

From LLM to NMT: Advancing Low-Resource Machine Translation with Claude

Maxim Enis and Mark Hopkins

Williams College
me4@williams.edu

Abstract

We show that Claude 3 Opus, a large language model (LLM) released by Anthropic in March 2024, exhibits stronger machine translation competence than other LLMs. Though we find evidence of data contamination with Claude on FLORES-200, we curate new benchmarks that corroborate the effectiveness of Claude for low-resource machine translation into English. We find that Claude has remarkable *resource efficiency* – the degree to which the quality of the translation model depends on a language pair’s resource level. Finally, we show that advancements in LLM translation can be compressed into traditional neural machine translation (NMT) models. Using Claude to generate synthetic data, we demonstrate that knowledge distillation advances the state-of-the-art in Yoruba-English translation, meeting or surpassing strong baselines like NLLB-54B and Google Translate.

1 Introduction

Large language models (LLMs) (Kaplan et al., 2020; Hoffmann et al., 2022) have emerged as a breakthrough technology for natural language processing. LLMs have demonstrated a remarkable ability to perform many downstream tasks (Brown et al., 2020), including machine translation. In fact, Zhu et al. (2023); Jiao et al. (2023); Robinson et al. (2023) have shown that the translation performance of GPT-4 (Achiam et al., 2023) is competitive with state-of-the-art neural machine translation (NMT) systems on several high-resource language pairs. However on low-resource language pairs, they showed that existing LLMs lag behind specialized systems like Meta AI’s NLLB-54B¹ (Team et al., 2022).

In this paper, we present evidence that this performance gap may be closing. On 25% of evaluated language pairs, we show that Claude 3 Opus,

an LLM produced by Anthropic, surpasses strong baselines, like Google Translate and NLLB-54B, when translating into English. Surprisingly, the source languages range from very low- to high-resource, indicating that Claude may have broader machine translation capabilities than prior LLMs. Our findings are based on newly created parallel corpora which are verifiably unseen by Claude, since an auxiliary finding of our paper is that Claude exhibits evidence of data contamination (Sainz et al., 2023) on existing benchmarks like FLORES-200 (Guzmán et al., 2019; Goyal et al., 2021; Team et al., 2022).

We also corroborate the findings of Zhu et al. (2023), who have shown that current LLMs are most effective at machine translation when English is the target language (i.e. they are better at xxx->eng translation than eng->xxx translation). Although Claude outperforms NLLB-54B on 55.6% of language pairs in the xxx->eng direction, it only outperforms NLLB-54B on 33.3% of language pairs in the eng->xxx direction, suggesting that supervised baselines still have an edge over LLMs when English is **not** the target language of the translation task.

However, the costs and inference time of massive LLMs like Claude limit the scope of their applicability for machine translation. In section 5, we show that distillation techniques (Hinton et al., 2015; Kim and Rush, 2016) can be applied productively to LLMs to create compact NMT models that outperform the state-of-the-art. We believe that further refinements and optimizations of our methods can result in even better performance, and that many more language pairs, whether currently supported by translation systems or not, are amenable to our approach.

In summary, we make the following contributions:

1. We find evidence of data contamination with Claude on the FLORES-200 benchmark.

¹Model card: <https://huggingface.co/facebook/nllb-moe-54b>.

2. By creating new and unseen evaluation benchmarks for 36 language pairs, we show that Claude nonetheless demonstrates state-of-the-art machine translation ability for many language pairs, including low- and very low-resource language pairs. We provide evidence that Claude’s translation performance (when English is the target language) has higher *resource efficiency* than other LLMs.
3. We show that when translating from English into a low-resource language, a large gap still exists between LLMs and state-of-the-art neural machine translation (NMT) systems on most languages. Even so, we show that Claude outperforms strong baselines for **two** such language pairs.
4. We demonstrate that translation abilities of Claude can be leveraged to advance the state-of-the-art in traditional neural machine translation (NMT) by generating a parallel corpus from Claude translations and fine-tuning the inexpensive model on this corpus. We describe an approach that leverages Claude’s context window to reduce distillation costs and improve translation quality, by ‘batching’ sentences from the same web-crawled document into the same prompt.

2 Background

LLM translation. Prior work has examined the translation abilities of large language models. [Robinson et al. \(2023\)](#) and [Zhu et al. \(2023\)](#) both run empirical studies assessing the translation ability of LLMs like GPT-4 on the FLORES-200 benchmark. Both works find that some LLMs are competitive with NLLB-54B on high-resource languages, but lag behind on low-resource languages. [Robinson et al. \(2023\)](#) find that the *number of Wikipedia pages in a given language* is the most important feature to predict the performance of GPT-4 translation. [Stap and Araabi \(2023\)](#) evaluate GPT-4 translation of low-resource indigenous American languages into Spanish. For all languages considered, the LLM underperforms a fine-tuned, supervised multilingual NMT model. We will reexamine these results in light of the release of Claude 3 Opus.

Dataset contamination. [Zhu et al. \(2023\)](#) examine dataset contamination for LLMs, finding that

FLORES-200 is unsuitable for evaluation on the BLOOMZ ([Muennighoff et al., 2023](#)) LLM. [Sainz et al. \(2023\)](#) highlight the dangers of evaluating closed-source LLMs on public benchmarks, arguing that data leakage should be a central concern for modern natural language processing researchers. [Carlini et al. \(2023\)](#) also consider the problem of quantifying data contamination on closed-source LLMs, using an *information extraction* approach to detect contamination.

Knowledge distillation with LLMs. To the best of our knowledge, [Li et al. \(2024\)](#) is the only work to examine knowledge distillation between LLMs and NMT systems. Their approach involves identifying corrections to translations generated by a student model, making those corrections, and then synthetically generating similar parallel sentences. Although their approach shows promise in the high-resource Chinese-English and English-German language pairs, we note that synthetic corpus generation from monolingual data is inefficient on high-resource pairs due to the large quantity of data needed to achieve an increase in model performance. In our paper, we examine knowledge distillation into low-resource language pairs, which is a substantially less data-hungry setting.

Relatedly, [Cho and Hariharan \(2019\)](#) studied knowledge distillation from very large models to much smaller models. By running knowledge distillation experiments on CNNs evaluated on the ImageNet dataset ([Deng et al., 2009](#)), the authors conclude that “bigger models are not better teachers”, and that higher teacher accuracy does not necessarily correspond to better child model performance. Given these negative results, we aim to find whether distillation is viable from LLMs to much smaller NMT models.

LLM-based document translation. We introduce a novel approach to generate training data by using sentence-aligned document translation. [Karpinska and Iyyer \(2023\)](#) have evaluated various approaches to document translation with LLMs, such as paragraph-level translation and sentence-level translation with the paragraph context. Neither of these approaches support our need for single-prompt sentence-by-sentence document translation. [Wang et al. \(2023\)](#) prompt ChatGPT to translate a document sentence-by-sentence by including sequential boundary tags, but they find that ChatGPT tends to translate “without adhering to

strict sequential boundaries”, making it difficult to extract parallel sentence pairs.

3 Experiments

Our main experiments involve testing a variety of different languages, from high- to low- to very low-resource, against a number of different datasets.

3.1 Languages

Following (Koisshckenov et al., 2023), we classify languages as very low-resource if they have less than 100k bitexts, low-resource if they have between 100k and 1m bitexts, and high-resource if they have more than 1m bitexts, according to (Team et al., 2022)². We experiment on English and a selection of 36 other languages, of which 15 are high-resource, 17 are low-resource, and 4 are very low-resource. All languages are supported by Google Translate, NLLB-200 (Guzmán et al., 2019), and are included in the FLORES-200 dataset. Every language is evaluated in both the eng->xxx and xxx->eng directions. We do not conduct experiments on non-English-centric language pairs. The full list of languages is provided in Table 2.

3.2 Datasets

We benchmark our model against the following datasets.

3.2.1 FLORES-200

FLORES-200 is a high-quality evaluation dataset containing human-curated translations between English and 204 different languages (Guzmán et al., 2019; Goyal et al., 2021; Team et al., 2022). The dataset serves as a universal benchmark for which to evaluate state-of-the-art in a number of different languages. Due to budget constraints³, we preselect 100 random sentences from the FLORES devtest split and test each language pair on this subset. Empirically, we find that NLLB-54B performance on this subset is consistent with the published metrics⁴.

3.2.2 BBC News

The FLORES datasets are high-quality, but might have both source and target sentences seen by the LLM. Thus data contamination is possible. We note that even the private FLORES-200 test set⁵

may be subject to this bias for LLM evaluation, since the English data comes from Wikipedia, on which the LLM has almost certainly been trained. These challenges highlight the importance of developing a machine translation benchmark for LLMs with unseen source and target sentences.

In order to check that this bias does not influence the results, we automatically create totally unseen datasets by aligning articles on BBC News. The concept of bitext mining is not novel (Heffernan et al., 2022); however, typically the quality of mined data is insufficient for evaluation. Using heuristics specific to BBC, we improve the confidence in aligned translations. For example, BBC articles that are translations of each other tend to include the same images, so we restrict parallel mining to documents with similar images according to Google Reverse Image Search⁶. A more detailed explanation of the mining approach is provided in Appendix A.3.

