

AFRIDOC-MT: Document-level MT Corpus for African Languages

Anonymous ACL submission

Abstract

This paper introduces AFRIDOC-MT, a document-level multi-parallel translation dataset covering English and five African languages: Amharic, Hausa, Swahili, Yorùbá, and Zulu. The dataset comprises 334 health and 271 information technology news documents, all human-translated from English to these languages. We conduct document-level translation benchmark experiments by evaluating the ability of neural machine translation (NMT) models and large language models (LLMs) to translate between English and these languages, at both the sentence and pseudo-document levels, the outputs being realigned to form complete documents for evaluation. Our results indicate that NLLB-200 achieves the best average performance among the standard NMT models, while GPT-4o outperforms general-purpose LLMs. Fine-tuning selected models leads to substantial performance gains, but models trained on sentences struggle to generalize effectively to longer documents. Furthermore, our analysis reveals that some LLMs exhibit issues such as under-generation, over-generation, repetition of words and phrases, and off-target translations, specifically for translation into African languages.

1 Introduction

The field of machine translation (MT) has seen notable progress in the past years, with neural machine translation (NMT) models achieving close to human performance in many high-resource language directions (Vaswani et al., 2017; Akhbardeh et al., 2021; Mohammadshahi et al., 2022; Yuan et al., 2023; Kocmi et al., 2023; NLLB Team et al., 2024). However, efforts have primarily been concentrated on sentence-level translation, without the use of inter-sentential context.

In recent years, there has been interest in document-level translation (i.e. the holistic trans-

lation of multiple sentences), where sentences are translated with their context rather than in isolation. Document-level translation is important in order to capture discourse relations (Bawden et al., 2018; Voita et al., 2018; Maruf et al., 2021), maintain consistency and coherence across sentences (Herold and Ney, 2023), particularly for technical domains, but poses unique challenges, such as how to handle longer documents (Wang et al., 2024b) given the limited context size of translation models. Current efforts have primarily focused on high-resource language directions, where document-level datasets are readily available (Lopes et al., 2020; Feng et al., 2022; Wu et al., 2023; Wang et al., 2023; Wu et al., 2024), and so far there has been no work on low-resource African languages. Developing and evaluating document-level MT systems for low-resource languages is a useful and under-studied direction, which requires the creation of datasets.

To fill this gap, we present AFRIDOC-MT, a document-level translation dataset for English from and into five African languages: Amharic, Hausa, Swahili, Yorùbá, and Zulu, created through the manual translation of English documents. It consists of 334 *health* documents and 271 *tech* documents. In addition, AFRIDOC-MT supports multi-way translation, allowing translations not only between English and the African languages but also between any two of the languages covered.

We conduct a comprehensive set of document translation benchmark experiments on AFRIDOC-MT, using sentence-level and pseudo-document translation due to most models' limited context length, and then realigning them to form complete documents. We evaluate performance using automatic metrics and compare the results of encoder-decoder models with decoder-only LLMs across both domains. Our results demonstrate that NLLB-200, both before and after fine-tuning on

Dataset	#Langs.	Multiway	Domain	Type	#Sents.	#Docs.
TICO-19 (Anastasopoulos et al., 2020)	12	✓	health	document-level	4k	30
MAFAND-MT (Adelani et al., 2022)	16	✗	news	sentence-level	4k-35k	-
FLORES-200 (NLLB Team et al., 2022)	42	✓	general	sentence-level	3k	-
NTREX-128 (Federmann et al., 2022)	24	✓	news	sentence-level	1.9k	-
AFRIDOC-MT (Ours)	5	✓	tech, health	document-level	10k	271-334

Table 1: Overview of highly related works, including for each dataset the number of African languages, the domain, the kind of MT task they can be used for and the range of the sentence numbers for each language direction.

AFRIDOC-MT, excels in sentence translation, surpassing all other models. GPT-4o performs equally well for sentences and pseudo-documents, while other decoder-only models lag behind. In addition to automatic metrics, we use GPT-4o as a judge, human evaluation, and qualitative assessment to compare documents translation carried out sentence-by-sentence and as pseudo-documents for selected models. The evaluation shows that GPT-4o is generally unreliable for assessing document translations into African languages. However, we observe agreement between other evaluation methods, all indicating that sentence-by-sentence translation results in better document-level translation into African languages. We conduct additional analyses of the models outputs to better understand their behavior and why they under-perform when translating pseudo-documents. They show that LLMs often under-generate, contain repetitions, and produce off-target translations, especially when translating into African languages.

2 Related Work

MT Datasets for African Languages Several MT datasets exist for African languages, including web-mined datasets such as WikiMatrix (Schwenk et al., 2021a) and CCMatrix (Schwenk et al., 2021b). However, they have been adjudged to be of poor quality for certain low-resource subsets, including African languages (Kreutzer et al., 2022). There are also well curated datasets for African languages including the Bible (McCarthy et al., 2020), JW300 (Agić and Vulić, 2019)¹ and MAFAND-MT (Adelani et al., 2022), which are from religious and news domains.

There exist several MT evaluation benchmark datasets for African languages. They can be categorized into two kinds. First, evaluation datasets specifically designed for translating into or from African languages (Ezeani et al., 2020; Azunre et al., 2021; Adelani et al., 2021, 2022, *inter alia*).

¹The dataset is no longer available for use.

Second, benchmark datasets covering many languages, including African languages. For example, TICO-19 (Anastasopoulos et al., 2020), NTREX-128 (Federmann et al., 2022), FLORES-101 (Goyal et al., 2022) and FLORES-200 (NLLB Team et al., 2024) are a few such datasets. However, most of these datasets are designed for sentence-level MT, primarily drawn from religious or news domains, although some consist of translated sentences originating from the same document. To the best of our knowledge, only TICO-19, a health domain translation benchmark, has the potential to be used for document-level MT, while it is restricted to topics related to COVID-19. Table 1 gives a comparison of the most related existing benchmarks.

Document-level Neural Machine Translation

Document-level NMT aims to overcome the limitations of sentence-level systems by translating an entire document as a whole. Similar to context-aware NMT, which involves translating segments with additional, localized context, it differs in that it involves in principle translating an entire document holistically. Both document-level and context-aware MT allow for the possibility of improving translation quality for context-dependent phenomena such as coreference resolution (Müller et al., 2018; Bawden et al., 2018; Voita et al., 2018; Herold and Ney, 2023), lexical disambiguation (Rios Gonzales et al., 2017; Martínez Garcia et al., 2019), and lexical cohesion (Wong and Kit, 2012; Garcia et al., 2014, 2017; Bawden et al., 2018; Voita et al., 2019). Various methods have been proposed to extend sentence-level models to capture document-level context (Tiedemann and Scherrer, 2017; Libovický and Helcl, 2017; Bawden et al., 2018; Miculicich et al., 2018; Sun et al., 2022). The emergence of LLMs, such as GPT-3 (Brown et al., 2020), Llama (Dubey et al., 2024) and Gemma (Gemma Team et al., 2024), has transformed NLP, including for MT (Zhu et al., 2024a,c; Lu et al., 2024). Pre-trained on vast amounts of text, LLMs can effectively manage long-range de-

Language	Classification	Spkrs. (M)
Amharic [amh]	Afro-Asiatic/Semitic	57.6
Hausa [hau]	Afro-Asiatic/Chadic	78.5
Swahili [swa]	Niger-Congo/Bantu	71.6
Yorùbá [yor]	Niger-Congo/Volta-Niger	45.9
isiZulu [zul]	Niger-Congo/Bantu	27.8

Table 2: Languages in the AFRIDOC-MT corpus, their classification and number of speakers (in millions).

dependencies, making them in principle well-suited for document-level translation. While these models have shown promising results for high-resource languages (Wu et al., 2023; Wang et al., 2023; Wu et al., 2024), research remains limited for low-resource languages (Ul Haq et al., 2020).

3 AFRIDOC-MT Corpus

Languages and their characteristics We cover five languages from the two most common African language families: Afro-Asiatic and Niger-Congo. Three languages belong to the Niger-Congo family: Swahili (North-East Bantu), Yorùbá (Volta-Niger) and isiZulu (Southern Bantu). The other two languages belong to the Afro-Asiatic family: Amharic (Semitic) and Hausa (Chadic). The choice of languages was based on geographical representation, speaking population, and web coverage (which we consider as a proxy for the potential performance of existing models on these languages). Details are in Table 2 and Appendix A.

Data Collection and Preprocessing We scraped English articles from the websites of Techpoint Africa² and the World Health Organization (WHO).^{3,4} The articles cover different topics of different lengths with an average length of 30 and 37 sentences for *health* and *tech* respectively. While our corpus is initially structured at the article level, we aim to make it suitable for sentence-level translation tasks as well. To achieve this, we segmented the raw articles into sentences using NLTK (Bird et al., 2009). To ensure high segmentation quality, we recruited a linguist and a professional translator to verify the correctness of the segmentation and made corrections as needed. Finally, we selected 334 and 271 English articles/documents from the *health* and *tech* domains respectively, which represents 10k sentences each per domain.

²<https://techpoint.africa/>

³<https://www.who.int/health-topics>

⁴<https://www.who.int/news-room/>

Domain	Train	Dev.	Test	Min/Max/Avg
Number of documents				
<i>health</i>	240	33	61	2/151/29.9
<i>tech</i>	187	25	59	8/247/36.9
Number of sentences				
<i>health</i>	7041	977	1982	-
<i>tech</i>	7048	970	1982	-

Table 3: The number of documents and sentences in AFRIDOC-MT, and (at the document level) minimum, maximum and average sentences per document.

Translation We translated the extracted 10k English sentences to the 5 African languages through 4 expert translators per language.⁵ The translators were recruited through a language coordinator who is also a native speaker of the language. The 10k sentences were distributed equally among the translators and the translations were done in-context (i.e. the translators translated on the sentence level but had access to the the whole document). Due to the domain-specific nature of the task, before starting the translation process, we conducted a translation workshop, during which three translation experts shared their experiences in creating terminologies and they also shared existing resources with the translators including a short translation guideline (Appendix B.1).

Quality Checks Quality control was conducted using automated quality estimation, followed by manual inspections by our language coordinators. We also used Quality Estimation (QE), specifically AfriCOMET (Wang et al., 2024a), to assess translation quality. Translations scoring below 0.65 were jointly reviewed by translators and language coordinators (see Appendix B.2).

AFRIDOC-MT data split We created train, development (dev), and test splits for each domain. To prevent data leakage, we first identified documents that shared sentences with the same English translation and assigned these documents to the training set. The dev and test sets are derived from the remaining documents. The dev set comprises 25 to 33 documents while the test set includes 59 to 61 documents. Table 3 shows some data statistics, and we provide more data statistics in Appendix B.⁶

⁵Each translator was paid \$1,250 for 2,500 sentences.

⁶Anonymous repo. Public release upon acceptance under CC BY-NC-SA 3.0 for *health*, and CC BY-NC 4.0 for *tech*.

4 Benchmark Experiments

Given the AFRIDOC-MT data, we conducted both sentence- and document-level translation, evaluating two types of models: encoder-decoder and decoder-only models. While the majority of these models are open-source, we also evaluated two proprietary models of the same type. Our evaluation primarily focuses on document-level translation, reflecting the availability of our document-level translation corpus. For completeness, we also conduct a series of sentence-level experiments, with the results presented in Appendix D.

4.1 Models

Encoder-Decoder Models We evaluate five kinds of open encoder-decoder model including Toucan (Elmadany et al., 2024; Adebara et al., 2024), M2M-100 (Fan et al., 2020), NLLB-200 (NLLB Team et al., 2024), MADLAD-400 (Kudugunta et al., 2023), and Aya-101 (Üstün et al., 2024). Toucan is an Afro-centric multilingual MT model supporting 150 African language pairs. In comparison, M2M-100, NLLB-200, and MADLAD-400 cover between 100 and 450 language pairs. Aya-101, an instruction-tuned mT5 model (Xue et al., 2021), supports 100 languages and can translate between various languages, including those considered in AFRIDOC-MT.

Decoder-only Models We also evaluate open and closed decoder-only models. Open models include LLaMA3.1 (Dubey et al., 2024), Gemma2 (Gemma Team et al., 2024), their instruction-tuned variants, and LLaMAX3 (Lu et al., 2024)—a LLaMA3-based model further pre-trained on 100+ languages, including several African ones. The closed models include OpenAI GPT models (GPT-3.5 Turbo and GPT-4o) (OpenAI, 2024), which have been shown to have document-level translation ability (Wang et al., 2023). While their language coverage is not well documented, they show some understanding of African languages (Adelani et al., 2024b; Bayes et al., 2024), though far below their performance in English, their primary training language.

