle	714,233
7		deu German	667,177
8		ber Berber	660,836
9		fra French	614,521
10		por Portuguese	432,384
11		spa Spanish	410,509
12		hun Hungarian	409,148
13		jpn Japanese	243,341
14		heb Hebrew	201,220
15		ukr Ukrainian	186,145
16		nld Dutch	185,628
17		fin Finnish	149,285
18		pol Polish	127,893
19		lit Lithuanian	108,016
20		ces Czech	79,393










404		kxi Keningau Murut	4
405		tso Tsonga	4
406		crk Plains Cree	4
407		hsn Xiang Chinese	4
408		hnj Hmong Njua (Green)	4
409		pfl Palatine German	3
410		syc Syriac	3
411		ayl Libyan Arabic	3
412		mni Meitei	3
413		hdn Northern Haida	3
414		gan Gan Chinese	3
415		osx Old Saxon	3
416		gaa Ga	3
417		urh Urhobo	2
418		aym Aymara	2
419		nys Nyungar	2
420		sot Southern Sotho	2
421		mnc Manchu	2
422		rel Rendille	1
423		hax Southern Haida	1
424		cyo Cuyonon	1

Figure 2: Tatoeba top 20 and bottom 20 languages based on sentences count [39]

No.	Language Pairs	Sentence Pairs Count
1	Chinese-English	68,814
2	Dutch-English	155,856
3	Finnish-English	102,202
4	French-English	405,088
5	German-English	501,145
6	Hebrew-English	172,082
7	Hungarian-English	171,698
8	Italian-English	624,160
9	Japanese-English	270,116
10	Polish-English	77,345
11	Russian-English	722,837
12	Spanish-English	265,253
13	Turkish-English	710,279
14	Ukrainian-English	214,244

Table 1: List of chosen languages for evaluation

translations are available for the same sentence in English, only the first sentence is taken while the rest is discarded. Therefore, there is only one reference for each candidate sentence.

Sentences typically consist of everyday phrases such as 'I have to go to sleep', 'That is intriguing', and 'Where do you live?'. They may also include single-word exclamations like 'Speak!', 'So what?', or 'Look!'. Additionally, multiple sentences such as 'You may write in any language you want. On Tatoeba, all languages are considered equal', and 'Guns don't kill people. People kill people' can be found inside the corpus. A few of them also include human names, 'Compare your answer with Tom's', 'Muiriel is 20 now'. All of the sentences are straightforward and literal, without the use of linguistic features such as metaphors or sarcasm. Therefore, machine translation process should be straightforward on this level. Figure 2 shows two examples of parallel sentences in the final dataset, across all 14 languages and the original English sentence.

Language	Sentence 1	Sentence 2
English	I have to go to sleep.	So what?
Chinese	我该去睡觉了。	那又怎?
Dutch	Ik moet gaan slapen.	Dus?
Finnish	Minun täytyy mennä nukkumaan.	Mitä sitten?
French	Je dois aller dormir.	Et alors?
German	Ich muss jetzt schlafen.	Na und?
Hebrew	<hidden-due-to-latex-incompatibility>	<hidden-due-to-latex-incompatibility>
Hungarian	Aludni kell mennem.	És akkor mi van?
Italian	Devo andare a dormire.	E allora?
Japanese	私は眠らなければなりません	だから何?
Polish	Muszę iść spać.	No i co?
Russian	Мне пора идти	Так что?
Spanish	Tengo que irme a dormir.	¿Entonces qué?
Turkish	Yatmaya gitmek zorundayım.	Öyleyse ne yapmalı?
Ukrainian	Маю пти спати.	Ну то що?

Table 2: A snippet of the dataset

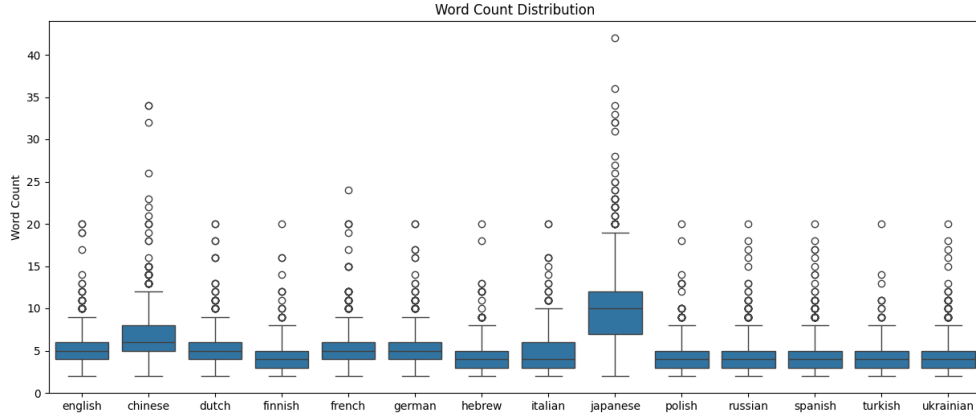


Figure 3: Dataset word count distribution per language, Chinese and Japanese sentences are counted per letter.

Figure 3 shows a box plot of sentences word count. The median for Dutch, Finnish, French, German, Hebrew, Italian, Polish, Russian, Spanish, Turkish, Ukrainian seems relatively low, with a moderate spread. A few to moderate count of outlier exist in all languages, suggesting occasional longer sentences, with the highest variability seen in Japanese. Note that for Chinese and Japanese the count is based on each Chinese/Japanese letter, which explains their higher interquartile range (IQR) and larger numbers of outliers.

3.2 Proposed Pre-Trained Models

3.2.1 OPUS-MT

OPUS-MT [45, 44] provides over 1,000 pre-trained models for translation between numerous language pairs. The architecture is based on MARIAN-NMT [21], based on a standard transformer setup: 6 self-attentive layers in both the encoder and decoder networks, each with 8 attention heads per layer [44].

While they provide both monolingual and multilingual PTMs, only the OPUS-MT monolingual models are used in this study.

3.2.2 mBART-50

mBART [25] is a sequence-to-sequence denoising auto-encoder model specifically designed for multilingual tasks. The mBART-50 variant supports many-to-many translations for over 50 languages.

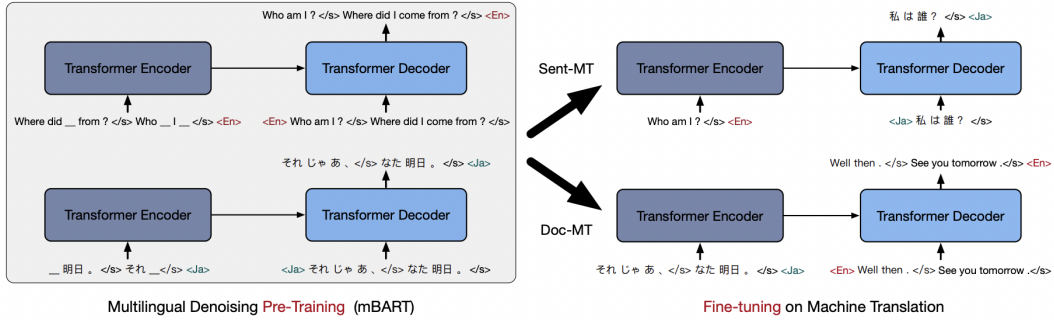


Figure 4: mBART framework for Multilingual Denoising Pre-training (left) and fine-tuning on downstream MT tasks (right) [25]

The mBART model is built as a typical sequence-to-sequence transformer architecture [46], using an embedding dimension of 1024 and 16 heads within 12 layers of encoder and 12 layers of decoder, totalling to around 680M parameters [25].