Through our parallel mining approach, we find parallel sentences from news articles between English and 36 other languages. Then, we filter all articles with publication date prior to the training cutoff of Claude, and finally evaluate on a subset of 100 sentences from each language.

3.2.3 Maltese Speech

In order to verify the LLM performance on different domains, we consider another totally unseen dataset: Maltese-English speech pairs taken from transcriptions and translations of the Common Voice, created for the shared IWSLT 2024 task (Hernandez Mena et al., 2020). We test on a random selection of 100 sentences from this dataset.

3.2.4 Evaluation Metrics

We evaluate on two automatic evaluation metrics: SentencePiece BLEU (Papineni et al., 2002), using the FLORES-200 tokenizer, and chrF++ (Popović, 2017)⁷. We choose both metrics for consistency with the NLLB-200 evaluation methodology (Team et al., 2022). We avoid using model-based automatic evaluation metrics such as COMET (Rei et al., 2020) since proper evaluation of very low-resource languages is critical to our methodology.

3.3 Baselines

We benchmark Claude performance against the strongest multilingual baselines, NLLB-200 and

²See <https://tinyurl.com/535f7ust>.

³The cost comes mainly from requests to the Claude API.

⁴See <https://tinyurl.com/nllb200moe54bmetrics>.

⁵Only the |dev| and |devtest| splits have been publicly released.

⁶images.google.com

⁷We conduct evaluation using the [sacreBLEU](#) library.

Google Translate. Prior work has examined the performance of various LLMs on the FLORES dataset (Zhu et al., 2023; Robinson et al., 2023). Previously, the NLLB-200 multilingual translation model (Team et al., 2022) has outperformed existing LLMs, especially on low-resource languages. Meanwhile, commercial translation systems (e.g. Google Translate) have outperformed NLLB, especially on medium to high resource languages (Zhu et al., 2023). We benchmark against the best NLLB model (NLLB-54B), and against the public-facing Google Translate API.

3.4 Claude Translation Methodology

To generate translations from Claude, we use the model Claude 3 Opus. All experiments were run between March 2024 and April 2024.

3.4.1 Prompt Tuning

Prior work has optimized LLM prompts for machine translation. Following the findings of (Zhang et al., 2023), we use a tuneable number of in-context sentence exemplars within the prompt, drawn from the dev split of the FLORES-200 dataset. We also introduce a new prompt used to “batch” translations from the same document together, leading to translation quality improvement and reduced API cost. The exact prompts are provided Appendix A.2, with examples shown in Table 3 and Table 4. We use 8 in-context exemplars for sentence-level prompts and 1 in-context exemplar for document-level prompts.

3.4.2 Temperature Tuning

The Claude API allows temperature as input, ranging from 0 to 1. Based on tuning experiments, we set temperature=0.7 for each evaluation.

4 Results

We begin by outlining the empirical performance of Claude on the FLORES datasets, and subsequently compare the results to unseen datasets such as the BBC News datasets and the Maltese speech dataset.

When translating into English, Claude surpasses the baselines on the majority of the considered language pairs in the FLORES-200 dataset. We showcase the full results on all 36 language pairs in Table 5 and Table 6.

The chrF++ score of Claude exceeds the baselines for 58% of language pairs in the xxx->eng

direction and for 11% of the evaluated language pairs in the eng->xxx direction. However, these results should not be taken as definitive evidence of Claude’s translation ability, as we will show that Claude demonstrates signs of data contamination on the FLORES-200 benchmark.

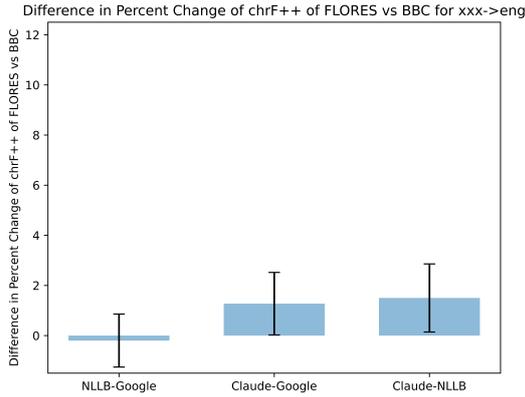
On unseen BBC datasets, Claude beats the SOTA on a number of both low and high-resource languages when translating into English. We evaluated all 36 language pairs in the eng->xxx and xxx->eng translation directions. We provide the complete results in Table 7 and Table 8.

In summary, Claude surpasses the state-of-the-art on Bengali-English, French-English, Kyrgyz-English, Korean-English, Nepali-English, Russian-English, Ukrainian-English, and Yoruba-English, English-Korean, and English-Thai, constituting 11 translation directions out of the 72 total. These account for 25% of language pairs in the xxx->eng direction and 5.5% of language pairs in the eng->xxx direction. The languages cover a wide array of scripts and language families, indicating that Claude-based translation is not necessarily biased towards languages similar to English. In Section 5, we will show that all of these languages may be amenable to knowledge distillation techniques in order to advance the state-of-the-art.

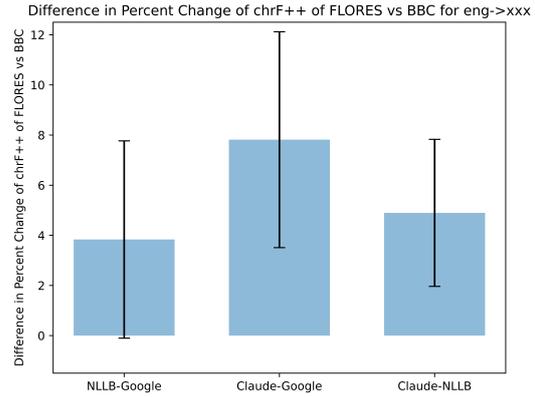
Claude shows signs of data contamination on the FLORES-200 dataset in both translation directions. It is likely that Claude has seen the FLORES data during its training, but it remains unclear whether this measurably affects Claude’s performance on the benchmark. We investigate this question by comparing the results on FLORES versus BBC News. Because the FLORES and BBC datasets may vary in difficulty and quality, we cannot directly compare the raw chrF++ scores of each model across the datasets. However, we expect that some model has no dataset contamination relative to another model if the relative performance between the models is similar between the dataset in question and unseen data.

In Figure 1, we visualize this difference. In the xxx->eng direction, we observe that Google and NLLB have very similar performance across the FLORES and BBC datasets, indicating little-to-no contamination of either dataset for either model.

However, we observe substantial increase in performance of Claude on FLORES compared to BBC



(a) Comparison for the xxx->eng direction.



(b) Comparison for the eng->xxx direction.

Figure 1: Comparison of relative performance of between FLORES and BBC datasets.

relative to either Google or NLLB, which suggests that Claude has overfit the FLORES dataset, with its performance overrepresented by 1-2 percentage points. This analysis calls into question the validity of evaluating Claude on FLORES.

In the eng->xxx direction, we observe somewhat more complicated behavior. As before, Claude also performs relatively worse on BBC than the other models, and to a significantly larger extent. Thus, LLMs may be more prone to contamination in the eng->xxx direction. However, the NLLB-Google column also suggests dataset contamination, such that either NLLB has overfit FLORES or Google has overfit BBC. We reject the first possibility since it is known that NLLB has not been trained on FLORES (Team et al., 2022); therefore, Google may be biased toward the BBC benchmark in the eng->xxx direction.

We offer a possible explanation of these findings: when BBC authors write translations of existing English articles into other languages, they might use Google Translate to provide candidate translations of the article before they edit the writing to improve fluency. This procedure introduces a bias towards Google-like translations in the eng->xxx direction. Thus, we avoid evaluating on Google in the eng->xxx BBC direction for the remainder of the analysis.

Claude’s translation ability is better when translating into English rather than out of English.

Zhu et al. (2023) observe that LLMs are better at translating into English, and we observe the same effect. Claude improves over the SOTA in a much smaller percentage of language pairs in the

eng->xxx direction.

According to Table 8, Claude exceeds NLLB-54B on 56% of language pairs in xxx->eng translation but only 33% of pairs in eng->xxx translation. Further, the mean improvement of Claude over NLLB is 0.81% in xxx->eng translation but -3.05% in eng->xxx translation. Thus, Claude has impressive translation into English but still struggles with out-of-English translation.

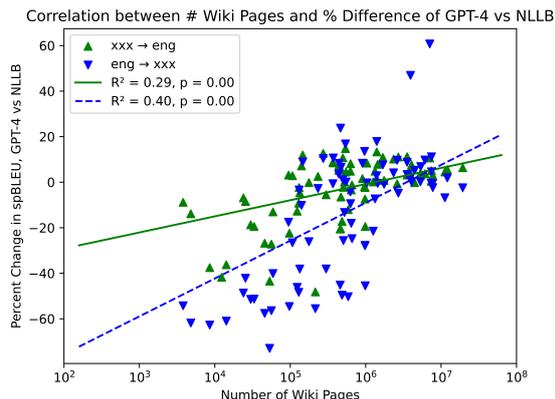
Unlike other evaluated LLMs, Claude has remarkable resource efficiency when translating into English.