We present the result of 12 models in total, including the 1.2B version of Toucan, 1.3B and 3.3B versions of NLLB-200, 3B and 13B versions of MADLAD-400 and Aya-101 respectively. We also have the 8B instruction tuned version of LLaMA3.1 (LLaMA3.1-IT), 9B version of Gemma-2 (Gemma2-

IT), and LLaMAX3-Alpaca.⁷ We provide more description of the models in Appendix C.1

Supervised fine-tuning of the models (SFT)

For sentence-level evaluation, we jointly fine-tune NLLB-200 with 1.3B parameters on the 30 language directions and on the two domains to make the models more specialized. Similarly, we did supervised fine-tuning on LLaMAX3 and LLaMA3.1 using the prompt augmentation approach from (Zhu et al., 2024b), and shown in Appendix C.4. We chose these two models because LLaMAX3 is already adapted to several languages including our languages of interest, and LLaMA3.1 because of its long context window. We perform SFT on LLaMAX3 and LLaMA3.1 for document-level translation, using pseudo-documents with $k=10$. We refer to each system as $\{\text{model_name}\}\text{-SFT}_k$.⁸

4.2 Experimental Setup

Sentence-level Evaluation Given that our created dataset can be used for sentence-level translation and as a baseline for document-level translation, we evaluate all models on the test splits for each domain. We evaluate the translation models (M2M-100, NLLB-200, and MADLAD-400) using the Fairseq (Ott et al., 2019) codebase for (M2M-100 and NLLB-200), and the Transformers (Wolf et al., 2020) codebase for MADLAD-400. However, for other models including Aya-101, we use the EleutherAI LM Evaluation Harness (lm-eval) tool (Biderman et al., 2024) using the three templates listed in Table 23 of Appendix C.4.

Document-level Evaluation We also perform document-level translation using a setup similar to the sentence-level experiment, but only with models that meet context length requirements. An initial analysis showed that some models were unable to process entire documents due to input length limits, which were exceeded by token counts in some languages (Amharic and Yorùbá). To address this, we adopted a similar approach to Lee et al. (2022), splitting documents into fixed-size chunks of k sentences to fit within token limits; the final chunk may contain fewer than k sentences. To select an appropriate chunk size, we conducted initial tests with $k = 1$ (sentence-level), 5, 10, and 25, choosing $k = 10$ for our experiments. We provide results from this analysis in Table 11.

⁷We refer to it as LLaMAX3-Alp in the results tables.

⁸We denote models finetuned on sentences as $\{\text{model_name}\}\text{-SFT}$ or $\{\text{model_name}\}\text{-SFT}_1$

Model	Size	$eng \rightarrow X$						$X \rightarrow eng$						AVG
		amh	hau	swa	yor	zul	Avg.	amh	hau	swa	yor	zul	Avg.	
Encoder-Decoder														
Toucan	1.2B	33.8 _{1.2}	57.6 _{1.4}	70.3 _{0.8}	36.0 _{1.5}	58.0 _{1.0}	51.2	54.7 _{1.0}	57.7 _{1.3}	65.2 _{0.9}	54.0 _{1.2}	59.9 _{0.8}	58.3	54.7
NLLB-200	1.3B	49.8 _{1.5}	64.7 _{2.2}	75.5 _{0.8}	45.1 _{1.0}	69.0 _{1.3}	60.8	69.4 _{1.3}	65.3 _{1.7}	75.3 _{0.8}	66.3 _{1.1}	73.2 _{0.9}	69.9	65.4
MADLAD-400	3B	36.5 _{0.9}	54.4 _{2.0}	74.2 _{0.9}	19.1 _{0.9}	57.1 _{1.4}	48.3	68.9 _{1.1}	63.8 _{1.6}	76.1 _{0.6}	51.4 _{1.8}	68.9 _{0.9}	65.8	57.0
NLLB-200	3.3B	53.0 _{1.9}	65.2 _{2.2}	76.7 _{0.7}	43.8 _{1.1}	70.7 _{1.3}	61.9	70.9 _{1.3}	66.5 _{1.7}	77.0 _{0.7}	67.6 _{1.1}	74.7 _{1.0}	71.3	66.6
Aya-101	13B	36.6 _{0.9}	56.4 _{1.5}	44.7 _{2.4}	31.2 _{1.4}	58.6 _{0.8}	45.5	64.6 _{1.1}	61.5 _{1.4}	70.8 _{0.8}	57.9 _{1.3}	67.4 _{0.8}	64.4	55.0
SFT on AFRI DOC-MT														
NLLB-SFT	1.3B	55.9_{1.6}	67.4_{1.9}	81.3_{0.7}	61.5_{1.0}	73.7_{1.6}	68.0	72.4_{1.2}	67.5_{1.6}	79.2_{0.7}	71.8_{1.1}	76.5_{0.9}	73.5	70.7
Decoder-only														
Gemma2-IT	9B	20.1 _{0.7}	56.4 _{1.4}	71.2 _{0.7}	21.0 _{0.6}	41.6 _{1.1}	42.1	61.6 _{0.9}	62.5 _{1.3}	74.2 _{0.7}	54.7 _{1.3}	63.9 _{0.9}	63.4	52.7
LLama3.1-IT	8B	19.6 _{0.5}	45.9 _{1.4}	63.7 _{0.9}	19.7 _{0.6}	28.5 _{0.7}	35.5	53.9 _{0.9}	59.8 _{1.3}	69.1 _{0.9}	53.4 _{1.3}	54.0 _{1.1}	58.0	46.8
LLaMAX3-Alp	8B	30.5 _{0.8}	56.3 _{1.5}	67.8 _{0.8}	19.3 _{0.8}	56.1 _{0.9}	46.0	63.3 _{1.0}	62.4 _{1.3}	71.7 _{0.8}	56.1 _{1.1}	65.3 _{0.9}	63.8	54.9
GPT-3.5	–	20.4 _{0.6}	44.3 _{0.9}	76.7 _{0.6}	21.3 _{0.9}	51.1 _{0.9}	42.8	48.3 _{0.9}	52.4 _{1.2}	75.0 _{0.6}	52.1 _{1.2}	59.5 _{0.9}	57.4	50.1
GPT-4o	–	36.7 _{0.8}	64.2 _{1.9}	79.8 _{0.6}	29.3 _{1.6}	69.0 _{1.3}	55.8	67.2 _{1.0}	66.5 _{1.5}	78.1 _{0.6}	69.1 _{1.1}	75.1 _{1.0}	71.2	63.5
SFT on AFRI DOC-MT														
LLaMAX3-SFT	8B	46.8 _{1.2}	62.5 _{1.4}	73.1 _{0.9}	57.5 _{1.0}	67.5 _{1.0}	61.5	66.6 _{1.2}	58.9 _{1.6}	73.1 _{1.1}	64.7 _{1.5}	70.5 _{1.0}	66.8	64.1
LLama3.1-SFT	8B	45.6 _{1.1}	61.8 _{1.5}	71.5 _{1.0}	57.0 _{1.1}	66.8 _{0.9}	60.6	64.3 _{1.2}	59.5 _{1.5}	72.1 _{0.8}	64.8 _{1.5}	69.0 _{1.0}	65.9	63.2

Table 4: Performance of the models in the Health domain, measured by d-CHRF at the sentence-level, realigned to the document-level. For each model and language, the best result from three prompt variations is reported.

Model	Size	$eng \rightarrow X$						$X \rightarrow eng$						AVG
		amh	hau	swa	yor	zul	Avg.	amh	hau	swa	yor	zul	Avg.	
Encoder-Decoder														
Toucan	1.2B	32.0 _{1.6}	59.5 _{1.7}	66.1 _{1.7}	37.1 _{2.0}	58.5 _{1.4}	50.7	54.0 _{1.6}	59.9 _{1.5}	64.1 _{1.4}	54.3 _{1.3}	59.6 _{1.2}	58.4	54.5
NLLB-200	1.3B	49.3 _{2.0}	65.7 _{2.2}	72.3 _{1.6}	43.0 _{1.3}	70.3 _{1.3}	60.1	69.5 _{1.0}	66.8 _{1.5}	72.0 _{1.4}	63.0 _{1.2}	71.5 _{1.2}	68.5	64.3
MADLAD-400	3B	37.3 _{1.3}	57.0 _{2.8}	62.1 _{2.9}	21.3 _{1.0}	58.5 _{1.8}	47.3	68.6 _{1.1}	66.0 _{1.4}	72.1 _{1.4}	53.1 _{1.4}	67.6 _{1.2}	65.5	56.4
NLLB-200	3.3B	52.2 _{2.4}	65.4 _{2.3}	72.8 _{1.5}	40.1 _{1.8}	71.6 _{1.3}	60.4	70.9 _{1.0}	67.7 _{1.5}	73.2 _{1.4}	63.9 _{1.1}	72.5 _{1.2}	69.6	65.0
Aya-101	13B	37.3 _{1.1}	58.9 _{2.3}	42.4 _{2.6}	31.4 _{1.4}	58.9 _{1.5}	45.8	65.2 _{1.2}	64.8 _{1.2}	69.1 _{1.1}	58.5 _{1.3}	67.1 _{1.1}	64.9	55.4
SFT on AFRI DOC-MT														
NLLB-SFT	1.3B	53.4_{2.4}	67.9_{2.2}	76.5_{1.6}	59.5_{1.3}	74.0_{1.5}	66.2	72.1_{1.0}	69.0 _{1.3}	74.1 _{1.4}	67.5_{1.1}	74.3_{1.1}	71.4	68.8
Decoder-only														
Gemma2-IT	9B	20.6 _{0.6}	58.3 _{1.5}	68.7 _{1.6}	23.9 _{1.3}	46.5 _{1.8}	43.6	61.1 _{1.3}	65.4 _{1.4}	71.5 _{1.2}	56.7 _{1.3}	63.8 _{1.1}	63.7	53.7
LLama3.1-IT	8B	19.5 _{0.9}	47.8 _{1.3}	63.4 _{1.5}	20.8 _{1.2}	30.4 _{1.3}	36.4	51.0 _{1.3}	61.0 _{1.4}	66.0 _{1.3}	53.5 _{1.2}	52.4 _{1.3}	56.8	46.6
LLaMAX3-Alp	8B	30.3 _{1.1}	58.9 _{1.9}	64.9 _{1.7}	22.0 _{0.8}	58.6 _{1.7}	46.9	63.4 _{1.4}	64.9 _{1.5}	69.1 _{1.1}	56.5 _{1.3}	65.7 _{1.2}	63.9	55.4
GPT-3.5	–	22.6 _{0.8}	49.2 _{1.5}	72.6 _{1.6}	23.0 _{1.0}	53.6 _{1.5}	44.2	47.4 _{1.5}	56.5 _{1.3}	71.5 _{1.4}	54.0 _{1.3}	59.9 _{1.1}	57.9	51.0
GPT-4o	–	36.9 _{1.2}	65.2 _{2.3}	75.3 _{1.6}	29.4 _{1.5}	71.1 _{1.4}	55.6	67.2 _{1.0}	69.1_{1.4}	74.4_{1.4}	66.4 _{1.1}	73.4 _{1.1}	70.1	62.8
SFT on AFRI DOC-MT														
LLaMAX3-SFT	8B	42.8 _{1.5}	62.4 _{1.9}	67.6 _{1.4}	55.2 _{1.5}	66.0 _{1.2}	58.8	63.0 _{1.2}	53.5 _{1.9}	67.5 _{1.2}	57.3 _{1.3}	66.8 _{1.3}	61.6	60.2
LLama3.1-SFT	8B	41.6 _{1.7}	61.8 _{2.0}	66.4 _{1.3}	54.9 _{1.4}	64.6 _{1.6}	57.9	62.0 _{1.2}	58.6 _{1.5}	67.1 _{1.2}	61.3 _{1.3}	65.6 _{1.3}	62.9	60.4

Table 5: Performance of the models in the Tech domain, measured by d-CHRF at the sentence-level, realigned to the document-level. For each model and language, the best result from three prompt variations is reported.

4.3 Evaluation Metrics

Evaluating document-level translation remains challenging, as existing automatic metrics struggle to capture improvements and identify discourse phenomena (Jiang et al., 2022; Dahan et al., 2024), and embedding-based metrics have been explored for African languages. We realigned sentence-level or pseudo-translation outputs into full documents, then computed BLEU and chrF to create document BLEU (d-BLEU) (Papineni et al., 2002) and document chrF (d-chrF) (Popović, 2015). Metrics were computed using SacreBLEU⁹ (Post, 2018) with bootstrap resampling ($n = 1000$) to report 95% confidence intervals. We report d-chrF scores for the best prompt per model and language direction in the main text, as chrF better captures the morphological richness of African languages (Ade-lani et al., 2022), with full results provided in Appendix D.

⁹case:mixed|eff:no|tok:13a|smooth:exp|v:2.3.1

We use GPT-4o as a judge to evaluate translation outputs, following recent work showing LLMs’ effectiveness in assessing translation quality and analyzing errors (Wu et al., 2024; Sun et al., 2025). Following Sun et al. (2025), we assess translated document’s fluency, content errors (CE), and cohesion errors—specifically lexical (LE) and grammatical (GE) errors—using GPT-4o, with evaluation limited to a few model outputs due to cost constraints (Appendix C.6). We also complement this with human evaluation for direct assessment scores (Appendix C.7) and qualitative analysis through manual inspection (Appendix C.8).