Figure 4 shows mBART’s key features: multilingual denoising pre-training and fine-tuning. The dataset is first corrupted by a noising function, which is applied to 35% of words of each span text [25]. During pre-training, the model will be made to denoise the texts, improving the model’s ability to generate coherent and accurate translations. After the pre-training step, the model is then able to be fine-tuned for downstream tasks. It is shown that the model achieved consistent performance gains through pre-training in low-to-medium resource sentence level MT [25]

3.2.3 M2M-100

M2M-100 [14] is designed to perform direct translation between any pair of 100 languages without relying on English as an intermediate language. Similarly, the architecture is of a transformer sequence-to-sequence with a typical encoder that outputs embeddings of the same length, while the decoder is an autoregressive component that sequentially produce the target sentence one token at a time. The parameters include 12 encoder layers and 12 decoder layers, 8192 FFN size and 1024 embedding dimension, totalling to a 1.2B parameter count [14].

Figure 5 show the inner working of the M2M model. The English-centric dataset (top left) includes training data exclusively involving translations to and from English, while the many-to-many multilingual setting (top right) involves direct translation data among multiple language pairs [14]. Finally, the model combines sparse language-specific parameters with dense parameters (bottom part of Figure 5),

3.2.4 NLLB-200

NLLB-200 [40] is built to handle translation tasks across a broad spectrum of languages, including many that are low-resource or underrepresented in existing datasets. It supports translations for 200 languages, including numerous underrepresented languages, and is currently one of the most extensive multilingual machine translation models. The model employs a Mixture of Experts (MoE) architecture

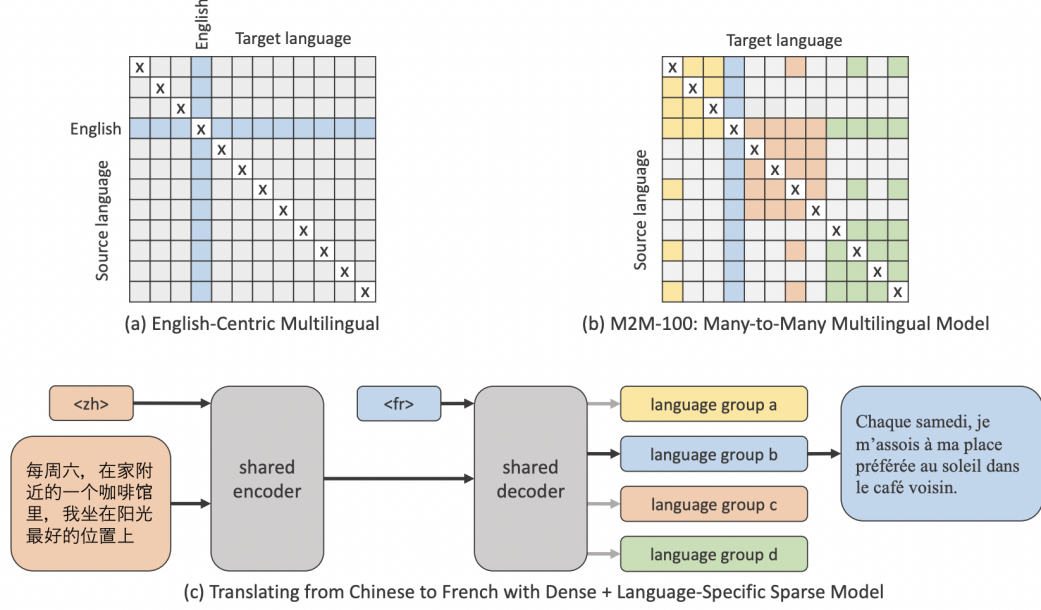


Figure 5: Summary of M2M dataset and multilingual model.

and achieves state-of-the-art (SoTA) results across many language pairs, even surpassing Meta’s previous model, M2M-100 [40].

The ‘nllb-200-distilled-600M’ variant has a total of 600M parameters, distilled from originally 54B parameters [40].

Figure 6 shows the Dense Transformer and MoE Transformer layers implemented within the model. MoE [28] utilises a gating mechanism to route different inputs to different subsets of experts (sub-networks), allowing the model to handle various linguistic phenomena efficiently. Additionally, the model is then trained on a diverse set of languages, complemented with back-translation and data augmentation to generate additional data for low-resource languages [40].

3.3 Evaluation Metrics

In this project, a variant of BLEU [33] called SacreBLEU [34] and METEOR [23] will be used in this experiment due to their popularity, simplicity, and ease-of-use.

While other metrics such as Crosslingual Optimized Metric for Evaluation of Translation (COMET) [37] and BERTScore [53] exist, they involve using a deep neural network or a transformer to evaluate, which naturally increase computational cost and slow down calculation greatly.

3.3.1 BLEU

Bilingual Evaluation Understudy (BLEU) [33] is the most commonly used metrics for machine translation (MT). It assesses how well a candidate translation matches the reference translation using precision metrics for n-grams and incorporates a brevity penalty to prevent overly short translations from achieving high scores.

The BLEU score is calculated as:

$$\text{BLEU} = BP \cdot \exp \left(\sum_{n=1}^N w_n \log p_n \right) \quad (1)$$

where BP is the brevity penalty, p_n is the precision for n-grams, and w_n is the weight for each n-gram (often uniformly distributed, so $w_n = \frac{1}{N}$).

The brevity penalty (BP) is calculated as:

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1 - \frac{r}{c})} & \text{if } c \leq r \end{cases} \quad (2)$$

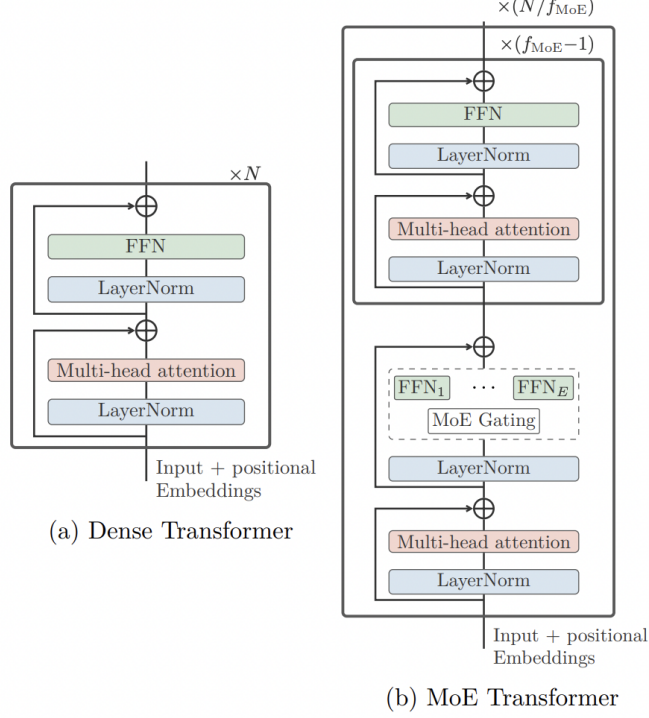


Figure 6: NLLB pipeline [40]

where c is the length of the candidate translation and r is the effective reference length.