Previous works (Zhu et al., 2023; Stap and Araabi, 2023; Robinson et al., 2023) have found that LLMs are more competitive with multilingual NMT models (like NLLB) when translating high-resource (rather than low-resource) language pairs. We dig deeper into this claim and evaluate the extent to which it is true for Claude-based translation.

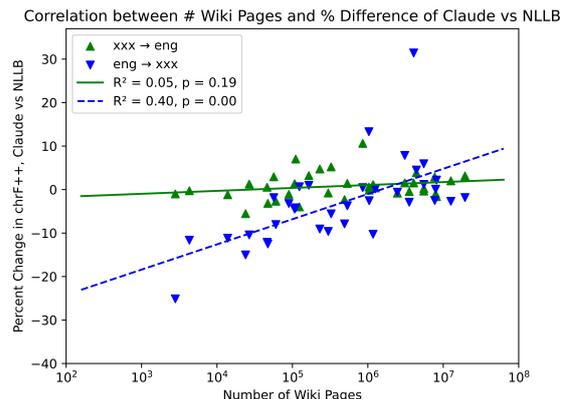
By *resource efficiency*, we mean the extent to which the performance of a multilingual translation engine depends on the resource level (e.g. high, low, very low) of the language pair. Previous work (Zhu et al., 2023) has shown that LLMs have inferior ability to translate low-resource languages relative to NMT models like NLLB-54B, indicating poor relative resource efficiency⁸.

To quantify the resource efficiency of a multilingual machine translation system *A* relative to a baseline system *B*, we perform a linear regres-

⁸Note that Zhu et al. (2023) claims that LLMs have good resource efficiency after finding that XGLM 7.5B can translate some *unresourced* languages. However, the same work finds that LLMs lag behind specialized NMT models when translating low-resource languages.



(a) Correlation between number of Wikipedia pages and GPT-4 relative performance compared to NLLB-1.3B.



(b) Correlation between number of Wikipedia pages and Claude relative performance compared to NLLB-54B.

Figure 2: Comparison of performance correlations between Claude and GPT-4, relative to NLLB models. Figure 2b uses our data and Figure 2a is generated from data from Zhu et al. (2023).

sion on the performance of A relative to B on a language pair against an independent variable measuring the resource level of the language pair. Following (Robinson et al., 2023), we use the number of Wikipedia articles⁹ as our independent variable, and measure the percent difference between our LLM of interest and the NLLB baseline model. We then run a t-test to assess the significance that the slope coefficient is nonzero¹⁰, using a significance cutoff of 0.05. A positive slope with significant p-value indicates that the resource-level positively predicts the translation quality of the LLM relative to NLLB, so the LLM has low resource efficiency. A slope close to zero with non-significant p-value indicates that the resource efficiency is close to NLLB. Finally, a significantly negative slope indicates that the LLM is *more* resource efficient than the supervised baseline.

Using this setup, we verify that 8 LLMs evaluated in Zhu et al. (2023), including LLAMA and GPT-4, all exhibit significant correlation on resource level with respect to comparative performance to NLLB, which indicates that NLLB is more resource efficient than these other LLMs. In Figure 2a, we plot the correlation for GPT-4 (separated by translation direction), where the data is collected from Zhu et al. (2023). In both translation directions, the GPT-4 model has significant correlation with respect to resource level on performance relative to NLLB. We provide similar plots for 7 other LLMs of interest in Appendix 4.

⁹https://en.wikipedia.org/wiki/List_of_Wikipedias

¹⁰We use the Wald Test with t-distribution of the test statistic (Wald, 1943), which is default in scipy.

However, there is one outlier: Claude. In Figure 2b, we show that in xxx->eng translation, Claude has comparable resource efficiency to NLLB. Claude may be the first LLM to demonstrate resource efficiency in machine translation versus strong NMT baselines. Thus, among current LLMs, Claude shows particular promise as a low-resource translator.

Claude outperforms NLLB on the IWSLT 2024 Maltese-English Shared Task dataset.

Our next dataset comes from the IWSLT 2024 Shared Task in low-resource machine translation. We consider the performance of Claude on the development split of the unseen, parallel MASRI-HEADSET dataset in the Maltese speech domain (Hernandez Mena et al., 2020)¹¹.

In Figure 3, we display the performance of Google, NLLB, and Claude on the MASRI-HEADSET dataset. In the Maltese-English direction, we observe unusual translation behavior: Google demonstrates nearly perfect BLEU and chrF++. Even a perfect translator usually should not have perfect scores against reference translations due to the intrinsic variations of natural language. The results are drastically different compared to the English-Maltese translation direction. One possible explanation is that the dataset translations were automatically translated from Maltese into English by Google Translate and then post-

¹¹Although the dataset was created before the training cutoff of Claude, it requires access to be granted in order to be downloaded. Furthermore, translations were not created until 2024.

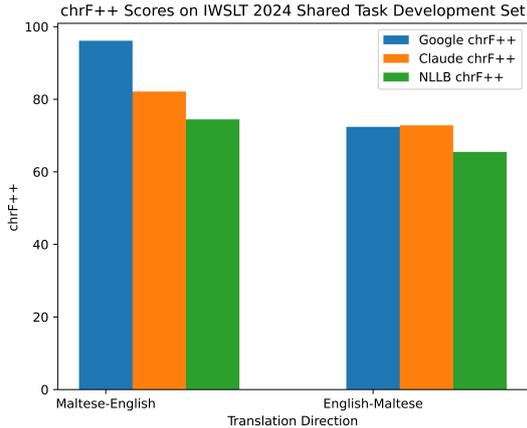


Figure 3: Translation performance, in chrF++, of the models on the development set of the IWSLT 2024 Maltese-English shared task.

edited.

Nevertheless, Claude outperforms NLLB in both the Maltese-English and English-Maltese directions. Thus, Claude demonstrates robust translation ability across multiple domains.

5 Knowledge Distillation

Although LLMs may achieve state-of-the-art results in certain translation directions, the cost, time, and energy use of computational inference limits their applicability as translators. For example, a system such as Google Translate needs cheap inference to support the billions of words translated on a daily basis (Turovsky, 2016). The task of compressing model performance into smaller models is known as knowledge distillation, and has been studied both in the broader deep learning literature (Hinton et al., 2015), and for machine translation (Kim and Rush, 2016). In this section, we devise LLM-based knowledge distillation methods.

5.1 Yoruba

Spoken primarily in the West African countries Nigeria, Benin, and Togo, Yoruba has over 44 million native speakers (Eberhard et al., 2024). However, Yoruba is low-resource and existing NMT models are low-quality. According to our results in Table 8, Claude may be able to translate Yoruba-English better than Google and NLLB-54B, showing promise as a potential case study for knowledge distillation.

5.2 Sequence-KD

Sequence-KD (Kim and Rush, 2016) is a knowledge distillation method which involves translating the source side of the teacher’s training corpus via beam search, and then training a student model on the translated corpus. Here, we generate translations of a monolingual Yoruba corpus using the LLM. Following the methodology of Dabre and Fujita (2020), we fine-tune a strong base model on the small forward-translated corpus.

5.3 Document Translation

We translate a new monolingual corpus of Yoruba news articles crawled from bbc.com/yoruba, ignoring all articles written after June 1, 2023 to ensure no overlap with the test corpus. To reduce number of inferences requested from Claude, and reuse in-context-exemplars across multiple sentences, we translate each article with one API request. Previous work has shown that LLMs have impressive document-level translation ability (Karpinska and Iyyer, 2023; Wang et al., 2023), although none have studied sentence-level document translation, where the aligned sentence translations must be recovered from the original document and the translated document. To achieve sentence-aligned document translation, we use the prompts described in Appendix A.2. An example prompt is shown in Table 4.

From the monolingual crawled corpus, we forward-translate 53193 sentences with NLLB-54B to create the base distillation model. Next, to create the Claude and Google distillation datasets, we forward-translate a selection of 431 Yoruba news documents totaling 7996 Yoruba sentences. The total API cost of parallel corpus generation using Claude 3 Opus was \$46.06.