5 Results

5.1 Sentence-level Evaluation

In Tables 4 and 5 we present d-chrF scores based on the realigned documents, created by merging the translated sentences into their corresponding documents. We highlight our main findings below, and sentence-level evaluation results using

Model	Size	$eng \rightarrow X$						$X \rightarrow eng$						AVG
		amh	hau	swa	yor	zul	Avg.	amh	hau	swa	yor	zul	Avg.	
Encoder-Decoder														
MADLAD-400	3B	27.5 _{1.8}	40.2 _{2.3}	46.6 _{3.4}	15.1 _{0.8}	43.6 _{2.6}	34.6	63.3 _{1.6}	62.5 _{2.0}	74.4 _{0.9}	44.2 _{1.6}	66.6 _{1.5}	62.2	48.4
Aya-101	13B	28.7 _{1.6}	48.5 _{2.3}	34.7 _{3.4}	18.7 _{1.3}	54.9 _{1.4}	37.1	61.6 _{1.7}	62.3 _{1.8}	71.2 _{0.9}	56.1 _{2.1}	69.0 _{1.0}	64.0	50.6
Decoder-only														
Gemma2-IT	9B	6.5 _{0.6}	37.0 _{3.4}	52.9 _{3.6}	6.4 _{0.5}	12.0 _{1.0}	23.0	36.5 _{3.0}	51.8 _{3.4}	65.0 _{3.0}	44.8 _{2.9}	56.1 _{3.3}	50.8	36.9
LLama3.1-IT	8B	7.5 _{0.5}	14.0 _{1.2}	43.2 _{3.9}	6.4 _{0.7}	8.7 _{0.6}	16.0	23.8 _{2.3}	49.3 _{4.1}	62.8 _{3.3}	31.7 _{3.9}	34.0 _{3.7}	40.3	28.1
LLaMAX3-Alp	8B	11.4 _{0.9}	28.9 _{2.9}	40.4 _{3.2}	9.2 _{0.8}	23.6 _{1.8}	22.7	29.2 _{2.1}	41.7 _{3.8}	55.4 _{4.9}	23.5 _{3.0}	40.5 _{4.7}	38.1	30.4
GPT-3.5	-	11.6 _{0.5}	23.1 _{2.0}	76.1 _{0.6}	10.1 _{0.9}	29.2 _{2.1}	30.0	41.6 _{2.3}	52.7 _{1.5}	77.7 _{0.6}	51.7 _{1.6}	61.1 _{1.1}	56.9	43.5
GPT-4o	-	29.6 _{1.7}	63.8 _{1.9}	80.2 _{0.6}	29.6 _{2.1}	69.5 _{1.6}	54.5	69.5 _{1.1}	69.3 _{1.7}	81.0 _{0.6}	73.8 _{1.0}	78.2 _{1.1}	74.4	64.4
SFT on AFRIDOC-MT														
LLaMAX3-SFT	8B	24.1 _{1.6}	29.0 _{3.2}	42.2 _{4.2}	33.8 _{2.8}	33.7 _{3.1}	32.6	22.6 _{1.8}	22.9 _{2.6}	33.1 _{4.4}	27.2 _{3.6}	31.5 _{6.7}	27.5	30.0
LLama3.1-SFT	8B	25.2 _{1.8}	31.9 _{4.0}	50.2 _{6.4}	33.8 _{2.8}	38.6 _{4.1}	35.9	24.2 _{3.7}	24.1 _{4.1}	33.7 _{5.4}	30.2 _{4.7}	29.3 _{6.2}	28.3	32.1
LLaMAX3-SFT ₁₀	8B	37.8 _{2.2}	51.9 _{5.0}	74.4 _{3.5}	52.2 _{3.3}	55.0 _{5.5}	54.2	64.0 _{3.4}	66.7 _{2.8}	77.8 _{0.7}	71.8 _{1.0}	74.1 _{0.9}	70.9	62.6
LLama3.1-SFT ₁₀	8B	27.6 _{2.4}	49.7 _{5.2}	64.1 _{5.6}	50.3 _{2.8}	47.0 _{4.8}	47.8	63.8 _{1.1}	61.7 _{3.5}	74.4 _{3.5}	68.9 _{3.4}	71.4 _{1.0}	68.0	57.9

Table 6: Performance results of various models on the pseudo-documents ($k=10$) translation task (Health domain), measured using d-CHRf. The best prompt was selected for each language after evaluating three different prompts.

Model	Size	$eng \rightarrow X$						$X \rightarrow eng$						AVG
		amh	hau	swa	yor	zul	Avg.	amh	hau	swa	yor	zul	Avg.	
Encoder-Decoder														
MADLAD-400	3B	29.5 _{2.1}	38.3 _{4.3}	31.7 _{4.6}	15.1 _{1.1}	44.1 _{3.6}	31.8	62.6 _{2.0}	63.5 _{2.2}	66.4 _{3.2}	45.9 _{2.4}	63.4 _{2.2}	60.3	46.0
Aya-101	13B	30.1 _{1.5}	55.0 _{3.2}	51.7 _{3.5}	22.3 _{1.7}	55.0 _{1.9}	42.8	62.5 _{1.4}	65.5 _{1.3}	68.8 _{1.8}	55.7 _{2.4}	68.4 _{1.0}	64.2	53.5
Decoder-only														
Gemma2-IT	9B	6.2 _{0.7}	42.1 _{3.9}	51.0 _{5.3}	6.6 _{0.8}	15.4 _{1.7}	24.3	35.9 _{4.8}	50.1 _{4.6}	57.7 _{3.7}	48.2 _{3.4}	51.7 _{3.7}	48.7	36.5
LLama3.1-IT	8B	7.4 _{0.9}	15.3 _{1.9}	43.3 _{4.4}	6.2 _{1.1}	8.8 _{0.7}	16.2	26.1 _{2.0}	48.7 _{3.4}	59.0 _{2.7}	34.4 _{3.2}	34.7 _{3.1}	40.6	28.4
LLaMAX3-Alp	8B	11.4 _{1.2}	32.5 _{4.4}	38.1 _{4.1}	12.0 _{1.4}	26.1 _{2.2}	24.0	29.4 _{2.9}	51.4 _{4.3}	62.4 _{2.5}	24.7 _{3.6}	48.8 _{5.3}	43.3	33.7
GPT-3.5	-	13.5 _{1.1}	29.7 _{2.5}	72.1 _{1.6}	12.7 _{1.2}	35.1 _{2.9}	32.6	38.5 _{4.0}	56.3 _{1.5}	73.5 _{1.4}	53.0 _{1.6}	61.2 _{1.3}	56.5	44.6
GPT-4o	-	31.3 _{1.9}	65.1 _{2.5}	75.1 _{1.6}	28.1 _{1.8}	70.7 _{1.5}	54.0	68.6 _{1.1}	71.6 _{1.4}	76.5 _{1.6}	70.1 _{1.1}	76.5 _{1.1}	72.7	63.3
SFT on AFRIDOC-MT														
LLaMAX3-SFT	8B	21.7 _{2.0}	29.9 _{3.2}	37.0 _{3.4}	30.5 _{2.7}	31.7 _{3.5}	30.2	24.2 _{2.6}	27.6 _{4.2}	32.3 _{4.5}	28.5 _{3.3}	29.8 _{5.4}	28.5	29.3
LLama3.1-SFT	8B	21.0 _{2.0}	30.8 _{3.2}	40.0 _{4.1}	33.4 _{3.8}	29.3 _{3.1}	30.9	23.9 _{2.5}	28.9 _{4.3}	36.9 _{5.8}	32.2 _{4.3}	32.3 _{5.2}	30.8	30.9
LLaMAX3-SFT ₁₀	8B	37.7 _{2.1}	58.6 _{5.1}	68.3 _{3.9}	49.3 _{4.1}	60.9 _{3.9}	55.0	65.4 _{1.4}	68.5 _{1.3}	73.1 _{1.2}	67.7 _{1.2}	71.6 _{1.2}	69.3	62.1
LLama3.1-SFT ₁₀	8B	23.7 _{1.9}	47.0 _{5.2}	58.6 _{5.6}	49.7 _{3.8}	43.8 _{4.5}	44.5	60.9 _{2.7}	65.4 _{2.5}	71.1 _{1.2}	66.3 _{1.2}	66.4 _{4.0}	66.0	55.3

Table 7: Performance results of various models on the pseudo-documents ($k=10$) translation task (Tech domain), measured using d-CHRf. The best prompt was selected for each language after evaluating three different prompts.

sentence-level metrics are reported in Appendix D.

NLLB-200 outperforms all other encoder-decoder models across languages and domains

On average the NLLB models obtain scores of 65.4/66.6 and 64.3/65.0 on *health* and *tech* domains respectively, with 3.3B outperforming 1.3B except when translating into Yorùbá. When translating to English, the worst performing model across the two domains is Toucan. However, it gives better results than MADLAD-400 and Aya-101 when translating to African languages. Furthermore, translating to African languages is significantly worse compared to translating to English for all the models.

GPT-4o outperforms other decoder-only counterparts

GPT-4o on average outperforms other decoder-only LMs, with average d-chrF scores of 63.5 and 62.8 for health and tech respectively. The next best performing decoder-only model is LLaMAX3-Alpaca, with d-chrF scores of 54.9 and 55.4. Unlike other open decoder-based LLMs, LLaMAX3-Alpaca was trained on African languages through continued pretraining and adapted via instruction tuning. It outperforms Gemma2-IT by +2.2 in the health domain and +1.7 in the *tech*

domain, particularly when translating into African languages. In contrast, GPT-3.5 and LLama3.1-IT are the worst performing models.

Fine-tuning models significantly improves translation quality

We obtain improved performance after fine-tuning NLLB-1.3B on AFRIDOC-MT, and the resulting model outperforms the 3.3B version without fine-tuning. Similarly, the SFT-based LLMs (LLaMAX3 and LLama3.1) become the best performing open LLMs and outperform their base-lines (LLaMAX3-Alpaca and LLama3.1-IT) but below GPT-4o. Overall, our fine-tuned NLLB-200 model is the state-of-the-art model, and our fine-tuned LLaMAX3 is competitive to GPT-4o.

5.2 Document-level Evaluation

In Tables 6 and 7 we present d-chrF scores based on the best prompt per language for the the translation output of the models when evaluated on the realigned documents from pseudo-documents with $k=10$ sentences per pseudo-document.

Pseudo-document translation is worse than sentence-level translation when translating into African languages Our results from pseudo-

Model	Setup	$eng \rightarrow X$					$X \rightarrow eng$				
		d-CHRf \uparrow	Fluency \uparrow	CE \downarrow	LE \downarrow	GE \downarrow	d-CHRf \uparrow	Fluency \uparrow	CE \downarrow	LE \downarrow	GE \downarrow
Aya-101	Sent	45.5	12.0	2.2	0.8	9.9	2.2	4.9	0.9	3.2	0.5
	Doc10	37.1	14.7	2.3	0.7	9.6	3.7	3.3	1.1	2.4	0.5
GPT-3.5	Sent	42.8	23.3	2.0	1.6	10.3	6.5	6.7	3.7	4.0	2.0
	Doc10	30.0	26.9	1.9	1.7	5.8	2.9	2.4	1.2	2.1	1.0
LLaMAX3-SFT ₁	Sent	61.5	10.0	3.5	0.3	11.3	1.4	4.2	1.0	3.1	0.7
	Doc10	32.6	6.7	2.5	0.5	9.0	0.9	2.6	0.6	2.1	0.3
LLaMAX3-SFT ₁₀	Sent	54.2	13.1	3.8	0.7	11.0	2.7	2.6	0.7	1.9	0.3
	Doc10	64.4	5.0	2.9	0.4	18.7	3.1	11.1	1.7	6.2	1.5
GPT-3.5	Sent	64.0	6.1	3.4	0.3	15.1	2.3	9.5	1.0	4.4	0.5
	Doc10	57.4	10.6	2.9	0.5	12.8	2.3	6.4	2.0	3.8	1.2
LLaMAX3-SFT ₁	Sent	56.9	13.5	4.2	0.3	8.5	1.7	3.9	1.3	2.1	0.6
	Doc10	66.8	5.5	3.4	0.4	11.5	1.1	6.4	1.6	3.1	0.4
LLaMAX3-SFT ₁₀	Sent	27.5	4.8	3.0	0.2	8.9	0.2	3.1	0.4	1.9	0.2
	Doc10	70.9	5.6	4.3	0.2	9.4	0.6	5.3	0.8	2.6	0.1

Table 8: Document-level evaluation in the *health* domain, judged by GPT-4o. Compares sentence- vs. document-level outputs on Fluency (1–5 scale), Content Errors (CE), Lexical (LE), and Grammatical Cohesion Errors (GE).