The n-gram precision, as presented in the original BLEU paper [33], is calculated as:

$$p_n = \frac{\sum_{C \in \{Candidates\}} \sum_{n\text{-gram} \in C} \text{Count}_{\text{clip}}(n\text{-gram})}{\sum_{C \in \{Candidates\}} \sum_{n\text{-gram} \in C} \text{Count}(n\text{-gram})} \quad (3)$$

where p_n is the precision for n-grams, $\sum_{C \in \{Candidates\}}$ denotes the summation over all candidate translations, $\sum_{n\text{-gram} \in C}$ denotes the summation over all n-grams in a candidate translation C . $\text{Count}_{\text{clip}}(n\text{-gram})$ is the clipped count of the n-gram, which is the count of the n-gram in the candidate translation limited by the maximum count of that n-gram in any reference translation. $\text{Count}(n\text{-gram})$ is the count of the n-gram in the candidate translation.

The machine translation community’s rely heavily BLEU score, however, it has several drawbacks. The metric has been reported to not correlate strongly with human judgement, showing variations in translation that could mean that a higher BLEU score does not necessarily indicate a true enhancement in translation quality [8].

Furthermore, It is challenging to directly compare BLEU scores between paper [34]. Thus, the author proposed a standardised variant called SacreBLEU [34], designed to ensure easy, consistent, and comparable evaluation across different implementations and research studies.

3.3.2 METEOR

Metric for Evaluation of Translation with Explicit ORDERing (METEOR) [23] assesses a translation by calculating a score that reflects explicit word-to-word matches between the reference and a candidate translation [1]. It is designed to address some limitations of the BLEU score, allowing matches between simple morphological variants and synonyms. The formula is defined as:

$$\text{METEOR} = (1 - \gamma \cdot \text{frag}) \cdot \frac{P \cdot R}{\alpha \cdot P + (1 - \alpha) \cdot R} \quad (4)$$

where P is the precision, R is the recall, and frag is the fragmentation penalty. γ is a parameter that controls the weight of the fragmentation penalty, commonly set to 0.5, and α is a parameter that controls the balance between precision and recall, commonly set to 0.9,

3.4 Inference Details

OPUS-MT models are taken from the Helsinki-NLP repository. The 'mbart-large-50-many-to-many-mmmt' is used, 'm2m100_418M', and 'nllb-200-distilled-600M'. All models are taken using the transformers [50] library from Hugging Face [13]. Batch size of 4 is used in all models. Translations are first generated and saved into CSV files, before calculating the metrics score.

The sacrebleu packages from pip is used to calculate SacreBLEU [34]. For BLEU and METEOR, the NLTK [7] package is used. Corpus BLEU is implemented for both SacreBLEU and BLEU in the calculations. Additionally, smoothing method 7 is used for BLEU, which interpolates two other smoothing functions: smoothing method 4, which divides by $\frac{1}{\ln(\ln(T))}$, where T is the length of the translation; and smoothing method 5, which averages the $n - 1$, n , and $n + 1$ gram matched counts [41].

All codes are written in Python 3.12.2, run locally on MacOS Sonoma 14.4.1 with M2 chip and 16 GB memory. GitHub: <https://github.com/stefanliemawan/eval-machine-translation>

4 Evaluation

Language	BLEU	SacreBLEU	METEOR
Chinese	0.0971	<i>59.5098</i>	0.7953
Dutch	0.0974	69.8803	0.8471
Finnish	0.0965	66.6267	0.8296
French	0.0967	69.8185	0.8357
German	0.0969	69.7422	0.8419
Hebrew	0.0977	66.5149	0.8229
Italian	0.0979	74.1298	0.8584
Japanese	0.0960	63.1435	<i>0.7893</i>
Polish	<i>0.0952</i>	61.9026	0.8425
Russian	0.0974	66.7046	0.8179
Spanish	0.0979	71.4174	0.8463
Turkish	0.0980	72.6551	0.8460
Ukrainian	0.0978	75.4447	0.8667

Table 3: OPUS-MT result.

Language	BLEU	SacreBLEU	METEOR
Chinese	0.0962	54.7322	0.7600
Dutch	0.0970	63.6482	0.8007
Finnish	0.0963	47.2194	0.7067
French	0.0968	57.2482	0.7598
German	0.0967	63.1666	0.8022
Hebrew	0.0972	58.0846	0.7590
Italian	0.0972	65.9415	0.8068
Japanese	0.0943	43.9547	0.7151
Polish	0.0978	62.7550	0.7923
Russian	0.0974	58.8820	0.7686
Spanish	<i>0.0883</i>	<i>35.1593</i>	0.7306
Turkish	0.0975	50.9377	<i>0.6982</i>
Ukrainian	0.0964	55.8637	0.7496

Table 4: mBART50 result.

Language	BLEU	SacreBLEU	METEOR
Chinese	0.0975	43.3490	0.6848
Dutch	0.0976	52.6920	0.7448
Finnish	0.0975	49.9149	0.7145
French	0.0969	51.4090	0.7288
German	0.0969	52.7877	0.7423
Hebrew	0.0981	50.3287	0.7198
Italian	0.0975	53.1644	0.7341
Japanese	<i>0.0959</i>	<i>41.6610</i>	<i>0.6511</i>
Polish	0.0984	52.5151	0.7301
Russian	0.0976	48.2725	0.7027
Spanish	0.0977	52.9323	0.7371
Turkish	0.0978	51.3009	0.7262
Ukrainian	0.0978	46.5555	0.6890

Table 5: M2M-100 result.

Language	BLEU	SacreBLEU	METEOR
Chinese	0.0949	50.4454	0.7180
Dutch	0.0915	20.8988	0.2899
Finnish	0.0921	32.3926	0.4006
French	0.0950	32.1389	0.4338
German	0.0903	7.1436	0.1528
Hebrew	0.0897	2.3455	0.1412
Italian	0.0926	24.2702	0.3059
Japanese	<i>0.0842</i>	<i>1.9569</i>	<i>0.0919</i>
Polish	0.0937	38.4395	0.4739
Russian	0.0918	20.6734	0.2977
Spanish	0.0929	26.3313	0.3578
Turkish	0.0945	45.2721	0.6187
Ukrainian	0.0914	9.0340	0.1613

Table 6: NLLB-200 result.