5.4 Experiments

Our approach is to train a strong base model by distilling NLLB-54B Yoruba-English into a smaller model. We use the NLLB-54B distillation dataset, as well as the train split of the MENYO-20k dataset (Adelani et al., 2021), which contains human-labeled Yoruba-English parallel sentences across multiple domains. Then, we fine-tune NLLB-1.3B¹² on a shuffling of these two datasets to create our base distillation model `base_distill`. To train the distillation models `claude_distill` and `google_distill`, we fine-tune `base_distill` on

¹²Model card: [facebook/nllb-200-distilled-1.3B](https://facebook.github.io/nllb-200-distilled-1.3B/)

| Metric | Baselines | | | | Distillation Models | | |
|--------|-----------|----------|-------------|--------------|---------------------|----------------|----------------|
| | google | nllb-54B | claude_sent | claude_doc | base_distill | google_distill | claude_distill |
| spBLEU | 21.09 | 22.51 | 21.15 | 26.17 | 22.82 | 20.81 | 22.61 |
| chrF++ | 41.99 | 43.26 | 43.78 | 47.06 | 43.23 | 42.29 | 44.15 |

Table 1: spBLEU, chrF++ scores on the Yoruba-English BBC News dataset. Bolded results are best in each category.

the respective synthetic forward-translated corpus, combined with a random sampling from the NLLB-54B distillation corpus and MENYO-20k train corpus in equal ratio 1:1:1.

The model `base_nllb_distill` is trained until validation loss is minimized, which took 1 epoch. The models `claude_distill` and `google_distill` are trained until spBLEU on the validation set is maximized, which took 2 epochs. We used an AdamW optimizer (Loshchilov and Hutter, 2019) with an initial learning rate of $5e-5$ for `base_distill` and $3e-5$ for `claude_distill` and `google_distill`. Each model is trained with batch size 6. To compute translations from each model, we use beam search with 5 beams and restrict repetition of n-grams of 3 or larger.

5.5 Distillation Results

The results of the experiments are displayed in Table 1. Under the ‘‘Baselines’’ category, the columns `claude_sent` and `claude_doc` refer to translating with Claude using the sentence-level prompt or the document-level prompt (see Appendix A.2).

The best model by far is `claude_doc`, surpassing all other models by over 3 spBLEU and chrF++. These results suggest that document-level context can substantially improve translation quality of Claude.

We observe that our model, `claude_distill`, attains higher chrF++ than all other non-Claude baselines, including Google Translate and NLLB-54B. The model also has considerably better performance than `google_distill`, which demonstrates the importance of data quality when augmenting a training corpus with synthetic data. Thus, by distilling on a relatively small dataset with Claude, we are able to match or exceed the performance of `claude_sent`, and construct a small model that outperforms the baselines.

6 Limitations

Due to a limited budget for the Claude API, our per-language dataset size was constrained. Moreover, all evaluated language pairs involved English as the source or target. Finally, since all unseen

data comes from BBC News, our evaluation strategy would not directly apply to languages that BBC News does not support. Since evaluating the translation performance of LLMs on published MT benchmarks can be problematic due to data contamination, an important question remains on how to evaluate Claude (and other closed-source LLMs) on a broader set of languages.

7 Conclusions

Our results point toward a future era of LLM-powered machine translation. Although we find that Claude shows signs of data contamination on FLORES-200, we also evaluate Claude on unseen datasets and find that Claude 3 Opus outperforms NLLB-54B on 44% of language pairs and Google Translate on 22%. Unlike prior LLM models, the spBLEU and chrF++ scores of Claude remain competitive, or even exceed, the baseline models on high, low, and very-low resource language pairs. In fact, among 8 other LLMs, we show that Claude uniquely demonstrates a *resource efficiency* comparable to NLLB-54B. Finally, in section 5, we show that state-of-the-art results from LLMs can be distilled into inexpensive machine translation models and we create a simple system that beats baselines on Yoruba-English for BBC News articles.

Due to the increasing capabilities of LLMs as models scale in size and efficiency, we expect that (potentially closed-source) LLMs will surpass the state-of-the-art in more and more language pairs. Our work demonstrates that these advancements can be harnessed by the MT community to improve under-resourced language pairs.

This work opens many interesting avenues of future research. Our evaluations are limited to English-centric translation, but our methods (and automatic dataset construction) should apply to any language pair. Zhu et al. (2023) has shown that LLMs struggle with non-English-centric machine translation - can we use the mined BBC dataset methods to evaluate new LLMs and prompting techniques on non-English-centric translation?

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- David Adelani, Dana Ruitter, Jesujoba Alabi, Damilola Adebajo, Adesina Ayeni, Mofe Adeyemi, Ayodele Esther Awokoya, and Cristina España-Bonet. 2021. [The effect of domain and diacritics in Yoruba-English neural machine translation](#). In *Proceedings of Machine Translation Summit XVIII: Research Track*, pages 61–75, Virtual. Association for Machine Translation in the Americas.
- 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*, 33:1877–1901.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2023. [Quantifying memorization across neural language models](#).
- Jang Hyun Cho and Bharath Hariharan. 2019. [On the efficacy of knowledge distillation](#). *CoRR*, abs/1910.01348.
- Raj Dabre and Atsushi Fujita. 2020. [Combining sequence distillation and transfer learning for efficient low-resource neural machine translation models](#). In *Proceedings of the Fifth Conference on Machine Translation*, pages 492–502, Online. Association for Computational Linguistics.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. [Imagenet: A large-scale hierarchical image database](#). In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255.
- David M. Eberhard, Gary F. Simons, and Charles D. Fennig. 2024. [Yoruba](#). In David M. Eberhard, Gary F. Simons, and Charles D. Fennig, editors, *Ethnologue: Languages of the World*, 27 edition. SIL International, Dallas, Texas.
- Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2021. The flores-101 evaluation benchmark for low-resource and multilingual machine translation.
- Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc’Aurelio Ranzato. 2019. Two new evaluation datasets for low-resource machine translation: Nepali-english and sinhala-english.
- Kevin Heffernan, Onur Çelebi, and Holger Schwenk. 2022. [Bitext mining using distilled sentence representations for low-resource languages](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2101–2112, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Carlos Daniel Hernandez Mena, Albert Gatt, Andrea DeMarco, Claudia Borg, Lonke van der Plas, Amanda Muscat, and Ian Padovani. 2020. [MASRI-HEADSET: A Maltese corpus for speech recognition](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6381–6388, Marseille, France. European Language Resources Association.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. [Distilling the knowledge in a neural network](#).
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. [Training compute-optimal large language models](#).
- Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. [Is chatgpt a good translator? yes with gpt-4 as the engine](#).
- 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](#).
- Marzena Karpinska and Mohit Iyyer. 2023. [Large language models effectively leverage document-level context for literary translation, but critical errors persist](#). In *Proceedings of the Eighth Conference on Machine Translation*, pages 419–451, Singapore. Association for Computational Linguistics.
- Yoon Kim and Alexander M. Rush. 2016. [Sequence-level knowledge distillation](#).
- Yeskendir Koishekenov, Alexandre Berard, and Vasilina Nikoulina. 2023. [Memory-efficient NLLB-200: Language-specific expert pruning of a massively multilingual machine translation model](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3567–3585, Toronto, Canada. Association for Computational Linguistics.
- Jiahuan Li, Shanbo Cheng, Shujian Huang, and Jiajun Chen. 2024. [Mt-patcher: Selective and extendable knowledge distillation from large language models for machine translation](#).
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#).

- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. [Crosslingual generalization through multitask finetuning](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Maja Popović. 2017. [chrF++: words helping character n-grams](#). In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. [COMET: A neural framework for MT evaluation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.
- 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*, pages 392–418, Singapore. Association for Computational Linguistics.
- Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. [NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10776–10787, Singapore. Association for Computational Linguistics.
- David Stap and Ali Araabi. 2023. [ChatGPT is not a good indigenous translator](#). In *Proceedings of the Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP)*, pages 163–167, Toronto, Canada. Association for Computational Linguistics.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti
- Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. [No language left behind: Scaling human-centered machine translation](#).
- Barak Turovsky. 2016. Ten years of google translate. <https://blog.google/products/translate/ten-years-of-google-translate/>. Accessed: 2024-04-11.
- Abraham Wald. 1943. Tests of statistical hypotheses concerning several parameters when the number of observations is large. *Transactions of the American Mathematical society*, 54(3):426–482.
- Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. [Document-level machine translation with large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore. Association for Computational Linguistics.
- Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. [Prompting large language model for machine translation: A case study](#).
- Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. [Multilingual machine translation with large language models: Empirical results and analysis](#).