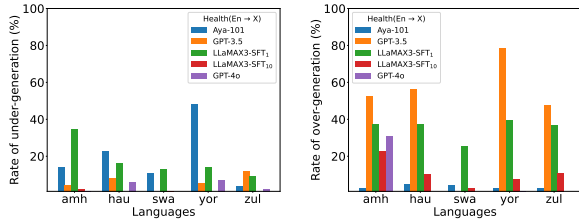


Figure 1: Rate of under-generation and over-generation in pseudo-document translation ($k=10$).

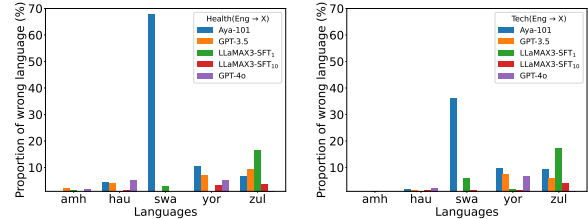


Figure 2: Rate of off-target translation ($k=10$).

document translation show a performance drop across different models compared to sentence-level translation, especially when translating into African languages. However, GPT-4o demonstrates similar and consistent performance in both setups and domains. Additionally, we observe that GPT-3.5 is the next best performing decoder-only LLM, which contrasts with its performance in sentence-level translation. Gemma2-IT outperforms LLaMAX3-Alpaca especially when translating into English, which also differs from the trends observed in the sentence-level setup.

LLMs trained on longer documents are better for long document translation Both Llama models trained via SFT on sentences (LLama3.1-SFT, and LLaMAX3-SFT) show a decline in performance in the pseudo-document setting compared to sentence-level translation. However, the same models trained via SFT on pseudo-documents with $k=10$ demonstrate significant improvements on pseudo-documents. Interestingly, the LLaMAX3-SFT₁₀ model performs consistently well, achieving results comparable to its sentence-level counterpart on sentence-level tasks, and also outperforming LLama3.1-SFT₁₀, particularly when translating into African languages.

5.3 GPT-4o based evaluation

Table 8 presents average GPT-4o evaluations of realigned outputs from sentence level and pseudo-document level tasks ($k=10$) across four models

in the *health* domain. Pseudo-document outputs are rated more fluent and show fewer content, lexical, and grammatical errors across both translation directions. However, these results often contradict d-chrF scores, especially when translating into African languages. For instance, GPT-3.5 has the lowest d-chrF yet the fewest errors—raising concerns about GPT-4o’s reliability. Hence, we focus on human evaluation going forward. Full GPT-4o results are provided in Appendix D.3.

5.4 Human evaluation

We report average direct assessment (DA) scores (on a scale from 0 to 100) from three¹⁰ annotators per language for the health domain in Table 9, when translating into four African languages. For each language, we used 30 documents across models and both domains to compute inter-annotator agreement. We obtained Krippendorff’s alpha values of ≥ 0.46 , which are relatively low due to the fine granularity of the evaluation scale. **Human evaluation results align closely with d-chrF**, which favors sentence-level translations over pseudo-document translations when translating into African languages. Among the models, LLaMAX3-SFT₁ receives higher ratings at the sentence level but is rated lower when translating pseudo-documents. In contrast, LLaMAX3-SFT₁₀ receives slightly lower ratings than LLaMAX3-SFT₁ at the sentence level but is rated higher in the pseudo-document setting. GPT-3.5 is gener-

¹⁰except for Swahili where we had just 2 annotators

Model	Setup	amh	hau	swh	yor
GPT-3.5	Sent	14.6	29.6	72.0	7.5
	Doc10	1.7	16.4	74.0	4.2
LLaMAX3-SFT ₁	Sent	64.5	81.5	68.8	65.1
	Doc10	27.4	45.7	50.2	44.3
LLaMAX3-SFT ₁₀	Doc10	38.5	76.7	67.4	64.9

Table 9: Average DA score (scale 0–100) from the human evaluators per language in the *health* domain.

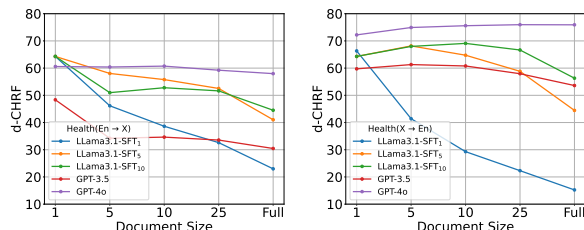


Figure 3: Comparison of Average d-chrF scores across models and pseudo-document lengths.

ally rated as the weakest model across languages, except for Swahili, where its performance is comparatively better (see Appendix D.4 for details).

5.5 Qualitative evaluation

Our qualitative analysis, based on author feedback, indicates that GPT-3.5 frequently over-generates in the pseudo-document setup by repeating words and phrases—except in the case of Swahili, where it ranks highest. However, for Yorùbá, it often uses inconsistent or partial diacritics, resulting in inaccuracies. LLaMAX3-SFT₁ also exhibits repetition in pseudo-document translations, likely due to a length generalization problem (Anil et al., 2022), and does so more than LLaMAX3-SFT₁₀. For the other four languages, LLaMAX3-SFT₁ with the sentence-level setup was rated higher than other models and configurations, owing to better context preservation and fewer repetitions. These observations are consistent with both d-chrF and DA scores, although d-chrF scores tend to be inflated.

6 Discussion and Analysis

To better understand model behavior, we analyze their pseudo-document ($k = 10$) translation outputs based on our findings and common issues in document-level MT with LLMs (Wu et al., 2024; Wang et al., 2024b). Additional results are provided in Appendix E.

Are the outputs generated by translation models of appropriate length? We compare model translations to the reference translations to identify empty outputs and cases of under- or over-generation. We found that all models rarely gen-

erate empty translations (refer to Appendix E), although GPT-3.5 and GPT-4o showed a slight tendency to do so for Yorùbá and Zulu, occurring in under 10% of cases. We defined under-generation as outputs $<70\%$ of reference length, and over-generation as $>143\%$. Consistent with our qualitative findings, GPT-3.5 tends to over-generate more in African languages except for Swahili, while LLaMAX3-SFT₁ often under-generates as a sentence-level model. Moreover, all models over-generated by about 20% for Amharic, likely due to its unique script.

Do LLMs generate translations in the correct target languages?

We evaluate whether these models understand the task by generating outputs in the target languages using language identification. Our results show that the models rarely generate outputs in the wrong language when translating to English. However, when translating to African languages, there is a higher likelihood of incorrect language translations, particularly with open models (see Figure 2).

What is the effect of document length on translation quality?

We compare average d-chrF scores of selected models, including GPT-3.5/4 and LLaMA3.1-SFT_k ($k = 1, 5, 10$), evaluated across pseudo-document lengths of 1, 5, 10, 25, and full length. As shown in Figure 3, d-chrF scores generally decline with increasing document length for African language translations. The reverse translation direction shows a similar trend, except for GPT-4o, which improves with length. Models trained on longer documents also generalize better to longer inputs than those trained on sentences.

7 Conclusion

We introduce AFRIDOC-MT, a document-level translation dataset in the *health* and *tech* domains for five African languages. We benchmarked various models, fine-tuning selected ones. Due to context length limits, documents were translated either sentence by sentence or as pseudo-documents. Outputs were evaluated using standard MT metrics, GPT-4o as a judge, and human direct assessment. NLLB-200 was the strongest built-in MT model, while GPT-4o outperformed general-purpose LLMs. However, our DA and qualitative analysis found GPT-4o’s judgments inconsistent for African languages, and sentence-by-sentence translation proved more effective for some languages.

8 Limitations

Choice of LLMs and Prompts We evaluated only a small subset of the numerous multilingual LLMs available. Our experiments were also limited by the context length of the LLMs, particularly for open LLMs. Except for LLama3.1, all other open LLMs have a context length of 8192 tokens, while encoder-decoder models were primarily based on T5. This makes it difficult to use the context length beyond a certain limit, making full document translation infeasible. Additionally, LLMs are prone to variance in performance based on the prompt. We therefore evaluated them for translation using three different prompts. However, it is possible that our prompts were not optimal.

Language Coverage Africa is home to thousands of indigenous languages, many of which exhibit unique linguistic properties. However, due to the high cost of translation using human translators and limited available funding, it is currently impossible to cover all languages. As a result, we focused on just five languages. We hope that future work will expand this dataset to include more languages and inspire the creation of additional datasets with broader coverage for document-level translation. Similarly, AFRIDOC-MT is a multi-way parallel dataset. However, due to the cost of running inference over three prompts and across all 30 translation directions for all the models evaluated, most of our analysis is limited to translation tasks between English and the five African languages. While we fine-tuned NLLB-200, LLama3.1 and LLaMAX3 on all 30 directions, we only provide results from NLLB-200 for all directions both before and after fine-tuning for sentence-level and pseudo-document tasks in the Appendix E.

Evaluation Metrics Quality evaluation in MT is an open and ongoing area of research, especially for document-level translation. Recent works have proposed embedding-based metrics for evaluation at both the sentence and document levels. While this has been well explored for high-resource language pairs, it remains underexplored for African languages, although there is a tool, AfriCOMET, that works for sentence-level evaluation in African languages. However, we evaluated the document-level translation outputs using *ModernBERT-base-long-context-qe-v1*¹¹, trained on the WMT human

evaluation dataset across 41 language pairs, including over 20 languages and three African languages (Hausa, Xhosa, and Zulu), two of which are included to our work. However, the scores were nearly identical across all models, offering no meaningful differentiation. Hence, for our document-level evaluation, in addition to lexical-based metrics, we incorporated three other evaluation approaches: using GPT-4o as a judge, human evaluation, and qualitative analysis. GPT-4o was employed to assess and rate the translation outputs of four models. While its ratings were consistent for translations into English, the same was not observed for translations into African languages, likely due to the model’s limited understanding of these languages. Therefore, we conducted a human evaluation for translations from English to African languages, comparing only three models due to cost constraints. However, we were unable to recruit annotators for Zulu.

Model Coverage and Evaluation Focus While we fine-tuned both NLLB-1.3B and LLaMAX3 models across all 30 language directions, due to computational constraints and the high cost of qualitative evaluation, our detailed analysis focuses only on translation between English and the 5 African languages. Nevertheless, we report quantitative results across all 30 directions for NLLB-1.3B. We will make all fine-tuned models publicly available to support future work, and we hope that further research will explore the remaining translation directions in greater depth.

Translationese and English-Centric Bias A potential limitation of our dataset is the influence of translationese (Koppel and Ordan, 2011). Since all source material translated originates in English, translated sentences in African languages may exhibit patterns such as unnatural syntax or overly literal phrasing. Although we have not conducted an analysis to quantify these effects, prior work suggests that they can affect MT model performance, generalization and evaluation including direct assessment (Freitag et al., 2019; Edunov et al., 2020). Furthermore, AFRIDOC-MT may reflect a bias toward English in terms of structure, semantics, and cultural framing. We leave a deeper investigation of these issues to future work.

¹¹<https://huggingface.co/yomoslem/ModernBERT-base-long-context-qe-v1>

Ethics Statement

AFRIDOC-MT was created with the utmost consideration for ethical standards. The English texts translated were sourced from publicly available and ethically sourced materials. The data sources were selected to represent different cultural perspectives, with a focus on minimizing any potential bias. Efforts were made to ensure the dataset does not include harmful, biased, or offensive content via manual inspection.

References

Ife Adebara, AbdelRahim Elmadany, and Muhammad Abdul-Mageed. 2024. [Cheetah: Natural language generation for 517 African languages](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12798–12823, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

David Adelani, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruitter, Dietrich Klakow, Peter Nabende, Ernie Chang, Tajudeen Gwadabe, Freshia Sackey, Bonaventure F. P. Dossou, Chris Emezue, Colin Leong, Michael Beukman, Shamsuddeen Muhammad, Guyo Jarso, Oreen Yousuf, Andre Niyongabo Rubungo, Gilles Hacheme, Eric Peter Wairagala, Muhammad Umair Nasir, Benjamin Ajibade, Tunde Ajayi, Yvonne Gitau, Jade Abbott, Mohamed Ahmed, Millicent Ochieng, Anuoluwapo Aremu, Perez Ogayo, Jonathan Mukiibi, Fatoumata Ouoba Kabore, Godson Kalipe, Derguene Mbaye, Allahsera Auguste Tapo, Victoire Memdjokam Koagne, Edwin Munkoh-Buabeng, Valencia Wagner, Idris Abdulmumin, Ayodele Awokoya, Happy Buzaaba, Blessing Sibanda, Andiswa Bukula, and Sam Manthalu. 2022. [A few thousand translations go a long way! leveraging pre-trained models for African news translation](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3053–3070, Seattle, United States. Association for Computational Linguistics.

David Adelani, Hannah Liu, Xiaoyu Shen, Nikita Vassilyev, Jesujoba Alabi, Yanke Mao, Haonan Gao, and En-Shiun Lee. 2024a. [SIB-200: A simple, inclusive, and big evaluation dataset for topic classification in 200+ languages and dialects](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 226–245, St. Julian's, Malta. Association for Computational Linguistics.

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.