Table 3, 4, 5, and 6, and shows the result from all four PTMs. All models perform similarly low on BLEU score across all languages, with NLLB-200 reaching 0.01 below more on Japanese translation. The BLEU scores are very low (below 10 BLEU) in all cases due to the sentences being short (find cite). This also highlights the superiority of SacreBLEU compared to the flawed, standard BLEU.

OPUS-MT model shows the best performance in all languages, consistently reaching mostly above 60 SacreBLEU and 0.8 METEOR. The highest-rated languages are Ukrainian-to-English translation with 75.44 SacreBLEU and 0.86 METEOR, while the lowest scores are achieved on Chinese-to-English translation: 59.50 SacreBLEU and Japanese-to-English: 0.78 METEOR.

The mBART-50 model is the second-best-performing PTMs and the best-performing multilingual model, achieving best SacreBLEU and METEOR score in Italian-to-English translation, while suffering the most in Spanish-to-English (SacreBLEU) and Turkish-to-English (METEOR) translation. Interestingly, Spanish received the worst score in SacreBLEU while it performs decently in METEOR, and

Turkish performs well in SacreBLEU score but has the worst in METEOR. Additionally, the model slightly outperforms OPUS-MT in Polish-to-English translation on SacreBLEU score (0.85 higher).

Similarly, SacreBLEU for Italian-to-English translation outperforms other languages with the M2M-100 model, although Dutch beat Italian by 0.01 in METEOR score. Japanese received the worst score in both SacreBLEU and METEOR for this model. M2M-100 performs worse in SacreBLEU compared to mBART-50, but comparably similar in METEOR scores.

NLLB-200 scores much lower compared to other models in both SacreBLEU and METEOR in almost all languages. It scores highest in Chinese-to-English translation in all metrics, while showing significant variability among other languages. However, Japanese-to-English translation suffers with extremely low scores of 1.95 SacreBLEU and 0.09 METEOR. Additionally, German, Hebrew, and Ukrainian also have low SacreBLEU and METEOR scores, among the other languages.

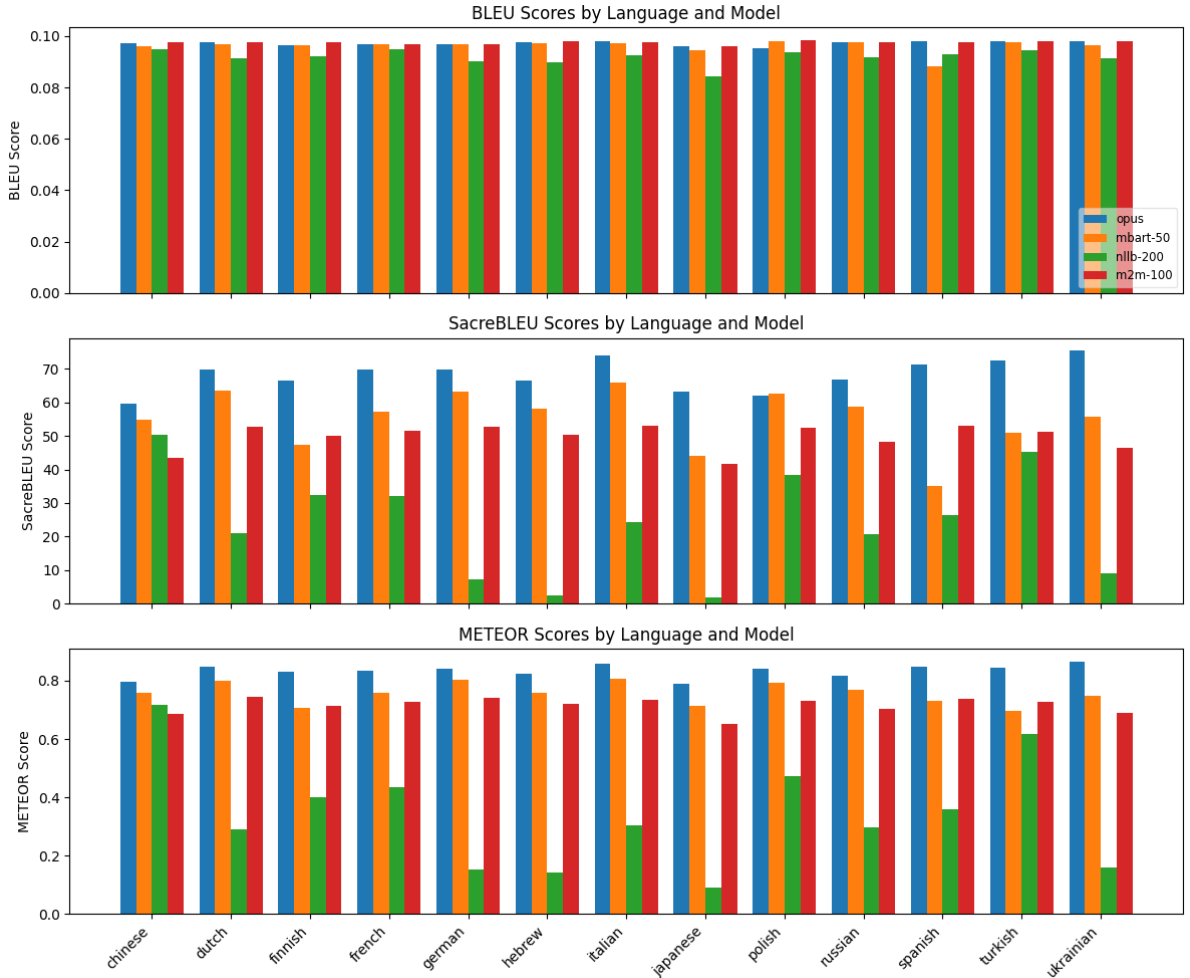


Figure 7: BLEU, SacreBLEU, and METEOR scores of each model on each language, divided by language

Figure 7 shows the visualisation of all scores, categorised per language.

Dutch and Italian translations generally score higher across all metrics, suggesting that these languages are well-supported by the pre-trained models. Japanese translations, however, pose significant challenges for all models, as indicated by the consistently low scores across all metrics.

Japanese-to-English translation seem to have the worst score on average across the four models, but particularly bad on multilingual models. In contrast, Dutch and German, have consistently high results on all models except for NLLB-200, often reaching closely to the highest-language-score SacreBLEU and METEOR score.