A Appendix

A.1 Language Table

| Language | Code | Script | ISO 639-1 | Family | Subgrouping | Res. |
|------------------------|----------|------------|-----------|----------------|-------------------|----------|
| Amharic | amh_Ethi | Ge'ez | am | Afro-Asiatic | Semitic | Low |
| Modern Standard Arabic | arb_Arab | Arabic | ar | Afro-Asiatic | Semitic | High |
| Azerbaijani | azj_Latn | Latin | az | Turkic | Common Turkic | Very Low |
| Bengali | ben_Beng | Bengali | bn | Indo-European | Indo-Aryan | High |
| English | eng_Latn | Latin | en | Indo-European | Germanic | High |
| French | fra_Latn | Latin | fr | Indo-European | Italic | High |
| Western Central Oromo | gaz_Latn | Latin | om | Afro-Asiatic | Cushitic | Low |
| Gujarati | guj_Gujr | Gujarati | gu | Indo-European | Indo-Aryan | Low |
| Hausa | hau_Latn | Latin | ha | Afro-Asiatic | Chadic | Low |
| Hindi | hin_Deva | Devanagari | hi | Indo-European | Indo-Aryan | High |
| Igbo | ibo_Latn | Latin | ig | Atlantic-Congo | Benue-Congo | Low |
| Indonesian | ind_Latn | Latin | id | Austronesian | Malayo-Polynesian | High |
| Japanese | jpn_Jpan | Japanese | ja | Japonic | Japanese | High |
| Kinyarwanda | kin_Latn | Latin | rw | Atlantic-Congo | Benue-Congo | Low |
| Kyrgyz | kir_Cyrl | Cyrillic | ky | Turkic | Common Turkic | Low |
| Korean | kor_Hang | Hangul | ko | Koreanic | Korean | High |
| Marathi | mar_Deva | Devanagari | mr | Indo-European | Indo-Aryan | Low |
| Burmese | mya_Mymr | Myanmar | my | Sino-Tibetan | Burmese-Lolo | Low |
| Nepali | npi_Deva | Devanagari | ne | Indo-European | Indo-Aryan | Very Low |
| Southern Pashto | pbt_Arab | Arabic | ps | Indo-European | Iranian | Very Low |
| Persian | pes_Arab | Arabic | fa | Indo-European | Iranian | High |
| Portuguese | por_Latn | Latin | pt | Indo-European | Italic | High |
| Russian | rus_Cyrl | Cyrillic | ru | Indo-European | Balto-Slavic | High |
| Sinhala | sin_Sinh | Sinhala | si | Indo-European | Indo-Aryan | Low |
| Somali | som_Latn | Latin | so | Afro-Asiatic | Cushitic | Very Low |
| Spanish | spa_Latn | Latin | es | Indo-European | Italic | High |
| Swahili | swh_Latn | Latin | sw | Atlantic-Congo | Benue-Congo | Low |
| Tamil | tam_Taml | Tamil | ta | Dravidian | South Dravidian | Low |
| Telegu | tel_Telu | Telegu | te | Dravidian | South Dravidian | Low |
| Thai | tha_Thai | Thai | th | Kra-Dai | Tai | Low |
| Tigrinya | tir_Ethi | Ge'ez | ti | Afro-Asiatic | Semitic | Low |
| Turkish | tur_Latn | Latin | tr | Turkic | Common Turkic | High |
| Ukrainian | ukr_Cyrl | Cyrillic | uk | Indo-European | Balto-Slavic | High |
| Urdu | urd_Arab | Arabic | ur | Indo-European | Indo-Aryan | Low |
| Northern Uzbek | uzn_Latn | Latin | uz | Turkic | Common Turkic | High |
| Vietnamese | vie_Latn | Latin | vi | Austroasiatic | Viet-Muong | High |
| Yoruba | yor_Latn | Latin | yo | Atlantic-Congo | Benue-Congo | Low |

Table 2: The list of 37 languages used for experimentation. All experiments are on English-centric translation, giving a total of 72 translation directions.

A.2 Claude Prompts

In this subsection, we formally specify the sentence-level prompt used in Section 4 and document-level prompt used in Section 5, and provide examples of both.

Sentence prompt We generate our sentence-level prompts in the following format:

```
{source}: {X1} {target}: {Y1}  
... {source}: {Xn} {target}: {Yn}\n  
{source}: {X'} {target}:
```

where source is the source language (e.g English, Spanish, etc.), target is the target language, X_i is the i 'th source sentence, Y_i is the i 'th target sentence (the gold translation of X_i), and X' is the desired sentence to translate. Here, n is a hyperparameter specifying the number of in-context exemplars. This prompt is taken from Zhang et al. (2023). An example prompt on a French-English translation task is provided in Table 3.

Document prompt We define a *single exemplar* of documents X and Y :

```
{source}: \n  
1. {X[1]}\n  
2. {X[2]}\n  
.  
.  
.  
{i}. {X[i]}\n  
\n  
Line-by-line {target} translations:\n  
1. {Y[1]}\n  
2. {Y[2]}\n  
...  
{i}. {Y[i]}
```

where X is a document with i sentences in the language given by source and Y is the *gold translation* document of i sentences in the language given by target (where each sentence $Y[j]$ is a translation of $X[j]$).

Then, we define the *query prompt*:

```
{source}: \n  
1. {X'[1]}\n  
2. {X'[2]}\n  
.  
.  
.  
{k}. {X'[k]}\n  
\n  
Line-by-line {target} translations:\n
```

where X' is the desired document with k sentences for which to generate line-by-line translations.

Then, the entire sentence-aligned document-level prompt is composed of n single exemplars followed by a query prompt, joined by double newlines. See Table 4 for an example of a document-level prompt and the respective Claude 3 output.

A.3 Creation of Mined BBC Datasets

Our dataset creation procedure involves six main steps:

1. **Monolingual BBC page collection** We begin by accessing the Web Archive API¹³ in order to find monolingual BBC webpages in the source language.
2. **Google Reverse Image Search proposals** We use Google Reverse Image Search to generate candidate translation pages for each source article. In some cases, multiple candidate pages or no candidate pages are proposed for each article.
3. **Sentence splitting** We create a multilingual sentence splitter and split the source and target articles into sentences.
4. **Per-document sentence alignment** We use the Facebook LASER library¹⁴ to align sentences between candidate source and target articles. For each source article, we use the intersection mining technique to find candidate sentences alignments across a given target article. If multiple target documents were proposed, we then select the document and sentence translations which maximize the product of LASER score across the given source document.
5. **Date filtration** We filter all source-target candidate translations where the English BBC document has date beyond September 1, 2023¹⁵ to ensure that the BBC articles are unseen by Claude.
6. **LASER score filtration** We filter all sentence alignments with LASER score less than 1.03. Then, we collect the sentence alignments into a dataset coming from source-target documents by the largest LASER score until we accumulate more than 100 sentences.

Note that due to the candidate article proposal step (Step 4), this procedure is very computationally inexpensive, and is bottlenecked only by HTTPS response time. Thus, we can use the technique to generate larger datasets, if required. We can also modify the date filtration step to collect unseen data for any given LLM with a training cutoff date sufficiently far in the past.

Finally, we briefly note that the sentence alignment procedure may create a dataset that bias towards NLLB, since NLLB has been trained on LASER-mined parallel datasets (Team et al., 2022). Empirically, we find in Figure 1 that such a bias is either very close to zero.

¹³<https://archive.org/developers/wayback-cdx-server.html>

¹⁴<https://github.com/facebookresearch/LASER>

¹⁵As of the publication date of the paper, this date is the beyond the training cutoff for Claude.

| | |
|--------|--|
| Prompt | French: Au fil des siècles, les gens ont soigneusement construit des terrasses sur le paysage accidenté et escarpé jusqu’aux falaises qui surplombent la mer. English: Over the centuries, people have carefully built terraces on the rugged, steep landscape right up to the cliffs that overlook the sea. French: 800 miles du réseau de pipelines Trans-Alaska ont été fermés suite à un déversement de milliers de barils de pétrole brut au sud de Fairbanks, en Alaska. English: 800 miles of the Trans-Alaska Pipeline System were closed down following a spill of thousands of barrels of crude oil south of Fairbanks, Alaska. French: On pense qu’il se présentera à la présidence en 2016. English: |
| Output | He is thought to be running for president in 2016. |

Table 3: An example of a sentence-level translation prompt in a French-English translation task on the FLORES-200 dataset, with 2 in-context exemplars originated from the dev split of FLORES-200. Note that in-context exemplars are not separated by newlines, according to (Zhang et al., 2023).

| | |
|--------|---|
| Prompt | <p>French:</p> <ol style="list-style-type: none"> 1. La manifestation a commencé vers 11 h heure locale (UTC+1) sur Whitehall, en face de l’entrée gardée par la police de Downing Street, la résidence officielle du Premier ministre. 2. À Whitehall, peu après 11 h 00, des manifestants ont bloqué la circulation des véhicules en direction du nord. 3. À 11 h 20, la police a demandé aux manifestants de retourner sur le trottoir, déclarant qu’ils devaient trouver un équilibre entre le droit de manifester et l’augmentation de la circulation. <p>Line-by-line English translations:</p> <ol style="list-style-type: none"> 1. The protest started around 11:00 local time (UTC+1) on Whitehall opposite the police-guarded entrance to Downing Street, the Prime Minister’s official residence. 2. Just after 11:00, protesters blocked traffic on the northbound carriage in Whitehall. 3. At 11:20, the police asked the protesters to move back on to the pavement, stating that they needed to balance the right to protest with the traffic building up. <p>French:</p> <ol style="list-style-type: none"> 1. Le Haut Karabakh déserté révèle les séquelles d’une défaite arménienne fulgurante 2. Dans un avant-poste arménien situé dans les montagnes du Haut-Karabakh, une marmite à moitié pleine se trouvait à côté d’une assiette de nourriture à moitié mangée. 3. Il y avait une cigarette à moitié fumée et un morceau de pain. 4. Dans un autre avant-poste plus petit, plus loin le long de l’ancienne ligne de front, un journal de bord arménien était abandonné dans l’herbe. <p>Line-by-line English translations:</p> |
| Output | <ol style="list-style-type: none"> 1. The deserted Nagorno-Karabakh reveals the aftermath of a lightning Armenian defeat 2. In an Armenian outpost located in the mountains of Nagorno-Karabakh, a half-full pot was next to a half-eaten plate of food. 3. There was a half-smoked cigarette and a piece of bread. 4. In another smaller outpost further along the former front line, an Armenian logbook was abandoned in the grass. |