David Ifeoluwa Adelani. 2022. [Natural language processing for african languages](#).

David Ifeoluwa Adelani, Jessica Ojo, Israel Abebe Azime, Jian Yun Zhuang, Jesujoba O. Alabi, Xuanli He, Millicent Ochieng, Sara Hooker, Andiswa Bukula, En-Shiun Annie Lee, Chiamaka Chukwunkeke, Happy Buzaaba, Blessing Sibanda, Godson Kalipe, Jonathan Mukiibi, Salomon Kabongo, Foutse Yuehgo, Mmasibidi Setaka, Lolwethu Ndolela, Nkiruka Odu, Rooweither Mabuya, Shamsuddeen Hassan Muhammad, Salomey Osei, Sokhar Samb, Tadesse Kebede Guge, and Pontus Stenetorp. 2024b. [Irokobench: A new benchmark for african languages in the age of large language models](#).

Željko Agić and Ivan Vulić. 2019. [JW300: A wide-coverage parallel corpus for low-resource languages](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3204–3210, Florence, Italy. Association for Computational Linguistics.

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondřej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina España-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. [Findings of the 2021 conference on machine translation \(WMT21\)](#). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–88, Online. Association for Computational Linguistics.

Jesujoba O. Alabi, David Ifeoluwa Adelani, Marius Mosbach, and Dietrich Klakow. 2022. [Adapting pre-trained language models to African languages via multilingual adaptive fine-tuning](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4336–4349, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Antonios Anastasopoulos, Alessandro Cattelan, Zi-Yi Dou, Marcello Federico, Christian Federmann, Dmitriy Genzel, Francisco Guzmán, Junjie Hu, Macduff Hughes, Philipp Koehn, Rosie Lazar, Will Lewis, Graham Neubig, Mengmeng Niu, Alp Öktem, Eric Paquin, Grace Tang, and Sylwia Tur. 2020. [TICO-19: the translation initiative for COvid-19](#). In *Proceedings of the 1st Workshop on NLP for COVID-19 (Part*

764	2) at EMNLP 2020, Online. Association for Computational Linguistics.	822
765		823
766	Cem Anil, Yuhuai Wu, Anders Johan Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Venkatesh Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. 2022. Exploring length generalization in large language models . In <i>Advances in Neural Information Processing Systems</i> .	824
767		825
768		826
769		827
770		828
771		829
772	Paul Azunre, Salomey Osei, Salomey Addo, Lawrence Asamoah Adu-Gyamfi, Stephen Moore, Bernard Adabankah, Bernard Opoku, Clara Asare-Nyarko, Samuel Nyarko, Cynthia Amoaba, Esther Dansoa Appiah, Felix Akwerh, Richard Nii Lante Lawson, Joel Budu, Emmanuel Debrah, Nana Boateng, Wisdom Ofori, Edwin Buabeng-Munkoh, Franklin Adjei, Isaac Kojo Essel Ampomah, Joseph Otoo, Reindorf Borkor, Standylove Birago Mensah, Lucien Mensah, Mark Amoako Marcel, Anokye Acheampong Amponsah, and James Ben Hayfron-Acquah. 2021. English-twi parallel corpus for machine translation . In <i>2nd AfricaNLP Workshop Proceedings, AfricaNLP@EACL 2021, Virtual Event, April 19, 2021</i> .	830
773		831
774		832
775		833
776		834
777		835
778		836
779		837
780		838
781		839
782		
783		840
784		841
785		842
786		843
787		844
788		
789		845
790		846
791		847
792		848
793		849
794		850
795		851
796		852
797		853
798		854
799		855
800		856
801		857
802		858
803		859
804		860
805		861
806		862
807		863
808		864
809		865
810		866
811		867
812		868
813		869
814		870
815		871
816		872
817		873
818		874
819		875
820		876
821		877
		878
		879
		880
		881
		882

883	Mathieu Rita, Maya Pavlova, Melanie Kambadur,	Hamid Shojanazeri, Han Zou, Hannah Wang, Han-	946
884	Mike Lewis, Min Si, Mitesh Kumar Singh, Mona	wen Zha, Haroun Habeeb, Harrison Rudolph, He-	947
885	Hassan, Naman Goyal, Narjes Torabi, Nikolay Bash-	len Suk, Henry Aspegren, Hunter Goldman, Igor	948
886	lykov, Nikolay Bogoychev, Niladri Chatterji, Olivier	Molybog, Igor Tufanov, Irina-Elena Veliche, Itai	949
887	Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan	Gat, Jake Weissman, James Geboski, James Kohli,	950
888	Zhang, Pengwei Li, Petar Vasic, Peter Weng, Pra-	Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff	951
889	jjwal Bhargava, Pratik Dubal, Praveen Krishnan,	Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizen-	952
890	Punit Singh Koura, Puxin Xu, Qing He, Qingxiao	stein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi	953
891	Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon	Yang, Joe Cummings, Jon Carvill, Jon Shepard,	954
892	Calderer, Ricardo Silveira Cabral, Robert Stojnic,	Jonathan McPhie, Jonathan Torres, Josh Ginsburg,	955
893	Roberta Raileanu, Rohit Girdhar, Rohit Patel, Ro-	Junjie Wang, Kai Wu, Kam Hou U, Karan Sax-	956
894	main Sauvestre, Ronnie Polidoro, Roshan Sumbaly,	ena, Karthik Prasad, Kartikay Khandelwal, Katay-	957
895	Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar	oun Zand, Kathy Matosich, Kaushik Veeraragha-	958
896	Hosseini, Sahana Chennabasappa, Sanjay Singh,	van, Kelly Michelena, Keqian Li, Kun Huang, Ku-	959
897	Sean Bell, Seohyun Sonia Kim, Sergey Edunov,	nul Chawla, Kushal Lakhotia, Kyle Huang, Lailin	960
898	Shaoliang Nie, Sharan Narang, Sharath Raparthy,	Chen, Lakshya Garg, Lavender A, Leandro Silva,	961
899	Sheng Shen, Shengye Wan, Shruti Bhosale, Shun	Lee Bell, Lei Zhang, Liangpeng Guo, Licheng	962
900	Zhang, Simon Vandenhende, Soumya Batra, Spencer	Yu, Liron Moshkovich, Luca Wehrstedt, Madian	963
901	Whitman, Sten Sootla, Stephane Collot, Suchin Gu-	Khabsa, Manav Avalani, Manish Bhatt, Maria Tsim-	964
902	rurangan, Sydney Borodinsky, Tamar Herman, Tara	poukelli, Martynas Mankus, Matan Hasson, Matthew	965
903	Fowler, Tarek Sheasha, Thomas Georgiou, Thomas	Lennie, Matthias Reso, Maxim Groshev, Maxim	966
904	Scialom, Tobias Speckbacher, Todor Mihaylov, Tong	Naumov, Maya Lathi, Meghan Keneally, Michael L.	967
905	Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor	Seltzer, Michal Valko, Michelle Restrepo, Mihir	968
906	Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent	Patel, Mik Vyatskov, Mikayel Samvelyan, Mike	969
907	Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-	Clark, Mike Macey, Mike Wang, Miquel Jubert Her-	970
908	vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-	moso, Mo Metanat, Mohammad Rastegari, Mun-	971
909	ney Meers, Xavier Martinet, Xiaodong Wang, Xiao-	ish Bansal, Nandhini Santhanam, Natascha Parks,	972
910	qing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei	Natasha White, Navyata Bawa, Nayan Singhal, Nick	973
911	Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine	Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev,	974
912	Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue	Ning Dong, Ning Zhang, Norman Cheng, Oleg	975
913	Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng	Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem	976
914	Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh,	Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pa-	977
915	Aaron Grattafori, Abha Jain, Adam Kelsey, Adam	van Balaji, Pedro Rittner, Philip Bontrager, Pierre	978
916	Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva	Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratan-	979
917	Goldstand, Ajay Menon, Ajay Sharma, Alex Boesen-	chandani, Pritish Yuvraj, Qian Liang, Rachad Alao,	980
918	berg, Alex Vaughan, Alexei Baeviski, Allie Feinstein,	Rachel Rodriguez, Rafi Ayub, Raghotham Murthy,	981
919	Amanda Kallet, Amit Sangani, Anam Yunus, An-	Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah	982
920	drei Lupu, Andres Alvarado, Andrew Caples, An-	Hogan, Robin Battey, Rocky Wang, Rohan Mah-	983
921	drew Gu, Andrew Ho, Andrew Poulton, Andrew	eswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu,	984
922	Ryan, Ankit Ramchandani, Annie Franco, Aparaj-	Samyak Datta, Sara Chugh, Sara Hunt, Sargun	985
923	jita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma,	986
924	Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	987
925	dan, Beau James, Ben Maurer, Benjamin Leonhardi,	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	988
926	Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi	Shengxin Cindy Zha, Shiva Shankar, Shuqiang	989
927	Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-	Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agar-	990
928	cock, Bram Wasti, Brandon Spence, Brani Stojkovic,	wal, Soji Sajuyigbe, Soumith Chintala, Stephanie	991
929	Brian Gamido, Britt Montalvo, Carl Parker, Carly	Max, Stephen Chen, Steve Kehoe, Steve Satterfield,	992
930	Burton, Catalina Mejia, Changan Wang, Changkyu	Sudarshan Govindaprasad, Sumit Gupta, Sungmin	993
931	Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu,	Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury,	994
932	Chris Cai, Chris Tindal, Christoph Feichtenhofer, Da-	Sydney Goldman, Tal Remez, Tamar Glaser, Tamara	995
933	mon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,	Best, Thilo Kohler, Thomas Robinson, Tianhe Li,	996
934	Danny Wyatt, David Adkins, David Xu, Davide Tes-	Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook	997
935	tuggine, Delia David, Devi Parikh, Diana Liskovich,	Shaked, Varun Vontimitta, Victoria Ajayi, Victoria	998
936	Didem Foss, Dingkan Wang, Duc Le, Dustin Hol-	Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal	999
937	land, Edward Dowling, Eissa Jamil, Elaine Mont-	Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu	1000
938	gomery, Eleonora Presani, Emily Hahn, Emily Wood,	Mihailescu, Vladimir Ivanov, Wei Li, Wenchen	1001
939	Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan	Wang, Wenwen Jiang, Wes Bouaziz, Will Constable,	1002
940	Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat	Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu,	1003
941	Ozgenel, Francesco Caggioni, Francisco Guzmán,	Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yan-	1004
942	Frank Kanayet, Frank Seide, Gabriela Medina Flo-	jun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin	1005
943	rez, Gabriella Schwarz, Gada Badeer, Georgia Swee,	Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu,	1006
944	Gil Halpern, Govind Thattai, Grant Herman, Grigory	Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach	1007
945	Sizov, Guangyi, Zhang, Guna Lakshminarayanan,	Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen,	1008

1009	Zhenyu Yang, and Zhiwei Zhao. 2024. The Llama 3 Herd of Models .	1065
1010		1066
1011	Sergey Edunov, Myle Ott, Marc’ Aurelio Ranzato, and Michael Auli. 2020. On the evaluation of machine translation systems trained with back-translation . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2836–2846, Online. Association for Computational Linguistics.	1067
1012		1068
1013		1069
1014		1070
1015		
1016		
1017		
1018	AbdelRahim Elmadany, Ife Adebara, and Muhammad Abdul-Mageed. 2024. Toucan: Many-to-many translation for 150 african language pairs . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 13189–13206, Bangkok, Thailand. Association for Computational Linguistics.	1071
1019		1072
1020		1073
1021		1074
1022		1075
1023		1076
1024	Ignatius Ezeani, Paul Rayson, Ikechukwu E. Onyenwe, Chinedu Uchekchukwu, and Mark Hepple. 2020. Igbo-english machine translation: an evaluation benchmark . In <i>1st AfricaNLP Workshop Proceedings, AfricaNLP@ICLR 2020, Virtual Conference, Formerly Addis Ababa Ethiopia, 26th April 2020</i> .	1077
1025		1078
1026		1079
1027		1080
1028		1081
1029		1082
1030	Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Man-deep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation .	1083
1031		1084
1032		1085
1033		1086
1034		1087
1035		1088
1036		1089
1037	Christian Federmann. 2018. Appraise evaluation framework for machine translation . In <i>Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations</i> , pages 86–88, Santa Fe, New Mexico. Association for Computational Linguistics.	1090
1038		1091
1039		1092
1040		1093
1041		1094
1042		1095
1043	Christian Federmann, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – news test references for MT evaluation of 128 languages . In <i>Proceedings of the First Workshop on Scaling Up Multilingual Evaluation</i> , pages 21–24, Online. Association for Computational Linguistics.	1096
1044		1097
1045		1098
1046		1099
1047		1100
1048		1101
1049	Yukun Feng, Feng Li, Ziang Song, Boyuan Zheng, and Philipp Koehn. 2022. Learn to remember: Transformer with recurrent memory for document-level machine translation . In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 1409–1420, Seattle, United States. Association for Computational Linguistics.	1102
1050		1103
1051		1104
1052		1105
1053		1106
1054		1107
1055		1108
1056	Markus Freitag, Isaac Caswell, and Scott Roy. 2019. APE at scale and its implications on MT evaluation biases . In <i>Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)</i> , pages 34–44, Florence, Italy. Association for Computational Linguistics.	1109
1057		1110
1058		1111
1059		1112
1060		1113
1061		1114
1062	Eva Martínez Garcia, Carles Creus, Cristina Espana-Bonet, and Lluís Màrquez. 2017. Using word embeddings to enforce document-level lexical consistency	1115
1063		1116
1064		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300

