

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

Stefan Liemawan Adjii

August 4, 2024

## 1 Introduction

According to Ethnologue [7], 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 [35]). In modern times, the need for translation services has surged due to the growing exchange of information across different regions that speak various languages [18].

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 [33], 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 [12]. 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' [13]. 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) [10, 11]. Subsequently, statistics-based approaches for MT re-emerged, bolstered by the recent success of probabilistic techniques in speech recognition. Statistical machine translation [17] dominated the domain between late 1990s through the early 2010s, before largely being surpassed by neural machine translation (NMT) [4, 26].

Since the introduction of Transformers in 2017 [32], 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 [6], GPT [19], and T5 [23], 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, namely mBART [16], mT5 [36], NLLB [28], M2M [8], and PolyLM [34]. Most of these models are multilingually trained, 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 [16].

Intento published 'The State of Machine Translation 2024' [29] 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.

Despite these advancements, pre-trained models are often evaluated using different set of benchmarks [16, 28, 8, 34], 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.

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 [5], no pre-training or fine-tuning is performed in this study for simplicity reasons. A dataset is curated from the Tatoeba repository [27], containing 1,241 parallel sentence pairs across source and target languages. The models includes a one-to-one PTMs: OPUS-MT [31], and multilingual PTMs: such as mBART-50 [16], NLLB-200 [28], and M2M-100 [8]. The performance of these models is evaluated using the BLEU score [22].

## 2 Literature Review

### 2.1 Parallel Corpora

Since Neural Machine Translation (NMT) systems require vast amounts of training data, the availability of parallel corpora is crucial for building effective models [14]. The lack of extensive parallel corpora, especially for low-resource languages, leads to suboptimal performance in NMT techniques compared to their high-resource counterparts [24]. OPUS [20] 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 [27] 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.

The Tatoeba Challenge [30], IWSLT [1], TED [37], Flores [9], and WMT [3] are among the biggest and most popular parallel corpora datasets, allowing for many-to-many translations.

Despite being used diligently by Meta AI papers for benchmarking [8, 28], these datasets are not seen used

### 2.2 Pre-Trained Models for Machine Translation

The encoder-decoder approach [4] remains as the foundation architecture for many sequence-to-sequence models in machine translation.

In terms of pre-trained models (PTMs) for machine translation, it can be divided into two categories: one-to-one models and many-to-many models (multilingual).

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. An example of this setup is OPUS-MT by Helsinki-NLP [31], which provides over 1,000 pre-trained models for translation between numerous language pairs.

With the advent of large language models and pre-trained language models, multilingual machine translation has gained prominence. This approach enables many-to-many translation, where a single model can translate between multiple source and target languages [2]. mBART [16] 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. M2M-100 [8] is designed to perform direct translation between 100 languages without relying on English as an intermediate language. NLLB-200: NLLB-200 [28] 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, encompassing numerous underrepresented languages. mBART, M2M-100, and NLLB-200 were all developed by Meta AI (formerly Facebook AI), showcasing the organisation’s significant impact in the machine translation field through pre-trained models (PTMs). These models represent a substantial advancement in multilingual translation capabilities.

mRASP2 [21]

### 2.3 Evaluation Metrics

Bilingual Evaluation Understudy (BLEU) [22] 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 n-gram precision, as presented in the original BLEU paper [22], is calculated as:

$$p_n = \frac{\sum_{C \in \{Candidates\}} \sum_{n\text{-gram} \in C} \text{Count}_{\text{clip}}(n\text{-gram})}{\sum_{C \in \{Candidates\}} \sum_{n\text{-gram} \in C} \text{Count}(n\text{-gram})} \quad (1)$$

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  $C$ .

- $\text{Count}_{\text{clip}}(\text{n-gram})$  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.
- $\text{Count}(\text{n-gram})$  is the count of the n-gram in the candidate translation.

Thus, the BLEU score is calculated as:

$$\text{BLEU} = BP \cdot \exp \left( \sum_{n=1}^N w_n \log p_n \right) \quad (2)$$

Where:

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

The brevity penalty (BP) is calculated as:

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1-\frac{r}{c})} & \text{if } c \leq r \end{cases} \quad (3)$$

Where:

- $c$  is the length of the candidate translation.
- $r$  is the effective reference length.

While other metrics exist for machine translation such as METEOR [15] and COMET [25], it is far more practical to implement BLEU due to its popular usage in other works.

## 3 Experiments

### 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 [27]. 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.

| No. | Language         |
|-----|------------------|
| 1   | Dutch            |
| 2   | Finnish          |
| 3   | French           |
| 4   | German           |
| 5   | Hebrew           |
| 6   | Hungarian        |
| 7   | Italian          |
| 8   | Japanese         |
| 9   | Mandarin Chinese |
| 10  | Polish           |
| 11  | Russian          |
| 12  | Spanish          |
| 13  | Turkish          |
| 14  | Ukrainian        |

Table 1: List of chosen languages for evaluation



Figure 1: Tatoeba's languages repository with 10,000+ sentences and 100,000+ sentences [27]

|    |  | 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 [27]

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. 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 1241 sentences in all 14 languages, this will be treated as a test set to evaluate the models performance on each language.

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!' 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 devices such as metaphors or sarcasm. Therefore, machine translation process should be straightforward on this level.