Table 4: An example of a sentence-aligned document-level translation prompt in a French-English translation task, with one in-context exemplar. The in-context exemplar comes from the FLORES-200 dataset and the translation task comes from <https://www.bbc.com/afrique/articles/c0kxzprpnqgo>.

| Language | Code | xxx -> eng spBLEU | | | eng->xxx spBLEU | | |
|------------------------|----------|-------------------|--------------|--------------|-----------------|--------------|--------------|
| | | Google | Claude | NLLB | Google | Claude | NLLB |
| Amharic | amh_Ethi | 43.18 | 42.95 | 41.74 | 32.98 | 26.05 | 29.30 |
| Modern Standard Arabic | arb_Arab | 48.21 | 49.04 | 49.91 | 48.18 | 44.40 | 42.60 |
| Azerbaijani | azj_Latn | 31.38 | 33.24 | 30.71 | 25.18 | 27.49 | 25.24 |
| Bengali | ben_Beng | 41.31 | 42.11 | 40.99 | 37.28 | 33.92 | 33.91 |
| French | fra_Latn | 49.96 | 50.82 | 48.94 | 55.52 | 54.46 | 51.17 |
| Western Central Oromo | gaz_Latn | 32.34 | 27.20 | 29.56 | 14.90 | 10.03 | 12.46 |
| Gujarati | guj_Gujr | 48.05 | 45.44 | 47.06 | 40.25 | 34.23 | 37.42 |
| Hausa | hau_Latn | 41.77 | 38.18 | 39.97 | 29.96 | 20.74 | 29.60 |
| Hindi | hin_Deva | 47.72 | 46.81 | 45.76 | 45.20 | 38.67 | 41.98 |
| Igbo | ibo_Latn | 39.03 | 32.90 | 37.18 | 23.26 | 17.71 | 20.55 |
| Indonesian | ind_Latn | 52.42 | 53.55 | 48.95 | 52.52 | 48.47 | 48.69 |
| Japanese | jpn_Jpan | 35.92 | 35.00 | 35.40 | 33.33 | 34.08 | 18.63 |
| Kinyarwanda | kin_Latn | 41.41 | 41.48 | 40.59 | 34.87 | 21.47 | 27.03 |
| Kyrgyz | kir_Cyrl | 32.99 | 32.07 | 27.79 | 29.57 | 27.76 | 27.97 |
| Korean | kor_Hang | 38.07 | 37.26 | 35.51 | 30.93 | 28.88 | 24.15 |
| Marathi | mar_Deva | 45.34 | 44.87 | 44.75 | 30.45 | 28.48 | 27.85 |
| Burmese | mya_Mymr | 34.27 | 25.57 | 34.67 | 24.66 | 23.62 | 18.26 |
| Nepali | npi_Deva | 48.72 | 50.53 | 48.59 | 38.27 | 34.49 | 29.83 |
| Southern Pashto | pbt_Arab | 41.68 | 43.30 | 40.62 | 23.71 | 21.11 | 25.59 |
| Persian | pes_Arab | 44.46 | 48.35 | 43.93 | 39.39 | 38.53 | 35.42 |
| Portuguese | por_Latn | 59.55 | 59.89 | 58.29 | 59.96 | 59.48 | 54.99 |
| Russian | rus_Cyrl | 43.74 | 43.89 | 43.23 | 47.05 | 44.16 | 42.59 |
| Sinhala | sin_Sinh | 43.49 | 43.88 | 40.43 | 39.65 | 35.35 | 35.37 |
| Somali | som_Latn | 37.94 | 37.97 | 33.95 | 18.70 | 17.31 | 18.65 |
| Spanish | spa_Latn | 38.37 | 41.59 | 41.08 | 35.71 | 34.92 | 33.41 |
| Swahili | swh_Latn | 53.98 | 55.47 | 50.78 | 43.88 | 40.03 | 36.47 |
| Tamil | tam_Taml | 40.74 | 41.62 | 41.88 | 39.84 | 33.88 | 37.53 |
| Telegu | tel_Telu | 47.66 | 45.82 | 47.54 | 46.05 | 36.99 | 41.78 |
| Thai | tha_Thai | 34.16 | 42.13 | 39.80 | 46.11 | 44.52 | 33.33 |
| Tigrinya | tir_Ethi | 27.48 | 27.34 | 28.93 | 17.34 | 13.99 | 18.71 |
| Turkish | tur_Latn | 47.79 | 49.40 | 45.52 | 45.21 | 42.12 | 41.30 |
| Ukrainian | ukr_Cyrl | 47.40 | 47.52 | 45.32 | 42.80 | 42.74 | 37.01 |
| Urdu | urd_Arab | 43.42 | 44.97 | 44.42 | 33.24 | 30.67 | 31.37 |
| Northern Uzbek | uzn_Latn | 45.25 | 46.43 | 39.51 | 36.33 | 34.50 | 31.09 |
| Vietnamese | vie_Latn | 42.17 | 41.12 | 40.46 | 45.29 | 41.65 | 42.48 |
| Yoruba | yor_Latn | 25.49 | 30.21 | 26.62 | 13.09 | 15.37 | 12.80 |

Table 5: spBLEU scores on the FLORES-200 evaluation set of the Google Translate model versus the Claude 3 Opus model versus the NLLB-54B model. **Warning:** There is evidence for data contamination of Claude on this evaluation set in Section 4.

| Language | Code | xxx->eng chrF++ | | | eng->xxx chrF++ | | |
|------------------------|----------|-----------------|--------------|--------------|-----------------|--------------|--------------|
| | | Google | Claude | NLLB | Google | Claude | NLLB |
| Amharic | amh_Ethi | 62.64 | 61.79 | 60.61 | 41.42 | 35.58 | 37.89 |
| Modern Standard Arabic | arb_Arab | 67.38 | 68.34 | 68.34 | 61.23 | 58.20 | 56.88 |
| Azerbaijani | azj_Latn | 55.55 | 57.76 | 54.82 | 44.45 | 46.20 | 43.33 |
| Bengali | ben_Beng | 62.29 | 62.98 | 61.09 | 50.43 | 48.08 | 46.99 |
| French | fra_Latn | 68.90 | 70.28 | 67.26 | 69.47 | 68.97 | 66.25 |
| Western Central Oromo | gaz_Latn | 52.73 | 49.46 | 50.41 | 39.86 | 34.91 | 37.28 |
| Gujarati | guj_Gujr | 67.66 | 64.53 | 66.79 | 55.77 | 51.01 | 53.18 |
| Hausa | hau_Latn | 59.55 | 56.97 | 57.64 | 52.56 | 45.67 | 52.10 |
| Hindi | hin_Deva | 66.50 | 65.00 | 65.51 | 60.24 | 55.35 | 57.92 |
| Igbo | ibo_Latn | 57.36 | 53.62 | 56.04 | 44.80 | 39.78 | 42.06 |
| Indonesian | ind_Latn | 70.12 | 70.39 | 67.66 | 71.35 | 68.79 | 68.77 |
| Japanese | jpn_Jpan | 58.32 | 58.04 | 56.26 | 35.89 | 33.07 | 26.77 |
| Kinyarwanda | kin_Latn | 59.74 | 60.07 | 58.46 | 56.82 | 45.87 | 49.79 |
| Kyrgyz | kir_Cyrl | 55.08 | 55.31 | 50.10 | 47.69 | 46.11 | 45.93 |
| Korean | kor_Hang | 59.46 | 59.20 | 56.65 | 38.70 | 36.29 | 33.89 |
| Marathi | mar_Deva | 65.48 | 64.04 | 64.77 | 49.15 | 47.55 | 45.87 |
| Burmese | mya_Mymr | 56.26 | 48.87 | 55.79 | 39.93 | 40.66 | 31.45 |
| Nepali | npi_Deva | 67.80 | 68.12 | 66.98 | 55.77 | 52.80 | 46.05 |
| Southern Pashto | pbt_Arab | 62.08 | 62.93 | 61.15 | 40.33 | 37.70 | 41.42 |
| Persian | pes_Arab | 64.94 | 66.94 | 63.62 | 53.96 | 55.48 | 49.61 |
| Portuguese | por_Latn | 75.14 | 75.39 | 74.11 | 73.71 | 73.49 | 70.33 |
| Russian | rus_Cyrl | 64.13 | 65.41 | 63.13 | 62.08 | 60.61 | 57.86 |
| Sinhala | sin_Sinh | 63.62 | 63.56 | 60.54 | 49.92 | 47.73 | 42.65 |
| Somali | som_Latn | 57.35 | 57.43 | 53.61 | 43.62 | 42.00 | 43.09 |
| Spanish | spa_Latn | 60.49 | 63.56 | 61.56 | 56.73 | 56.40 | 54.50 |
| Swahili | swl_Latn | 70.34 | 70.77 | 67.80 | 63.87 | 61.14 | 58.24 |
| Tamil | tam_Taml | 61.59 | 61.49 | 61.14 | 56.04 | 52.62 | 54.39 |
| Telegu | tel_Telu | 66.02 | 64.40 | 65.89 | 59.52 | 52.26 | 56.06 |
| Thai | tha_Thai | 58.25 | 62.87 | 60.36 | 50.80 | 50.26 | 43.32 |
| Tigrinya | tir_Ethi | 51.68 | 50.60 | 50.69 | 25.70 | 23.00 | 26.05 |
| Turkish | tur_Latn | 66.62 | 68.06 | 64.40 | 60.91 | 58.87 | 57.76 |
| Ukrainian | ukr_Cyrl | 65.72 | 66.45 | 63.97 | 59.21 | 59.86 | 54.95 |
| Urdu | urd_Arab | 64.02 | 65.79 | 63.96 | 50.58 | 49.11 | 49.56 |
| Northern Uzbek | uzn_Latn | 64.64 | 66.63 | 60.02 | 55.26 | 54.31 | 51.46 |
| Vietnamese | vie_Latn | 61.65 | 61.82 | 60.54 | 60.77 | 59.21 | 59.00 |
| Yoruba | yor_Latn | 45.99 | 49.23 | 46.57 | 38.27 | 30.19 | 38.25 |