Figure 8 shows the SacreBLEU and METEOR scores for each language by model. The OPUS-MT and mBART-50 models are again shown as the clear winners in almost all languages. M2M-100 also performs competitively with mBART-50 in some languages, surpassing mBART-50 in Spanish-to-English translation. All four models seem to perform consistently against each other in Chinese-to-English and

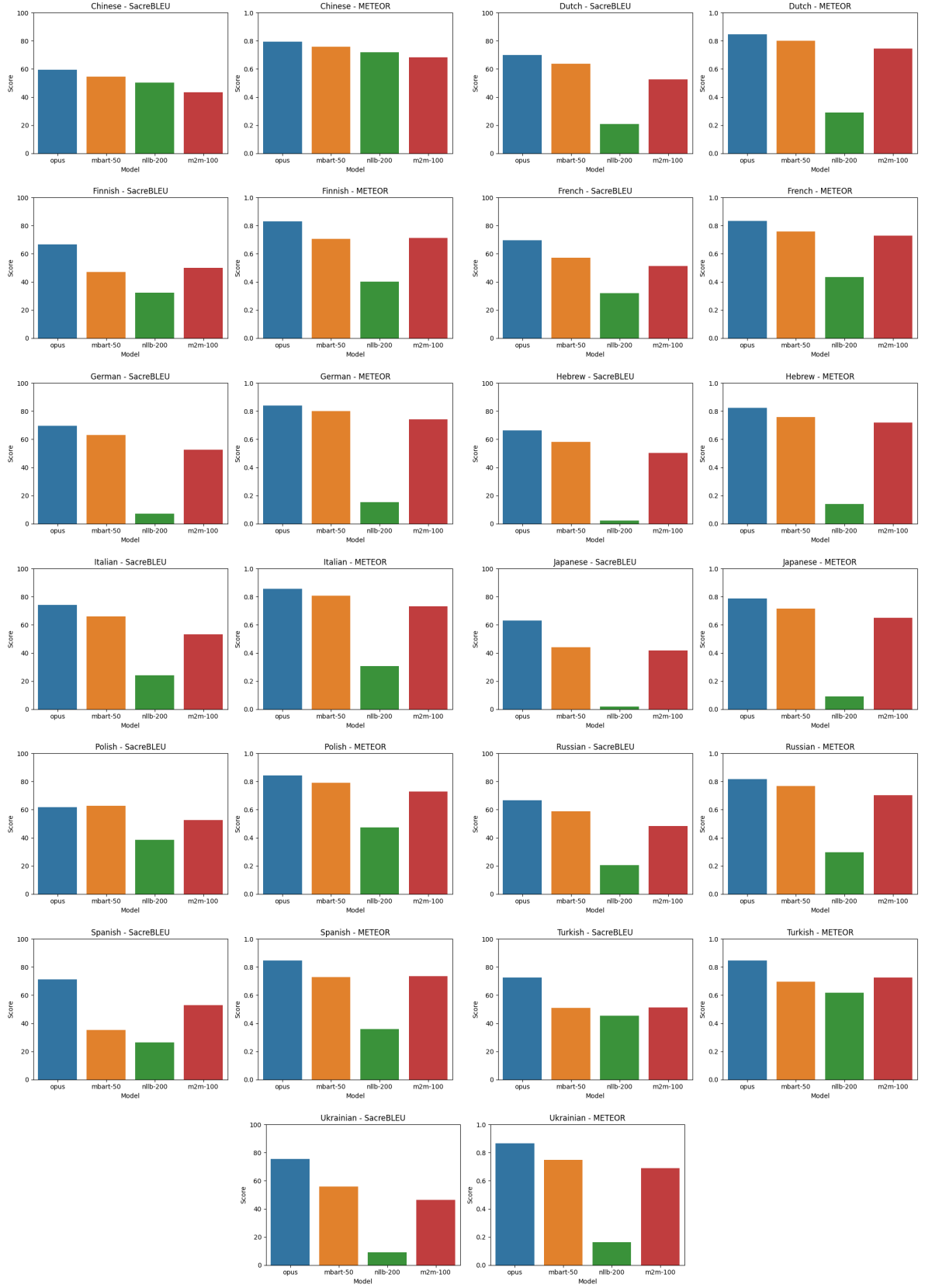


Figure 8: Performance of every language.

Turkish-to-English translations.

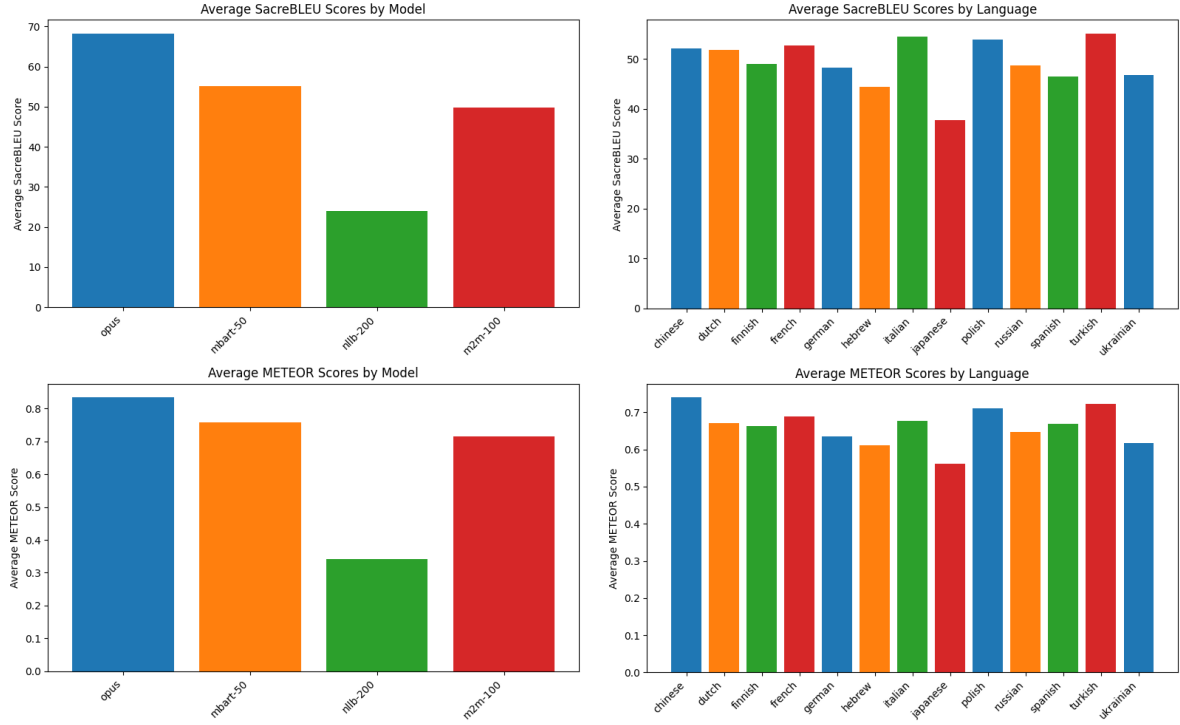


Figure 9: The average BLEU, SacreBLEU, and METEOR scores by model (left) and by language (right)

Figure 9 visualise the average BLEU, SacreBLEU, and METEOR, for each model and each language, highlighting how the models and languages perform independently. It can be seen that mBART-50 perform more closely to OPUS-MT on METEOR score than SacreBLEU score. Language-wise, Italian, Polish, and Turkish achieves the best average SacreBLEU scores while Chinese, Polish and Turkish are shown as the top three in average METEOR scores.