1127	Giang, Ludovic Peran, Tris Warkentin, Eli Collins,	1186
1128	Joelle Barral, Zoubin Ghahramani, Raia Hadsell,	1187
1129	D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov,	1188
1130	Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray	1189
1131	Kavukcuoglu, Clement Farabet, Elena Buchatskaya,	1190
1132	Sebastian Borgeaud, Noah Fiedel, Armand Joulin,	1191
1133	Kathleen Kenealy, Robert Dadashi, and Alek An-	1192
1134	dreev. 2024. Gemma 2: Improving open language	1193
1135	models at a practical size.	1194
1136	Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-	1196
1137	Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Kr-	1197
1138	ishnan, Marc Aurelio Ranzato, Francisco Guzmán,	1198
1139	and Angela Fan. 2022. The Flores-101 evaluation	1199
1140	benchmark for low-resource and multilingual ma-	1200
1141	chine translation. <i>Transactions of the Association for</i>	1201
1142	<i>Computational Linguistics</i> , 10:522–538.	1202
1143	Christian Herold and Hermann Ney. 2023. Improving	1203
1144	long context document-level machine translation. In	1204
1145	<i>Proceedings of the 4th Workshop on Computational</i>	1205
1146	<i>Approaches to Discourse (CODI 2023)</i> , pages 112–	1206
1147	125, Toronto, Canada. Association for Computational	1207
1148	Linguistics.	
1149	Yuchen Jiang, Tianyu Liu, Shuming Ma, Dongdong	1208
1150	Zhang, Jian Yang, Haoyang Huang, Rico Sennrich,	1209
1151	Ryan Cotterell, Mrinmaya Sachan, and Ming Zhou.	1210
1152	2022. BlonDe: An automatic evaluation metric for	1211
1153	document-level machine translation. In <i>Proceedings</i>	1212
1154	<i>of the 2022 Conference of the North American Chap-</i>	1213
1155	<i>ter of the Association for Computational Linguistics:</i>	1214
1156	<i>Human Language Technologies</i> , pages 1550–1565,	1215
1157	Seattle, United States. Association for Computational	1216
1158	Linguistics.	1217
1159	Tom Kocmi, Eleftherios Avramidis, Rachel Bawden,	1218
1160	Ondřej Bojar, Anton Dvorkovich, Christian Fed-	1219
1161	ermann, Mark Fishel, Markus Freitag, Thamme	1220
1162	Gowda, Roman Grundkiewicz, Barry Haddow,	1221
1163	Philipp Koehn, Benjamin Marie, Christof Monz,	1222
1164	Makoto Morishita, Kenton Murray, Makoto Nagata,	1223
1165	Toshiaki Nakazawa, Martin Popel, Maja Popović,	1224
1166	and Mariya Shmatova. 2023. Findings of the 2023	1225
1167	conference on machine translation (WMT23): LLMs	1226
1168	are here but not quite there yet. In <i>Proceedings of the</i>	1227
1169	<i>Eighth Conference on Machine Translation</i> , pages	1228
1170	1–42, Singapore. Association for Computational Lin-	1229
1171	guistics.	1230
1172	Tom Kocmi, Vilém Zouhar, Eleftherios Avramidis,	1231
1173	Roman Grundkiewicz, Marzena Karpinska, Maja	1232
1174	Popović, Mrinmaya Sachan, and Mariya Shmatova.	1233
1175	2024. Error span annotation: A balanced approach	1234
1176	for human evaluation of machine translation. In	1235
1177	<i>Proceedings of the Ninth Conference on Machine</i>	1236
1178	<i>Translation</i> , pages 1440–1453, Miami, Florida, USA.	
1179	Association for Computational Linguistics.	
1180	Moshe Koppel and Noam Ordan. 2011. Translationese	1237
1181	and its dialects. In <i>Proceedings of the 49th An-</i>	1238
1182	<i>annual Meeting of the Association for Computational</i>	1239
1183	<i>Linguistics: Human Language Technologies</i> , pages	1240
1184	1318–1326, Portland, Oregon, USA. Association for	1241
1185	Computational Linguistics.	1242
	Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab,	1186
	Daan van Esch, Nasanbayar Ulzii-Orshikh, Allah-	1187
	sera Tapo, Nishant Subramani, Artem Sokolov, Clay-	1188
	tone Sikasote, Monang Setyawan, Supheakmungkol	1189
	Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, An-	1190
	nette Rios, Isabel Papadimitriou, Salomey Osei, Pe-	1191
	dro Ortiz Suarez, Iro-ro Orife, Kelechi Ogueji, An-	1192
	dre Niyongabo Rubungo, Toan Q. Nguyen, Math-	1193
	ias Müller, André Müller, Shamsuddeen Hassan	1194
	Muhammad, Nanda Muhammad, Ayanda Mnyak-	1195
	eni, Jamshidbek Mirzakhlov, Tapiwanashe Matan-	1196
	gira, Colin Leong, Nze Lawson, Sneha Kudugunta,	1197
	Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaven-	1198
	ture F. P. Dossou, Sakhile Dlamini, Nisansa de Silva,	1199
	Sakine Çabuk Ballı, Stella Biderman, Alessia Bat-	1200
	tisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar,	1201
	Israel Abebe Azime, Ayodele Awokoya, Duygu Ata-	1202
	man, Orevaoghene Ahia, Oghenefego Ahia, Sweta	1203
	Agrawal, and Mofetoluwa Adeyemi. 2022. Quality	1204
	at a glance: An audit of web-crawled multilingual	1205
	datasets. <i>Transactions of the Association for Compu-</i>	1206
	<i>tational Linguistics</i> , 10:50–72.	1207
	Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier	1208
	Garcia, Christopher A. Choquette-Choo, Katherine	1209
	Lee, Derrick Xin, Aditya Kusupati, Romi Stella,	1210
	Ankur Bapna, and Orhan Firat. 2023. Madlad-400:	1211
	A multilingual and document-level large audited	1212
	dataset.	1213
	Chia-Hsuan Lee, Aditya Siddhant, Viresh Ratnakar, and	1214
	Melvin Johnson. 2022. DOCmT5: Document-level	1215
	pretraining of multilingual language models. In <i>Find-</i>	1216
	<i>ings of the Association for Computational Linguistics:</i>	1217
	<i>NAACL 2022</i> , pages 425–437, Seattle, United States.	1218
	Association for Computational Linguistics.	1219
	Jindřich Libovický and Jindřich Helcl. 2017. Attention	1220
	strategies for multi-source sequence-to-sequence	1221
	learning. In <i>Proceedings of the 55th Annual Meeting</i>	1222
	<i>of the Association for Computational Linguistics (Vol-</i>	1223
	<i>ume 2: Short Papers)</i> , pages 196–202, Vancouver,	1224
	Canada. Association for Computational Linguistics.	1225
	António Lopes, M. Amin Farajian, Rachel Bawden,	1226
	Michael Zhang, and André F. T. Martins. 2020.	1227
	Document-level neural MT: A systematic compar-	1228
	ison. In <i>Proceedings of the 22nd Annual Conference</i>	1229
	<i>of the European Association for Machine Translation</i> ,	1230
	pages 225–234, Lisboa, Portugal. European Associa-	1231
	tion for Machine Translation.	1232
	Yinquan Lu, Wenhao Zhu, Lei Li, Yu Qiao, and Fei	1233
	Yuan. 2024. Llamax: Scaling linguistic horizons of	1234
	llm by enhancing translation capabilities beyond 100	1235
	languages. <i>arXiv preprint arXiv:2407.05975.</i>	1236
	Eva Martínez Garcia, Carles Creus, and Cristina España-	1237
	Bonet. 2019. Context-aware neural machine transla-	1238
	tion decoding. In <i>Proceedings of the Fourth Work-</i>	1239
	<i>shop on Discourse in Machine Translation (DiscoMT</i>	1240
	<i>2019)</i> , pages 13–23, Hong Kong, China. Association	1241
	for Computational Linguistics.	1242

1243	Sameen Maruf, Fahimeh Saleh, and Gholamreza Haffari.	Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In <i>Proceedings of NAACL-HLT 2019: Demonstrations</i> .	1300
1244	2021. A survey on document-level neural machine translation: Methods and evaluation. <i>ACM Computing Surveys (CSUR)</i> , 54(2):1–36.		1301
1245			1302
1246			
1247	Arya D. McCarthy, Rachel Wicks, Dylan Lewis, Aaron Mueller, Winston Wu, Oliver Adams, Garrett Nicolai, Matt Post, and David Yarowsky. 2020. The Johns Hopkins University Bible corpus: 1600+ tongues for typological exploration. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 2884–2892, Marseille, France. European Language Resources Association.	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	1303
1248			1304
1249			1305
1250			1306
1251			1307
1252			1308
1253			1309
1254			
1255	Lesly Miculicich, Dhananjay Ram, Nikolaos Pappas, and James Henderson. 2018. Document-level neural machine translation with hierarchical attention networks. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2947–2954, Brussels, Belgium. Association for Computational Linguistics.	Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In <i>Proceedings of the Tenth Workshop on Statistical Machine Translation</i> , pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.	1310
1256			1311
1257			1312
1258			1313
1259			1314
1260			
1261		Matt Post. 2018. A call for clarity in reporting BLEU scores. In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 186–191, Brussels, Belgium. Association for Computational Linguistics.	1315
1262	Alireza Mohammadshahi, Vassilina Nikoulina, Alexandre Berard, Caroline Brun, James Henderson, and Laurent Besacier. 2022. SMaLL-100: Introducing shallow multilingual machine translation model for low-resource languages. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 8348–8359, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.		1316
1263			1317
1264			1318
1265			1319
1266		Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i> , 21(140):1–67.	1320
1267			1321
1268			1322
1269			1323
1270			1324
1271	Mathias Müller, Annette Rios, Elena Voita, and Rico Sennrich. 2018. A large-scale test set for the evaluation of context-aware pronoun translation in neural machine translation. In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 61–72, Brussels, Belgium. Association for Computational Linguistics.	Annette Rios Gonzales, Laura Mascarell, and Rico Sennrich. 2017. Improving word sense disambiguation in neural machine translation with sense embeddings. In <i>Proceedings of the Second Conference on Machine Translation</i> , pages 11–19, Copenhagen, Denmark. Association for Computational Linguistics.	1325
1272			1326
1273			1327
1274			1328
1275			1329
1276			1330
1277			1331
1278	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, Semaarley 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.	Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021a. WikiMatrix: Mining 135M parallel sentences in 1620 language pairs from Wikipedia. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 1351–1361, Online. Association for Computational Linguistics.	1332
1279			1333
1280			1334
1281			1335
1282			1336
1283			1337
1284			1338
1285			1339
1286		Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela Fan. 2021b. CCMatrix: Mining billions of high-quality parallel sentences on the web. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6490–6500, Online. Association for Computational Linguistics.	1340
1287			1341
1288			1342
1289			1343
1290			1344
1291			1345
1292			1346
1293	NLLB Team et al. 2024. Scaling neural machine translation to 200 languages. <i>Nature</i> , 630(8018):841.	Yirong Sun, Dawei Zhu, Yanjun Chen, Erjia Xiao, Xinghao Chen, and Xiaoyu Shen. 2025. Fine-grained and multi-dimensional metrics for document-level machine translation. In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 4: Student Research Workshop)</i> , pages 1–17, Albuquerque, USA. Association for Computational Linguistics.	1347
1294			1348
1295	OpenAI. 2024. Introducing ChatGPT. https://openai.com/index/chatgpt/ . [Accessed 01-06-2024].		1349
1296			1350
1297			1351
1298	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael		1352
1299			1353
			1354
			1355
			1356
			1357