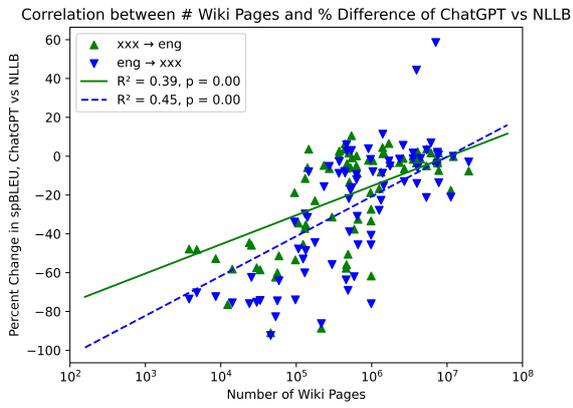
Table 6: chrF++ scores on the FLORES-200 evaluation set of the Google Translate model versus the Claude 3 Opus model versus the NLLB-54B model. **Warning:** There is evidence for data contamination of Claude on this evaluation set as argued in Section 4.

| Language | Code | xxx->eng spBLEU | | | eng->xxx spBLEU | | |
|------------------------|----------|-----------------|--------------|--------------|-----------------|--------------|--------------|
| | | Google | Claude | NLLB | Google | Claude | NLLB |
| Amharic | amh_Ethi | 17.13 | 17.55 | 17.33 | 16.49 | 12.04 | 16.05 |
| Modern Standard Arabic | arb_Arab | 45.49 | 44.59 | 46.49 | 71.04 | 49.48 | 47.12 |
| Azerbaijani | azj_Latn | 35.42 | 34.41 | 34.66 | 48.31 | 33.80 | 38.03 |
| Bengali | ben_Beng | 25.28 | 26.92 | 25.37 | 35.39 | 24.64 | 32.06 |
| French | fra_Latn | 56.51 | 57.30 | 54.91 | 72.90 | 61.41 | 66.71 |
| Western Central Oromo | gaz_Latn | 20.02 | 17.07 | 20.55 | 11.63 | 7.48 | 11.67 |
| Gujarati | guj_Gujr | 31.06 | 28.48 | 30.37 | 32.59 | 22.69 | 25.36 |
| Hausa | hau_Latn | 27.43 | 25.73 | 26.76 | 28.51 | 16.76 | 23.35 |
| Hindi | hin_Deva | 27.64 | 28.59 | 26.89 | 22.65 | 20.88 | 21.74 |
| Igbo | ibo_Latn | 41.95 | 33.82 | 38.20 | 58.13 | 27.16 | 34.90 |
| Indonesian | ind_Latn | 50.63 | 48.37 | 48.76 | 71.74 | 47.76 | 54.90 |
| Japanese | jpn_Jpan | 24.96 | 21.66 | 21.72 | 18.75 | 14.82 | 11.00 |
| Kinyarwanda | kin_Latn | 28.24 | 28.70 | 30.39 | 14.76 | 12.18 | 17.26 |
| Kyrgyz | kir_Cyrl | 21.49 | 21.56 | 19.44 | 23.11 | 18.69 | 21.69 |
| Korean | kor_Hang | 20.49 | 20.62 | 22.67 | 19.82 | 20.29 | 17.70 |
| Marathi | mar_Deva | 18.89 | 16.94 | 18.26 | 24.14 | 15.17 | 19.53 |
| Burmese | mya_Mymr | 27.41 | 25.37 | 23.55 | 26.28 | 14.91 | 19.86 |
| Nepali | npi_Deva | 25.93 | 27.24 | 26.44 | 25.35 | 18.23 | 23.69 |
| Southern Pashto | pbt_Arab | 28.98 | 30.68 | 34.29 | 37.20 | 24.06 | 29.24 |
| Persian | pes_Arab | 38.12 | 36.55 | 37.56 | 43.36 | 35.96 | 32.04 |
| Portuguese | por_Latn | 54.85 | 53.59 | 53.55 | 67.49 | 58.45 | 56.99 |
| Russian | rus_Cyrl | 34.10 | 37.14 | 37.34 | 44.03 | 33.65 | 34.63 |
| Sinhala | sin_Sinh | 38.62 | 35.06 | 36.17 | 39.62 | 28.85 | 33.17 |
| Somali | som_Latn | 49.72 | 46.12 | 44.30 | 74.11 | 32.50 | 40.53 |
| Spanish | spa_Latn | 52.77 | 46.36 | 49.89 | 63.85 | 53.09 | 54.56 |
| Swahili | swh_Latn | 56.87 | 50.71 | 47.63 | 83.18 | 50.85 | 50.60 |
| Tamil | tam_Taml | 29.56 | 25.32 | 27.61 | 47.29 | 28.41 | 35.80 |
| Telegu | tel_Telu | 19.83 | 21.06 | 18.56 | 23.14 | 16.24 | 19.56 |
| Thai | tha_Thai | 19.21 | 19.18 | 20.07 | 26.58 | 27.31 | 23.64 |
| Tigrinya | tir_Ethi | 36.45 | 31.28 | 31.53 | 58.09 | 17.01 | 27.84 |
| Turkish | tur_Latn | 28.93 | 30.64 | 31.41 | 26.94 | 24.74 | 26.10 |
| Ukrainian | ukr_Cyrl | 38.45 | 39.67 | 37.60 | 53.47 | 44.83 | 42.71 |
| Urdu | urd_Arab | 38.18 | 36.64 | 39.78 | 48.87 | 38.14 | 40.15 |
| Northern Uzbek | uzn_Latn | 31.55 | 32.75 | 26.72 | 0.12 | 0.09 | 0.04 |
| Vietnamese | vie_Latn | 40.92 | 39.44 | 38.06 | 54.43 | 40.24 | 42.77 |
| Yoruba | yor_Latn | 20.59 | 21.15 | 21.53 | 13.75 | 15.49 | 18.27 |

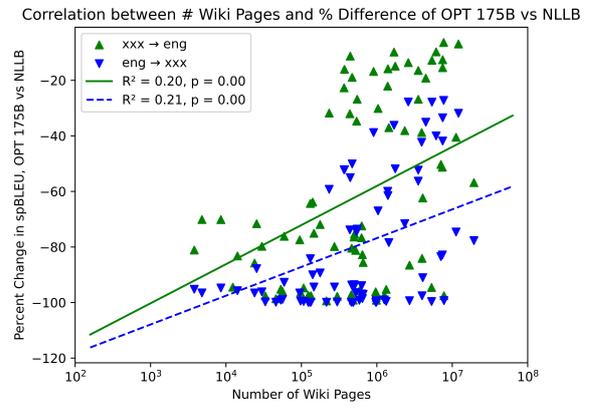
Table 7: spBLEU scores on the BBC evaluation set of the Google Translate model versus the Claude 3 Opus model versus the NLLB-54B model.