Figure 10 visualises the difference between each OPUS-MT models reported BLEU scores compared to the SacreBLEU result of this study. The reported BLEU is evaluated on the Tatoeba dataset [44]. It can be seen that the scores are quite aligned, with Ukrainian and Italian reaching the two highest score while Chinese, Polish, and Japanese suffers.

This study experiments also resulted in moderately higher overall SacreBLEU score, however, the dataset used here is also significantly less in size than the full Tatoeba challenge dataset [43].

Language	SacreBLEU	Reported BLEU
Chinese	59.5	36.1
Dutch	69.8	60.9
Finnish	66.6	53.4
French	69.8	57.5
German	69.7	55.4
Hebrew	66.5	52.0
Italian	74.1	70.9
Japanese	63.1	41.7
Polish	61.9	54.9
Russian	66.7	61.1
Spanish	71.4	59.6
Turkish	72.6	63.5
Ukrainian	75.4	64.1

Table 7: OPUS-MT result compared with reported BLEU result

Additionally, Table 7 shows the full scores of both results SacreBLEU and reported BLEU from each OPUS-MT models. The reported BLEU scores are taken from each model cards in the HuggingFace

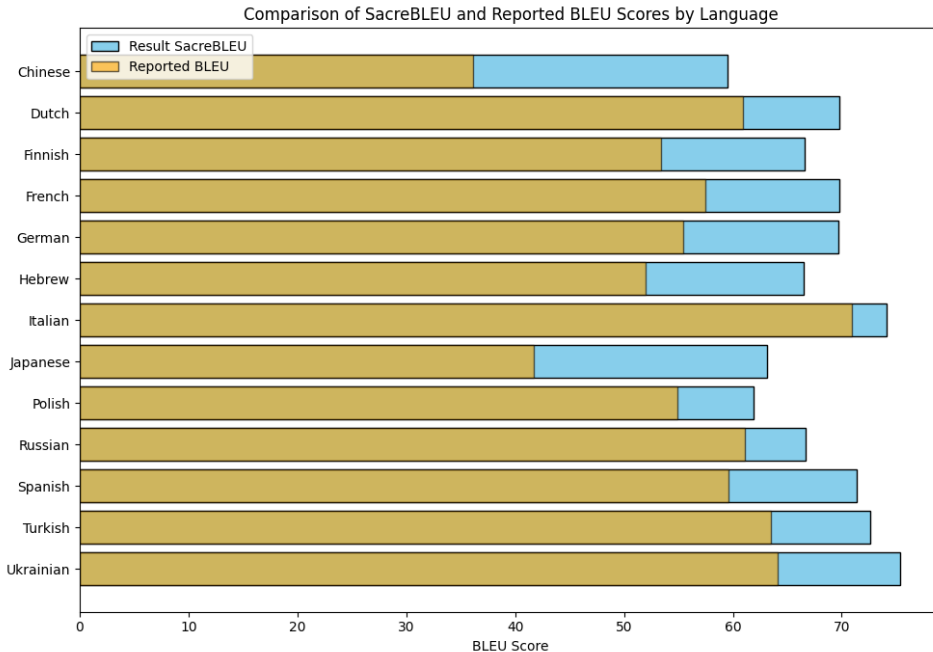


Figure 10: Comparison between reported OPUS-MT models performance and performance in this experiment.

repository [13].

Lastly, Figure 11 shows the time taken to translate for each language across experimented models. The OPUS-MT model is the fastest among the other three multilingual models, which perform similarly in timings. It is quite unclear why mBART50 Spanish-to-English translation suffers from a significantly higher timing. Re-running the experiment shows that this timing is consistent and not an anomaly.

5 Conclusion

This study findings show that while multilingual models have been shown to achieve higher performance compared to one-to-one translation models [25], this does not hold on short sentences, many-to-English translations across the 14 source languages.

OPUS-MT achieves the best performance across all 14 languages. mBART-50 is the second-best-performing model, beating OPUS-MT only in Polish SacreBLEU. Nevertheless, both mBART-50 and M2M-100 performs generally well on all languages. NLLB-200, however, achieves much lower scores in both SacreBLEU and METEOR for all languages except for Chinese, slightly outperforming M2M-100 but still slightly below mBART-50.

NLLB-200 should then be avoided for use without any further fine-tuning on these 14 languages.

Chinese, Italian, Polish, Turkish have the best average scores over all four models, Ukrainian by OPUS-MT has the highest SacreBLEU and METEOR score across all the other models and languages.

Limitations, only the first translation of the same phrases are taken, which may disrupts the possibility for synonyms, especially for BLEU and SacreBLEU metrics. Additionally, no fine-tuning is applied in this experiment, which may affect the translation results, as fine-tuning models for specific languages can significantly improve models performance [54].

6 Discussion

Future works should be done by incorporating more Pre-Trained Models (PTMs) and implementing fine-tuning for specific languages on languages-to-English translation. Evaluating with bigger datasets and longer texts can also be beneficial, as well as using language other than English as the target language.

Furthermore, it may be beneficial to consider more general multilingual PTMs such as mBERT [51] and PolyLM [48], which are trained for a wide arrange of NLP tasks instead of just machine translation.

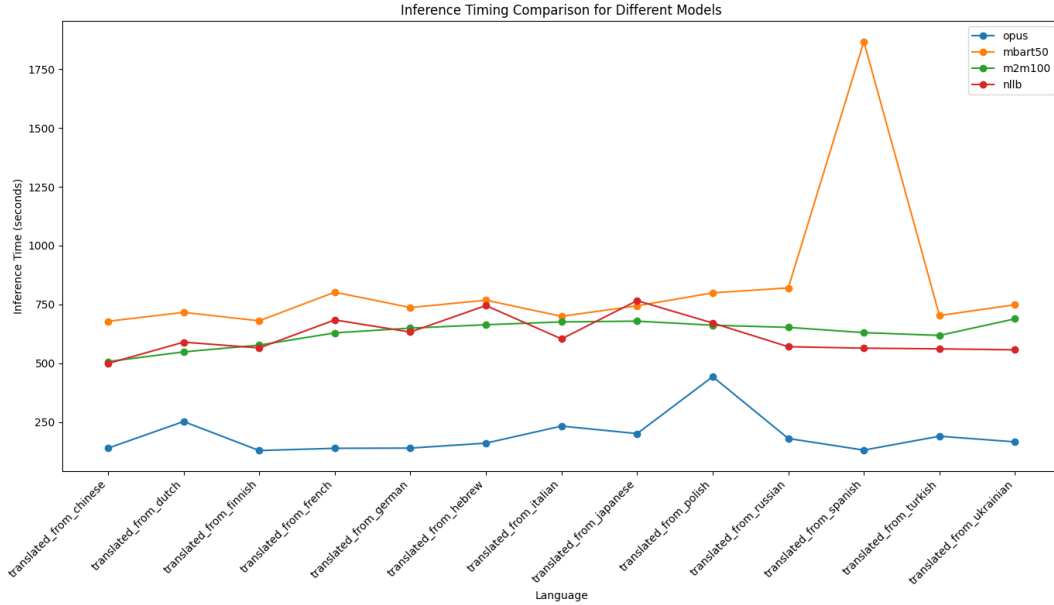


Figure 11: Inference timings for every model.