1358	Zewei Sun, Mingxuan Wang, Hao Zhou, Chengqi Zhao,	2024a. Afrimte and africomet: Enhancing comet to	1416
1359	Shujian Huang, Jiajun Chen, and Lei Li. 2022. Re-	embrace under-resourced african languages.	1417
1360	thinking document-level neural machine translation.		
1361	In <i>Findings of the Association for Computational</i>	Longyue Wang, Zefeng Du, Wenxiang Jiao, Chenyang	1418
1362	<i>Linguistics: ACL 2022</i> , pages 3537–3548, Dublin,	Lyu, Jianhui Pang, Leyang Cui, Kaiqiang Song,	1419
1363	Ireland. Association for Computational Linguistics.	Derek Wong, Shuming Shi, and Zhaopeng Tu. 2024b.	1420
		Benchmarking and improving long-text translation	1421
1364	Jörg Tiedemann and Yves Scherrer. 2017. Neural ma-	with large language models. In <i>Findings of the As-</i>	1422
1365	chine translation with extended context. In <i>Proceed-</i>	<i>sociation for Computational Linguistics ACL 2024</i> ,	1423
1366	<i>ings of the Third Workshop on Discourse in Machine</i>	pages 7175–7187, Bangkok, Thailand and virtual	1424
1367	<i>Translation</i> , pages 82–92, Copenhagen, Denmark.	meeting. Association for Computational Linguistics.	1425
1368	Association for Computational Linguistics.		
1369	Sami Ul Haq, Sadaf Abdul Rauf, Arsalan Shaukat, and	Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang,	1426
1370	Abdullah Saeed. 2020. Document level NMT of low-	Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023.	1427
1371	resource languages with backtranslation. In <i>Proceed-</i>	Document-level machine translation with large lan-	1428
1372	<i>ings of the Fifth Conference on Machine Translation</i> ,	guage models. In <i>Proceedings of the 2023 Confer-</i>	1429
1373	pages 442–446, Online. Association for Computa-	<i>ence on Empirical Methods in Natural Language Pro-</i>	1430
1374	tional Linguistics.	<i>cessing</i> , pages 16646–16661, Singapore. Association	1431
		for Computational Linguistics.	1432
1375	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	1433
1376	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	Chaumond, Clement Delangue, Anthony Moi, Pier-	1434
1377	Kaiser, and Illia Polosukhin. 2017. Attention is all	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	1435
1378	you need. In <i>Advances in Neural Information Pro-</i>	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	1436
1379	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	1437
		Teven Le Scao, Sylvain Gugger, Mariama Drame,	1438
1380	Elena Voita, Rico Sennrich, and Ivan Titov. 2019. When	Quentin Lhoest, and Alexander Rush. 2020. Trans-	1439
1381	a good translation is wrong in context: Context-aware	formers: State-of-the-art natural language processing.	1440
1382	machine translation improves on deixis, ellipsis, and	In <i>Proceedings of the 2020 Conference on Empirical</i>	1441
1383	lexical cohesion. In <i>Proceedings of the 57th An-</i>	<i>Methods in Natural Language Processing: System</i>	1442
1384	<i>Annual Meeting of the Association for Computational</i>	<i>Demonstrations</i> , pages 38–45, Online. Association	1443
1385	<i>Linguistics</i> , pages 1198–1212, Florence, Italy. Asso-	for Computational Linguistics.	1444
1386	ciation for Computational Linguistics.		
1387	Elena Voita, Pavel Serdyukov, Rico Sennrich, and Ivan	Billy T. M. Wong and Chunyu Kit. 2012. Extending	1445
1388	Titov. 2018. Context-aware neural machine trans-	machine translation evaluation metrics with lexical	1446
1389	lation learns anaphora resolution. In <i>Proceedings of the</i>	cohesion to document level. In <i>Proceedings of the</i>	1447
1390	<i>56th Annual Meeting of the Association for</i>	<i>2012 Joint Conference on Empirical Methods in Natu-</i>	1448
1391	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	<i>ral Language Processing and Computational Natural</i>	1449
1392	pages 1264–1274, Melbourne, Australia. Association	<i>Language Learning</i> , pages 1060–1068, Jeju Island,	1450
1393	for Computational Linguistics.	Korea. Association for Computational Linguistics.	1451
1394	Jiayi Wang, David Ifeoluwa Adelani, Sweta Agrawal,	Minghao Wu, George Foster, Lizhen Qu, and Gholam-	1452
1395	Marek Masiak, Ricardo Rei, Eleftheria Briakou,	reza Haffari. 2023. Document flattening: Beyond	1453
1396	Marine Carpuat, Xuanli He, Sofia Bourhim, An-	concatenating context for document-level neural ma-	1454
1397	diswa Bukula, Muhidin Mohamed, Temitayo Ola-	chine translation. In <i>Proceedings of the 17th Confer-</i>	1455
1398	toye, Tosin Adewumi, Hamam Mokayed, Christine	<i>ence of the European Chapter of the Association for</i>	1456
1399	Mwase, Wangui Kimotho, Foutse Yuehgoh, An-	<i>Computational Linguistics</i> , Dubrovnik, Croatia.	1457
1400	uoluwapo Aremu, Jessica Ojo, Shamsuddeen Has-	Minghao Wu, Thuy-Trang Vu, Lizhen Qu, George Fos-	1458
1401	san Muhammad, Salomey Osei, Abdul-Hakeem	ter, and Gholamreza Haffari. 2024. Adapting large	1459
1402	Omotayo, Chiamaka Chukwuneke, Perez Ogayo,	language models for document-level machine trans-	1460
1403	Oumaima Hourrane, Salma El Anigri, Lolwethu	lation.	1461
1404	Ndolela, Thabiso Mangwana, Shafie Abdi Mohamed,	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale,	1462
1405	Ayinde Hassan, Oluwabusayo Olufunke Awoyomi,	Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and	1463
1406	Lama Alkhaled, Sana Al-Azzawi, Naome A. Etori,	Colin Raffel. 2021. mT5: A massively multilingual	1464
1407	Millicent Ochieng, Clemencia Siro, Samuel Njoroge,	pre-trained text-to-text transformer. In <i>Proceedings</i>	1465
1408	Eric Muchiri, Wangari Kimotho, Lyse Naomi Wamba	<i>of the 2021 Conference of the North American Chap-</i>	1466
1409	Momo, Daud Abolade, Simbiat Ajao, Iyanuoluwa	<i>ter of the Association for Computational Linguistics:</i>	1467
1410	Shode, Ricky Macharm, Ruqayya Nasir Iro, Sa-	<i>Human Language Technologies</i> , pages 483–498, On-	1468
1411	heed S. Abdullahi, Stephen E. Moore, Bernard	line. Association for Computational Linguistics.	1469
1412	Opoku, Zainab Akinjobi, Abeebe Afolabi, Nnaemeka	Fei Yuan, Yinquan Lu, Wenhao Zhu, Lingpeng Kong,	1470
1413	Obiefuna, Onyekachi Raphael Ogbu, Sam Brian,	Lei Li, Yu Qiao, and Jingjing Xu. 2023. Lego-MT:	1471
1414	Verrah Akinyi Otiende, Chinedu Emmanuel Mbonu,		
1415	Sakayo Toadoum Sari, Yao Lu, and Pontus Stenetorp.		

1472	Learning detachable models for massively multilingual machine translation.	B	More details about AFRIDOC-MT	1527
1473	In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> ,			
1474	pages 11518–11533, Toronto, Canada. Association		Table 10 shows the average number of white-space-	1528
1475	for Computational Linguistics.		separated tokens for sentences across various do-	1529
1476			domains and their corresponding translations in all the	1530
1477	Dawei Zhu, Pinzhen Chen, Miaoran Zhang, Barry Had-		languages including English. The <i>health</i> domain	1531
1478	dow, Xiaoyu Shen, and Dietrich Klakow. 2024a.		has more tokens on average than <i>tech</i> . Hausa and	1532
1479	Fine-tuning large language models to translate: Will		Yorùbá have more tokens on average than English,	1533
1480	a touch of noisy data in misaligned languages suf-		possibly because they are descriptive languages,	1534
1481	fice?		while Swahili has a comparably similar length to	1535
1482	In <i>Proceedings of the 2024 Conference on</i>		English. However, Amharic and Zulu have rela-	1536
1483	<i>Empirical Methods in Natural Language Processing</i> ,		tively shorter average lengths, demonstrating inter-	1537
1484	pages 388–409, Miami, Florida, USA. Association		esting linguistic properties.	1538
	for Computational Linguistics.			
1485	Dawei Zhu, Sony Trenous, Xiaoyu Shen, Dietrich		B.1 Translation guidelines	1539
1486	Klakow, Bill Byrne, and Eva Hasler. 2024b. A		The translation guidelines, aside the details shared	1540
1487	preference-driven paradigm for enhanced translation		at the workshop on translation and terminology	1541
1488	with large language models.		creation can be found below.	1542
1489	In <i>Proceedings of the 2024 Conference of the North American Chapter of</i>			
1490	<i>the Association for Computational Linguistics: Hu-</i>		• Thank you for agreeing to work on this project.	1543
1491	<i>man Language Technologies (Volume 1: Long Pa-</i>		Below is the link to access the data for transla-	1544
1492	<i>pers)</i> , pages 3385–3403, Mexico City, Mexico. Asso-		tion. The files are in .csv format, and you can	1545
1493	ciation for Computational Linguistics.		open them using Google Sheets or Microsoft	1546
			Excel (for offline work).	1547
1494	Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu,		• Each file contains 2500 sentences, and they	1548
1495	Shujian Huang, Lingpeng Kong, Jiajun Chen, and		are named in the format of a serial number	1549
1496	Lei Li. 2024c. Multilingual machine translation with		followed by your first name.	1550
1497	large language models: Empirical results and anal-			
1498	ysis.		• Please do not delete double empty rows, as	1551
1499	In <i>Findings of the Association for Computa-</i>		they serve to separate paragraphs. Also, avoid	1552
1500	<i>tional Linguistics: NAACL 2024</i> , pages 2765–2781,		deleting any rows, columns, or provided text.	1553
1501	Mexico City, Mexico. Association for Computational			
	Linguistics.		• Use the language field to input the transla-	1554
1502	Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-		tions. It is essential not to rely on translation	1555
1503	Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel		engines, as our quality assurance process can	1556
1504	Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid,		detect this. Depending on such tools may re-	1557
1505	Freddie Vargus, Phil Blunsom, Shayne Longpre,		sult in potential issues that you would need to	1558
1506	Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer,		address, leading to additional work on your	1559
1507	and Sara Hooker. 2024. Aya model: An instruction		part.	1560
1508	finetuned open-access multilingual language model.		• We will provide a list of extracted terminolo-	1561
1509	<i>arXiv preprint arXiv:2402.07827.</i>		gies soon so that you can harmonize how ter-	1562
			minologies are translated.	1563
1510	A More about the languages and their		• Thank you for your attention to these guide-	1564
1511	characteristics		lines. Should you have any questions, con-	1565
1512	We selected at least one language from each sub-		cerns, or suggestions, feel free to contact us	1566
1513	region of Sub-Saharan Africa: Hausa (North-West		or reach out to your language coordinator.	1567
1514	Africa), Yorùbá (West Africa), Amharic (East			
1515	Africa), Swahili (East Africa), and Zulu (South-		B.2 Quality evaluation of the translations	1568
1516	ern Africa). Each of these languages has over 20		As part of the human translation process, we con-	1569
1517	million speakers. All of them use the Latin script		ducted quality estimation to assess the transla-	1570
1518	except for Amharic, which uses the Ge’ez script.			
1519	The Latin-script languages use the Latin alphabet			
1520	with the omission of some letters and the addition			
1521	of new ones, and the use of diacritics (e.g., Yorùbá).			
1522	The languages are tonal, except for Amharic and			
1523	Swahili. Just like English, all languages follow			
1524	the subject-verb-object word order. Refer to Ade-			
1525	lani (2022) for a comprehensive overview of the			
1526	characteristics of these languages.			

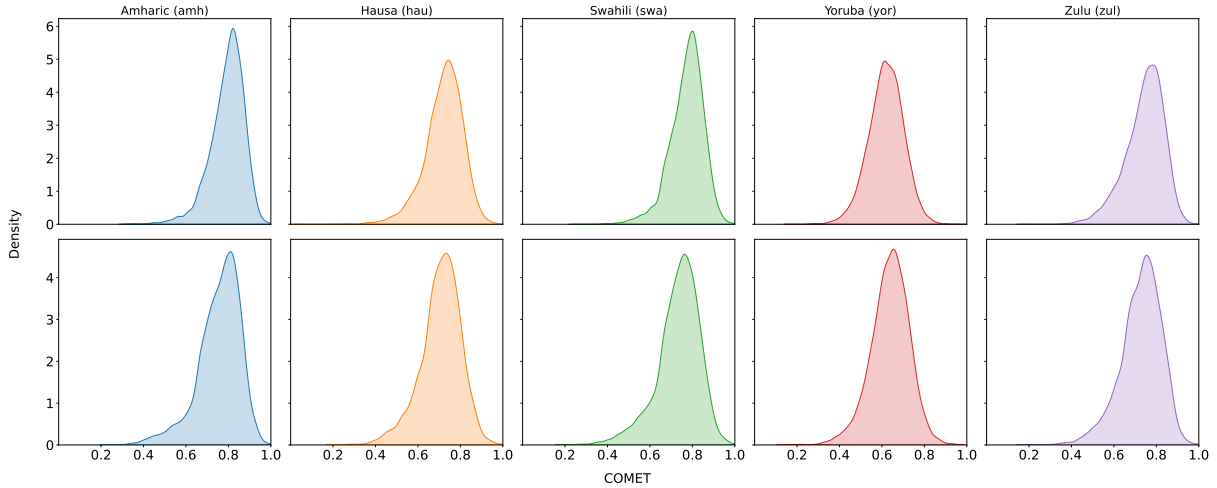


Figure 4: Distribution of the quality estimation of the translated sentences using COMET scores for the *health* (top), *tech* (bottom).