| Language | Code | xxx->eng chrF++ | | | eng->xxx chrF++ | | |
|------------------------|----------|-----------------|--------------|--------------|-----------------|--------------|--------------|
| | | Google | Claude | NLLB | Google | Claude | NLLB |
| Amharic | amh_Ethi | 38.84 | 38.08 | 37.88 | 24.96 | 21.32 | 24.26 |
| Modern Standard Arabic | arb_Arab | 66.44 | 65.23 | 66.31 | 77.77 | 60.86 | 59.55 |
| Azerbaijani | azj_Latn | 56.72 | 56.56 | 55.77 | 63.03 | 53.44 | 55.50 |
| Bengali | ben_Beng | 45.55 | 46.66 | 46.14 | 50.14 | 42.77 | 47.66 |
| French | fra_Latn | 71.51 | 71.90 | 70.49 | 80.63 | 74.87 | 76.93 |
| Western Central Oromo | gaz_Latn | 41.35 | 40.25 | 40.36 | 34.83 | 30.12 | 34.08 |
| Gujarati | guj_Gujr | 50.28 | 48.43 | 50.44 | 49.27 | 41.89 | 41.61 |
| Hausa | hau_Latn | 47.62 | 45.88 | 47.38 | 49.30 | 40.07 | 45.78 |
| Hindi | hin_Deva | 48.88 | 48.62 | 48.44 | 44.69 | 43.63 | 43.56 |
| Igbo | ibo_Latn | 58.07 | 52.31 | 55.37 | 69.16 | 45.10 | 53.06 |
| Indonesian | ind_Latn | 66.65 | 64.95 | 65.25 | 85.70 | 72.25 | 74.40 |
| Japanese | jpn_Jpan | 46.42 | 44.76 | 44.09 | 26.41 | 23.21 | 17.66 |
| Kinyarwanda | kin_Latn | 48.11 | 49.30 | 49.89 | 38.14 | 35.79 | 40.29 |
| Kyrgyz | kir_Cyrl | 41.90 | 42.49 | 39.70 | 43.26 | 40.61 | 42.38 |
| Korean | kor_Hang | 47.77 | 47.83 | 47.11 | 26.65 | 27.13 | 25.15 |
| Marathi | mar_Deva | 38.13 | 37.26 | 37.56 | 45.51 | 38.17 | 42.23 |
| Burmese | mya_Mymr | 49.91 | 48.47 | 46.28 | 37.77 | 31.13 | 34.23 |
| Nepali | npi_Deva | 49.13 | 49.27 | 48.61 | 42.80 | 38.54 | 40.35 |
| Southern Pashto | pbt_Arab | 51.37 | 53.16 | 54.63 | 49.67 | 40.20 | 43.71 |
| Persian | pes_Arab | 58.22 | 58.17 | 58.38 | 56.31 | 50.66 | 47.82 |
| Portuguese | por_Latn | 71.25 | 70.84 | 70.66 | 79.05 | 72.73 | 71.86 |
| Russian | rus_Cyrl | 55.69 | 57.45 | 55.83 | 56.11 | 50.04 | 51.32 |
| Sinhala | sin_Sinh | 57.64 | 54.50 | 55.11 | 48.54 | 39.45 | 40.75 |
| Somali | som_Latn | 66.15 | 63.27 | 62.48 | 82.03 | 54.20 | 60.50 |
| Spanish | spa_Latn | 69.02 | 66.22 | 66.89 | 75.94 | 69.66 | 69.60 |
| Swahili | swh_Latn | 70.28 | 66.25 | 64.18 | 89.28 | 69.40 | 68.71 |
| Tamil | tam_Taml | 49.62 | 47.15 | 48.26 | 63.75 | 50.69 | 55.02 |
| Telegu | tel_Telu | 40.09 | 41.35 | 39.30 | 38.41 | 33.79 | 35.78 |
| Thai | tha_Thai | 40.52 | 40.39 | 40.16 | 33.16 | 37.67 | 33.24 |
| Tigrinya | tir_Ethi | 55.28 | 50.52 | 51.03 | 63.25 | 25.98 | 34.69 |
| Turkish | tur_Latn | 50.80 | 51.41 | 51.84 | 49.22 | 47.91 | 48.19 |
| Ukrainian | ukr_Cyrl | 56.21 | 57.70 | 55.55 | 64.55 | 58.56 | 56.04 |
| Urdu | urd_Arab | 59.28 | 59.80 | 59.89 | 61.58 | 53.71 | 55.11 |
| Northern Uzbek | uzn_Latn | 51.36 | 51.14 | 46.23 | 2.27 | 2.21 | 2.20 |
| Vietnamese | vie_Latn | 59.66 | 59.05 | 57.26 | 66.94 | 58.06 | 59.13 |
| Yoruba | yor_Latn | 41.77 | 43.78 | 42.53 | 26.03 | 27.95 | 28.47 |

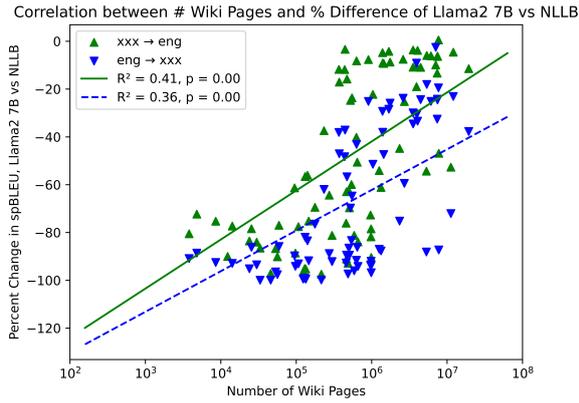
Table 8: chrF++ scores on the BBC evaluation set of the Google Translate model versus the Claude 3 Opus model versus the NLLB-54B model.



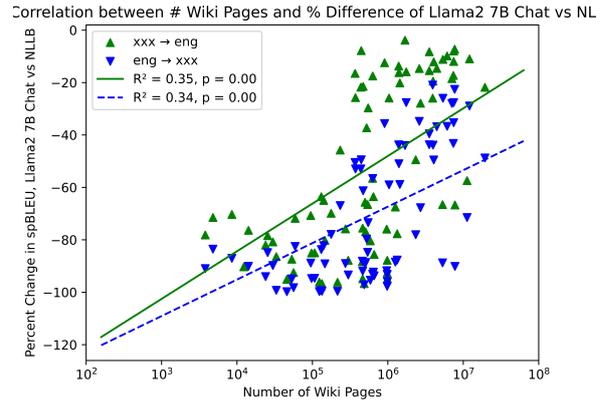
(a) Correlations with ChatGPT translation performance.



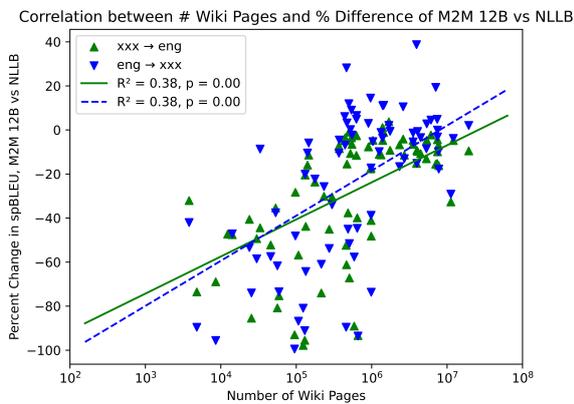
(b) Correlations with OPT 175B translation performance.



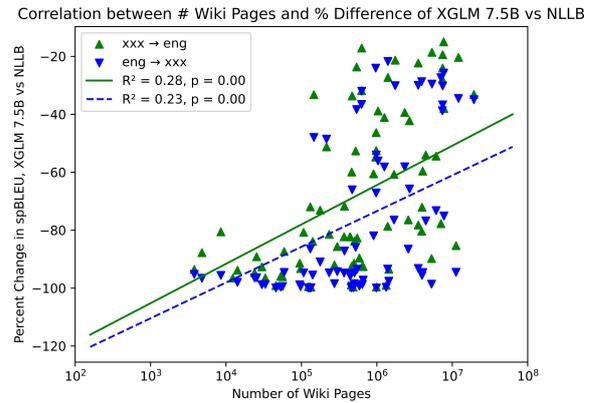
(c) Correlations with LLAMA2 7B translation performance.



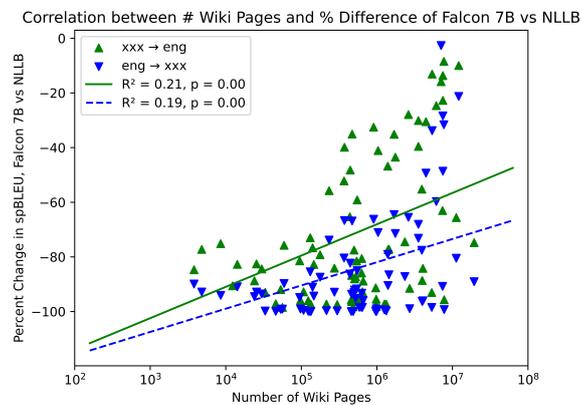
(d) Correlations with LLAMA2 7B Chat translation performance.



(e) Correlations with M2M 12B translation performance.



(f) Correlations with XGLM 7.5B LLM translation performance.



(g) Correlations with Falcon 7B LLM translation performance.

Figure 4: Correlation plots for seven different LLMs, with data adapted from Zhu et al. (2023).