References

- [1] Abhaya Agarwal and Alon Lavie. “METEOR, M-BLEU and M-TER: evaluation metrics for high-correlation with human rankings of machine translation output”. In: *Proceedings of the Third Workshop on Statistical Machine Translation*. StatMT ’08. Columbus, Ohio: Association for Computational Linguistics, 2008, 115–118. ISBN: 9781932432091.
- [2] Milind Agarwal et al. “Findings of the IWSLT 2023 Evaluation Campaign”. In: *Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023)*. 20th International Conference on Spoken Language Translation. IWSLT 2023 (Toronto, Kanada, July 13–14, 2023). Association for Computational Linguistics (ACL), 2023, 1–61.
- [3] Roei Aharoni, Melvin Johnson, and Orhan Firat. *Massively Multilingual Neural Machine Translation*. 2019. arXiv: 1903.00089 [cs.CL]. URL: <https://arxiv.org/abs/1903.00089>.
- [4] Hanna Pięta Alexandra Assis Rosa and Rita Bueno Maia. “Theoretical, methodological and terminological issues regarding indirect translation: An overview”. In: *Translation Studies* 10.2 (2017), pp. 113–132. DOI: 10.1080/14781700.2017.1285247. eprint: <https://doi.org/10.1080/14781700.2017.1285247>. URL: <https://doi.org/10.1080/14781700.2017.1285247>.
- [5] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. *Neural Machine Translation by Jointly Learning to Align and Translate*. 2016. arXiv: 1409.0473 [cs.CL]. URL: <https://arxiv.org/abs/1409.0473>.
- [6] Loïc Barrault et al. “Findings of the 2020 Conference on Machine Translation (WMT20)”. In: *Proceedings of the Fifth Conference on Machine Translation*. Ed. by Loïc Barrault et al. Online: Association for Computational Linguistics, Nov. 2020, pp. 1–55. URL: <https://aclanthology.org/2020.wmt-1.1>.
- [7] Steven Bird, Edward Loper, and Ewan Klein. *Natural Language Processing with Python*. O’Reilly Media Inc., 2009.
- [8] Chris Callison-Burch, Miles Osborne, and Philipp Koehn. “Re-evaluating the Role of Bleu in Machine Translation Research”. In: *11th Conference of the European Chapter of the Association for Computational Linguistics*. Ed. by Diana McCarthy and Shuly Wintner. Trento, Italy: Association for Computational Linguistics, Apr. 2006, pp. 249–256. URL: <https://aclanthology.org/E06-1032>.

- [9] Kyunghyun Cho et al. “On the Properties of Neural Machine Translation: Encoder–Decoder Approaches”. In: *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*. Ed. by Dekai Wu et al. Doha, Qatar: Association for Computational Linguistics, Oct. 2014, pp. 103–111. DOI: 10.3115/v1/W14-4012. URL: <https://aclanthology.org/W14-4012>.
- [10] Asa Cooper Stickland, Xian Li, and Marjan Ghazvininejad. “Recipes for Adapting Pre-trained Monolingual and Multilingual Models to Machine Translation”. In: *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Association for Computational Linguistics, 2021. DOI: 10.18653/v1/2021.eacl-main.301. URL: <http://dx.doi.org/10.18653/v1/2021.eacl-main.301>.
- [11] Jacob Devlin et al. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. 2019. arXiv: 1810.04805 [cs.CL]. URL: <https://arxiv.org/abs/1810.04805>.
- [12] David M. Eberhard, Gary F. Simons, and Charles D. Fennig, eds. *Ethnologue: Languages of the World*. 27th. Dallas, Texas: SIL International, 2024. URL: <http://www.ethnologue.com>.
- [13] Hugging Face. *Hugging Face: Natural Language Processing Made Easy*. <https://huggingface.co>. Accessed: 2024-08-07. 2024.
- [14] Angela Fan et al. *Beyond English-Centric Multilingual Machine Translation*. 2020. arXiv: 2010.11125 [cs.CL]. URL: <https://arxiv.org/abs/2010.11125>.
- [15] Xavier Garcia et al. *A Multilingual View of Unsupervised Machine Translation*. 2020. arXiv: 2002.02955 [cs.CL]. URL: <https://arxiv.org/abs/2002.02955>.
- [16] Xu Han et al. “Pre-trained models: Past, present and future”. In: *AI Open 2* (2021), pp. 225–250. ISSN: 2666-6510. DOI: <https://doi.org/10.1016/j.aiopen.2021.08.002>. URL: <https://www.sciencedirect.com/science/article/pii/S2666651021000231>.
- [17] John Hutchins. “Research methods and system designs in machine translation: a ten-year review, 1984-1994”. In: *BCS International Academic Conference*. 1994. URL: <https://api.semanticscholar.org/CorpusID:15952756>.
- [18] John Hutchins. “The development and use of machine translation systems and computer-based translation tools in Europe, Asia, and North America”. In: 1998. URL: <https://api.semanticscholar.org/CorpusID:18918684>.
- [19] John Hutchins. “The first public demonstration of machine translation : the Georgetown-IBM system, 7th January 1954”. In: 2006. URL: <https://api.semanticscholar.org/CorpusID:132677>.
- [20] William J. Hutchins. “Machine translation over fifty years”. In: 2001. URL: <https://api.semanticscholar.org/CorpusID:6196527>.
- [21] Marcin Junczys-Dowmunt et al. “Marian: Fast Neural Machine Translation in C++”. In: *Proceedings of ACL 2018, System Demonstrations*. Melbourne, Australia: Association for Computational Linguistics, 2018, pp. 116–121. URL: <http://www.aclweb.org/anthology/P18-4020>.
- [22] Philipp Koehn and Rebecca Knowles. *Six Challenges for Neural Machine Translation*. 2017. arXiv: 1706.03872 [cs.CL]. URL: <https://arxiv.org/abs/1706.03872>.
- [23] Alon Lavie and Abhaya Agarwal. “Meteor: an automatic metric for MT evaluation with high levels of correlation with human judgments”. In: *Proceedings of the Second Workshop on Statistical Machine Translation*. StatMT ’07. Association for Computational Linguistics, 2007, 228–231.
- [24] Marie-Aude Lefer. “Parallel Corpora”. In: *A Practical Handbook of Corpus Linguistics*. Ed. by Magali Paquot and Stefan Th. Gries. Cham: Springer International Publishing, 2020, pp. 257–282. ISBN: 978-3-030-46216-1. DOI: 10.1007/978-3-030-46216-1_12. URL: https://doi.org/10.1007/978-3-030-46216-1_12.
- [25] Yinhan Liu et al. *Multilingual Denoising Pre-training for Neural Machine Translation*. 2020. arXiv: 2001.08210 [cs.CL]. URL: <https://arxiv.org/abs/2001.08210>.
- [26] Adam Lopez. “Statistical machine translation”. In: *ACM Comput. Surv.* 40.3 (Aug. 2008). ISSN: 0360-0300. DOI: 10.1145/1380584.1380586. URL: <https://doi.org/10.1145/1380584.1380586>.
- [27] Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. *Effective Approaches to Attention-based Neural Machine Translation*. 2015. arXiv: 1508.04025 [cs.CL]. URL: <https://arxiv.org/abs/1508.04025>.