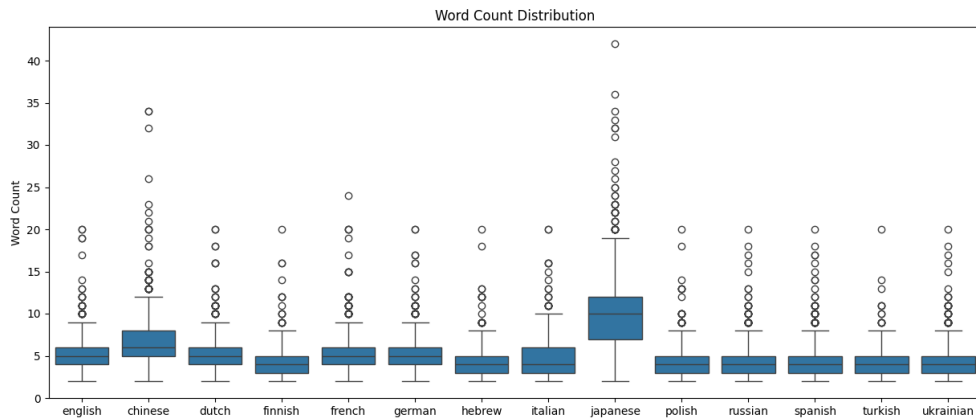


Figure 3: Dataset word count distribution, consisting of mostly short sentences.

Figure 3 shows a box plot of sentences word count. It can be seen that the majority of sentences are short.

Chinese and Japanese are counted per letter.

While most languages cluster around 4-6 words per sentence, there are notable exceptions like Japanese, which exhibits much less variability. This analysis can be useful for understanding language-specific characteristics in sentence structure, which could inform tasks like translation, text processing, or linguistic studies.

## 4 Evaluation

## 5 Conclusion

## References

- [1] 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.
- [2] 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>.
- [3] 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>.

- [4] 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>.
- [5] 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>.
- [6] 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>.
- [7] 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>.
- [8] Angela Fan et al. *Beyond English-Centric Multilingual Machine Translation*. 2020. arXiv: 2010.11125 [cs.CL]. URL: <https://arxiv.org/abs/2010.11125>.
- [9] Naman Goyal et al. “The Flores-101 Evaluation Benchmark for Low-Resource and Multilingual Machine Translation”. In: *Transactions of the Association for Computational Linguistics* 10 (May 2022), pp. 522–538. ISSN: 2307-387X. DOI: 10.1162/tac1\_a\_00474. eprint: [https://direct.mit.edu/tac1/article-pdf/doi/10.1162/tac1\\_a\\_00474/2020699/tac1\\_a\\_00474.pdf](https://direct.mit.edu/tac1/article-pdf/doi/10.1162/tac1_a_00474/2020699/tac1_a_00474.pdf). URL: [https://doi.org/10.1162/tac1\\_a\\_00474](https://doi.org/10.1162/tac1_a_00474).
- [10] 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>.
- [11] 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>.
- [12] 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>.
- [13] William J. Hutchins. “Machine translation over fifty years”. In: 2001. URL: <https://api.semanticscholar.org/CorpusID:6196527>.
- [14] Philipp Koehn and Rebecca Knowles. *Six Challenges for Neural Machine Translation*. 2017. arXiv: 1706.03872 [cs.CL]. URL: <https://arxiv.org/abs/1706.03872>.
- [15] 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.
- [16] 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>.
- [17] Adam Lopez. “Statistical machine translation”. In: *ACM Comput. Surv.* 40.3 (2008). ISSN: 0360-0300. DOI: 10.1145/1380584.1380586. URL: <https://doi.org/10.1145/1380584.1380586>.
- [18] Margaret Dumebi Okpor. “Machine Translation Approaches: Issues and Challenges”. In: 2014. URL: <https://api.semanticscholar.org/CorpusID:11483090>.
- [19] OpenAI et al. *GPT-4 Technical Report*. 2024. arXiv: 2303.08774 [cs.CL]. URL: <https://arxiv.org/abs/2303.08774>.
- [20] OPUS. *OPUS: The Open Parallel Corpus*. <https://opus.nlpl.eu/>. Accessed: 2024-07-29. 2024.
- [21] Xiao Pan et al. *Contrastive Learning for Many-to-many Multilingual Neural Machine Translation*. 2021. arXiv: 2105.09501 [cs.CL]. URL: <https://arxiv.org/abs/2105.09501>.
- [22] 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. Philadelphia, Pennsylvania: Association for Computational Linguistics, 2002, 311–318. DOI: 10.3115/1073083.1073135. URL: <https://doi.org/10.3115/1073083.1073135>.

- [23] 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>.
- [24] 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>.
- [25] Ricardo Rei et al. *COMET: A Neural Framework for MT Evaluation*. 2020. arXiv: 2009.09025 [cs.CL]. URL: <https://arxiv.org/abs/2009.09025>.
- [26] 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>.
- [27] Tatoeba Community. *About Tatoeba*. Accessed: 2024-07-06. 2024. URL: <https://tatoeba.org/en>.
- [28] 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>.
- [29] *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](https://try.inten.to/mt_report_2020).
- [30] 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>.
- [31] 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.
- [32] 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](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- [33] Warren Weaver. “Memorandum on Translation”. In: *MT News International* 22 (1999), pp. 5–6, 15.
- [34] 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>.
- [35] 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](https://en.wikipedia.org/w/index.php?title=Google_Translate&oldid=1232822378).
- [36] Linting Xue et al. *mT5: A massively multilingual pre-trained text-to-text transformer*. 2021. arXiv: 2010.11934 [cs.CL]. URL: <https://arxiv.org/abs/2010.11934>.
- [37] Qi Ye et al. “When and Why are pre-trained word embeddings useful for Neural Machine Translation”. In: *HLT-NAACL*. 2018.