- [28] Saeed Masoudnia and Reza Ebrahimpour. “Mixture of experts: a literature survey”. In: *Artificial Intelligence Review* 42.2 (May 2012), 275–293. ISSN: 1573-7462. DOI: 10.1007/s10462-012-9338-y. URL: <http://dx.doi.org/10.1007/s10462-012-9338-y>.
- [29] Margaret Dumebi Okpor. “Machine Translation Approaches: Issues and Challenges”. In: 2014. URL: <https://api.semanticscholar.org/CorpusID:11483090>.
- [30] OpenAI et al. *GPT-4 Technical Report*. 2024. arXiv: 2303.08774 [cs.CL]. URL: <https://arxiv.org/abs/2303.08774>.
- [31] OPUS. *OPUS: The Open Parallel Corpus*. <https://opus.nlpl.eu/>. Accessed: 2024-07-29. 2024.
- [32] Myle Ott et al. “fairseq: A Fast, Extensible Toolkit for Sequence Modeling”. In: *Proceedings of NAACL-HLT 2019: Demonstrations*. 2019.
- [33] Kishore Papineni et al. “BLEU: a method for automatic evaluation of machine translation”. In: *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*. ACL ’02. Association for Computational Linguistics, 2002, 311–318. DOI: 10.3115/1073083.1073135. URL: <https://doi.org/10.3115/1073083.1073135>.
- [34] Matt Post. *A Call for Clarity in Reporting BLEU Scores*. 2018. arXiv: 1804.08771 [cs.CL]. URL: <https://arxiv.org/abs/1804.08771>.
- [35] Colin Raffel et al. *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*. 2023. arXiv: 1910.10683 [cs.LG]. URL: <https://arxiv.org/abs/1910.10683>.
- [36] Surangika Ranathunga et al. “Neural Machine Translation for Low-resource Languages: A Survey”. In: *ACM Comput. Surv.* 55.11 (2023). ISSN: 0360-0300. DOI: 10.1145/3567592. URL: <https://doi.org/10.1145/3567592>.
- [37] Ricardo Rei et al. *COMET: A Neural Framework for MT Evaluation*. 2020. arXiv: 2009.09025 [cs.CL]. URL: <https://arxiv.org/abs/2009.09025>.
- [38] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. *Sequence to Sequence Learning with Neural Networks*. 2014. arXiv: 1409.3215 [cs.CL]. URL: <https://arxiv.org/abs/1409.3215>.
- [39] Tatoeba Community. *About Tatoeba*. Accessed: 2024-07-06. 2024. URL: <https://tatoeba.org/en>.
- [40] NLLB Team et al. *No Language Left Behind: Scaling Human-Centered Machine Translation*. 2022. arXiv: 2207.04672 [cs.CL]. URL: <https://arxiv.org/abs/2207.04672>.
- [41] NLTK Team. *NLTK Documentation*. 2023. URL: https://www.nltk.org/api/nltk.translate.bleu_score.html.
- [42] *The State of Machine Translation 2020*. Independent multi-domain evaluation of commercial Machine Translation engines. Intento, 2020. URL: https://try.inten.to/mt_report_2020.
- [43] Jörg Tiedemann. “The Tatoeba Translation Challenge – Realistic Data Sets for Low Resource and Multilingual MT”. In: *Proceedings of the Fifth Conference on Machine Translation*. Online: Association for Computational Linguistics, Nov. 2020, pp. 1174–1182. URL: <https://www.aclweb.org/anthology/2020.wmt-1.139>.
- [44] Jörg Tiedemann and Santhosh Thottingal. “OPUS-MT — Building open translation services for the World”. In: *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation (EAMT)*. Lisbon, Portugal, 2020.
- [45] Jörg Tiedemann et al. “Democratizing neural machine translation with OPUS-MT”. In: *Language Resources and Evaluation* 58.2 (Dec. 2023), 713–755. ISSN: 1574-0218. DOI: 10.1007/s10579-023-09704-w. URL: <http://dx.doi.org/10.1007/s10579-023-09704-w>.
- [46] Ashish Vaswani et al. “Attention is All you Need”. In: *Advances in Neural Information Processing Systems*. Ed. by I. Guyon et al. Vol. 30. Curran Associates, Inc., 2017. URL: https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- [47] Warren Weaver. “Memorandum on Translation”. In: *MT News International* 22 (1999), pp. 5–6, 15.
- [48] Xiangpeng Wei et al. *PolyLM: An Open Source Polyglot Large Language Model*. 2023. arXiv: 2307.06018 [cs.CL]. URL: <https://arxiv.org/abs/2307.06018>.

- [49] Wikipedia contributors. *Google Translate — Wikipedia, The Free Encyclopedia*. [Online; accessed 6-July-2024]. 2024. URL: https://en.wikipedia.org/w/index.php?title=Google_Translate&oldid=1232822378.
- [50] Thomas Wolf et al. “Transformers: State-of-the-Art Natural Language Processing”. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Online: Association for Computational Linguistics, Oct. 2020, pp. 38–45. URL: <https://www.aclweb.org/anthology/2020.emnlp-demos.6>.
- [51] Shijie Wu and Mark Dredze. “Are All Languages Created Equal in Multilingual BERT?” In: *Proceedings of the 5th Workshop on Representation Learning for NLP*. Ed. by Spandana Gella et al. Online: Association for Computational Linguistics, July 2020, pp. 120–130. DOI: 10.18653/v1/2020.repl4nlp-1.16. URL: <https://aclanthology.org/2020.repl4nlp-1.16>.
- [52] Qi Ye et al. “When and Why are pre-trained word embeddings useful for Neural Machine Translation”. In: *HLT-NAACL*. 2018.
- [53] Tianyi Zhang et al. *BERTScore: Evaluating Text Generation with BERT*. 2020. arXiv: 1904.09675 [cs.CL]. URL: <https://arxiv.org/abs/1904.09675>.
- [54] Xuan Zhang et al. “Machine Translation with Large Language Models: Prompting, Few-shot Learning, and Fine-tuning with QLoRA”. In: *Proceedings of the Eighth Conference on Machine Translation*. Ed. by Philipp Koehn et al. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 468–481. DOI: 10.18653/v1/2023.wmt-1.43. URL: <https://aclanthology.org/2023.wmt-1.43>.