Domain	eng	amh	hau	swa	yor	zul
Sentence						
<i>health</i>	21.6	19.3	28.1	23.2	27.9	16.7
<i>tech</i>	17.8	15.6	22.2	18.0	23.7	13.4
Document						
<i>health</i>	647.3	576.7	841.7	695.4	834.8	500.1
<i>tech</i>	658.2	575.0	821.6	665.4	873.4	495.9

Table 10: The average number of tokens in AFRIDOC-MT, both at sentence and document level.

tions. For this purpose, we used AfriCOMET.¹² Given a translated sentence in any African language and its corresponding source English sentence, AfriCOMET generates a score between 0 and 1, where 0 indicates poor quality and higher values signify better quality. The translators, in collaboration with the language coordinators, were tasked with reviewing instances that had quality estimation scores below 0.65. This step was essential to identify and correct low-quality translations.

Figure 4 illustrate the distribution of the final quality scores for the five languages and both domains. Our manual check indicates that QE scores below 0.65 are not necessarily indicative of poor translations, which is consistent with the findings of Adelani et al. (2024b). We attribute this observation to factors such as domain shift, translation length, and other potential influences, which warrant further investigation in future research.

¹²<https://huggingface.co/masakhane/africomet-qe-stl>

B.3 Creation of pseudo-documents for AFRIDOC-MT

Given that the translated documents vary in length in terms of sentences and tokens, and considering the maximum token length limitations of the different LLMs used, we adopted a chunking approach for document-level evaluation. In this approach, documents were divided into smaller pseudo-documents that fit within the maximum length constraints of the models. To establish an appropriate chunk size, each document was divided into fixed-size chunks of k sentences, with the possibility that the final chunk may contain fewer than k sentences. These sentence groups, referred to as pseudo-documents, were used for document-level translation.

We conducted an initial analysis, testing different values for k (5, 10, and 25), with $k=1$ serving as our sentence-level setup. Table 11 presents the resulting number of parallel pseudo-documents, as well as the average number of tokens per pseudo-document per language for the various model tokenizers, including the 95th percentile token count. Our analysis revealed that Amharic and Yorùbá—languages with unique characteristics such as non-Latin scripts and diacritics, respectively—had the largest average token counts across the tokenizers. Additionally, the domain with the highest number of average tokens for pseudo-document varies from language to language.

To accommodate both languages in our experiments, we chose pseudo-documents with $k=10$. However, for the SFT models described in Appendix C.2, we used both $k=5$ and $k=10$.

Languages/Split	Models	Full		25 sent.		10 sent.		5 sent.	
		Health	Tech	Health	Tech	Health	Tech	Health	Tech
Sizes of data splits in AFRIDOC-MT pseudo-document									
Train		240	187	402	369	812	789	1506	1483
Dev		33	25	56	48	112	106	209	204
Test		61	59	108	106	224	227	417	418
Statistics of LLM tokens in AFRIDOC-MT pseudo-document training splits									
en	NLLB-200	923.7/2017.6	941.9/1982.1	551.5/951.7	477.4/758.8	273.0/430.9	223.2/343.6	147.2/233.8	118.8/184.9
	MADLAD-400	971.0/2095.2	991.4/2100.1	579.7/1017.1	502.4/797.8	287.0/449.3	235.0/362.0	154.7/245.0	125.0/196.9
	Aya-101	1008.2/2183.5	1020.5/2184.3	601.9/1038.0	517.2/820.2	298.0/463.4	241.9/372.6	160.7/255.0	128.7/199.0
	LLaMA3	801.4/1788.0	842.5/1798.4	478.5/833.8	427.0/664.0	236.9/372.9	199.7/304.2	127.8/203.0	106.3/166.0
	Gemma-2	802.9/1820.1	857.9/1857.6	479.3/841.0	434.8/689.6	237.3/375.0	203.4/314.0	128.0/205.0	108.2/169.0
ModernBERT	801.0/1800.7	862.7/1819.3	478.3/837.8	437.3/685.6	236.8/373.4	204.5/311.0	127.8/204.0	108.9/171.0	
am	NLLB-200	1304.4/2785.8	1376.3/2888.7	778.8/1329.9	697.5/1130.8	385.6/592.0	326.2/520.0	207.9/328.0	173.5/282.9
	MADLAD-400	1624.8/3393.6	1685.0/3487.4	970.0/1684.2	853.9/1380.4	480.2/750.0	399.4/640.2	258.9/413.8	212.5/342.9
	Aya-101	1887.4/3937.9	1934.7/4126.9	1126.8/1931.8	980.5/1598.0	557.9/855.4	458.5/722.0	300.8/477.8	244.0/390.0
	LLaMA3	6798.0/13986.2	6829.6/14750.9	4058.5/6971.8	3461.1/5584.8	2009.3/3084.4	1618.7/2560.8	1083.3/1716.0	861.2/1379.9
	Gemma-2	2817.9/5857.5	2868.4/6227.4	1682.1/2896.4	1453.2/2342.4	832.4/1267.8	679.3/1071.6	448.5/710.0	361.0/575.0
ModernBERT	7347.8/15045.1	7386.4/15952.3	4386.4/7544.1	3742.8/6035.8	2171.1/3331.2	1749.9/2760.4	1170.2/1851.0	930.6/1485.9	
ha	NLLB-200	1204.4/2713.7	1171.4/2463.0	719.0/1252.8	593.6/962.6	356.0/554.0	277.6/430.6	191.9/306.8	147.7/232.0
	MADLAD-400	1297.1/2849.4	1260.5/2643.7	774.4/1359.7	638.8/1042.0	383.4/606.4	298.8/465.6	206.7/329.0	158.9/251.0
	Aya-101	1614.9/3497.4	1535.3/3241.9	964.1/1672.3	778.0/1254.6	477.3/742.6	363.9/563.2	257.4/410.8	193.6/306.0
	LLaMA3	1916.7/4012.9	1822.6/3917.9	1144.3/1988.8	923.7/1513.6	566.6/882.4	432.1/674.6	305.5/488.8	230.0/365.9
	Gemma-2	1642.4/3568.9	1581.3/3373.4	980.6/1716.7	801.4/1297.8	485.5/757.4	374.8/584.0	261.8/417.8	199.4/317.8
ModernBERT	1998.5/4207.5	1916.8/4139.7	1193.1/2057.8	971.5/1575.8	590.8/916.9	454.4/701.0	318.6/510.8	241.8/382.9	
sw	NLLB-200	1100.8/2494.8	1094.8/2187.5	657.2/1145.9	554.8/896.4	325.4/517.0	259.5/409.6	175.4/280.0	138.1/218.0
	MADLAD-400	1177.3/2629.9	1155.3/2293.9	702.8/1227.6	585.5/938.6	348.0/547.0	273.8/436.0	187.6/297.0	145.7/231.9
	Aya-101	1345.3/2925.0	1311.0/2667.8	803.2/1390.9	664.4/1076.2	397.6/627.9	310.7/487.4	214.4/339.0	165.3/261.0
	LLaMA3	1668.1/3605.0	1619.4/3364.9	995.9/1735.4	820.7/1330.0	493.1/771.4	383.9/599.8	266.0/418.0	204.3/323.0
	Gemma-2	1413.3/3097.3	1377.1/2770.0	843.8/1467.7	697.9/1126.2	417.8/658.9	326.4/513.0	225.3/356.8	173.7/277.9
ModernBERT	1757.9/3753.4	1719.7/3594.1	1049.5/1822.8	871.6/1421.0	519.7/810.0	407.7/632.0	280.2/441.0	217.0/342.8	
yo	NLLB-200	1702.6/3854.7	1724.8/3577.1	1016.5/1857.2	874.1/1428.6	503.2/814.7	408.8/644.6	271.3/443.8	217.5/348.9
	MADLAD-400	1983.6/4470.9	1990.4/4136.7	1184.3/2137.5	1008.7/1650.2	586.3/939.4	471.7/742.2	316.1/512.0	251.0/401.9
	Aya-101	2729.2/5832.3	2659.8/5549.7	1629.4/2956.4	1347.9/2211.6	806.7/1292.4	630.4/988.0	434.9/704.0	335.4/544.0
	LLaMA3	2945.8/6322.4	2880.0/5995.5	1758.6/3203.9	1459.4/2400.4	870.5/1406.0	682.5/1077.6	469.3/677.8	363.0/585.9
	Gemma-2	2620.4/5745.5	2593.5/5406.9	1564.3/2867.7	1314.3/2143.8	774.4/1245.4	614.6/965.6	417.4/678.0	327.0/530.0
ModernBERT	3648.3/7780.9	3595.2/7600.6	2178.1/4002.0	1822.0/3020.4	1078.3/1761.4	852.1/1339.8	581.4/957.2	453.3/733.9	
zu	NLLB-200	1201.8/2513.3	1230.4/2555.7	717.5/1233.0	623.5/1016.6	355.2/554.3	291.6/461.2	191.5/300.0	155.1/250.0
	MADLAD-400	1215.2/2524.0	1230.7/2519.6	725.5/1284.8	623.7/1007.2	359.2/557.8	291.7/465.6	193.7/305.5	155.2/251.0
	Aya-101	1491.3/3012.2	1485.2/3180.8	890.3/1521.8	752.7/1213.0	440.8/688.9	352.0/554.4	237.7/372.8	187.3/298.9
	LLaMA3	1921.7/3822.6	1834.3/3933.4	1147.3/1963.9	929.7/1512.4	568.1/885.4	434.9/689.2	306.4/475.8	231.5/373.0
	Gemma-2	1787.5/3573.5	1703.0/3666.1	1067.2/1834.8	863.0/1416.2	528.3/819.4	403.6/637.6	284.9/447.8	214.8/343.9
ModernBERT	2073.7/4134.2	1965.8/4239.3	1238.1/2138.4	996.3/1625.6	613.0/956.3	466.1/737.0	330.6/515.8	248.0/399.0	

Table 11: AFRIDOC-MT Pseudo-document statistics. The number of translation instances in the data AFRIDOC-MT pseudo-document splits. average and 95th percentile (average/95 percentile) of the AFRIDOC-MT document train split tokenization statistics using the different LLM tokenizers.

C Experimental details

C.1 Evaluated Models

C.1.1 Translation Models

M2M-100 (Fan et al., 2020) is a transformer-based multilingual neural translation model from Meta, trained to translate between 100 languages, including several African languages. It has three variants of different sizes: 400M parameters, 1.2B parameters, and 12B parameters. For our experiments, we evaluated the 400M and 1.2B variants.

NLLB (NLLB Team et al., 2024) is a model similar to M2M-100, with broader coverage, trained to translate between just over 200 languages, including more than 50 African languages. It also has different sizes: 600M, 1.3B, 3.3B, and 54B parameters. For this work, we evaluated the first three variants.

MADLAD-400 (Kudugunta et al., 2023) is a multilingual translation model based on the T5 architecture (Raffel et al., 2020), covering 450 languages, including many African languages. It was trained on data collected from the Common Crawl

dataset. The dataset underwent a thorough self-audit to filter out noisy content and ensure its quality for training MT models.

Toucan (Elmadany et al., 2024; Adebara et al., 2024) is another multilingual but Afro-centric translation model based on the T5 architecture, covering 150 language pairs of African languages. It was first pre-trained on large multilingual texts covering over 500 African languages and then finetuned on translation task covering over 100 language pairs.

C.1.2 Large Language Models

Aya-101 (Üstün et al., 2024) is an instruction-tuned mT5 model (Xue et al., 2021) designed to handle both discriminative and generative multilingual tasks. With 13B parameters, it covers 100 languages and is capable of translating between a wide range of languages, including African languages.

Gemma2 (Gemma Team et al., 2024) is a decoder-only LLM trained on billions of tokens sourced from the web. The training data primarily consists of English-language text, but it also

includes code and mathematical content. While Gemma2 has an English-centric focus, it also possesses multilingual capabilities. We evaluate the base Gemma2 model with 9B parameters, as well as its instruction-tuned version.

LLama3.1 (Dubey et al., 2024) is another decoder-only LLM trained on trillions of tokens across multiple languages. It was fine-tuned using existing instruction datasets as well as synthetically generated instruction data to create its instruction-tuned version. One advantage LLama3.1 has over other models is its context window of 128K tokens, the largest among all models considered in this work, making it particularly suitable for document-based tasks such as document-level translation. We evaluate the base LLama3.1 model with 8B parameters, as well as its instruction-tuned version.

LLaMAX3 (Lu et al., 2024) is a multilingual LLM built on the LLama3 with 8B parameters as its base. It was trained on 102 languages, including several African languages, through continued pretraining. Using an English instruction dataset (Alpaca), it was further fine-tuned to create LLaMAX3-Alpaca. We evaluated both models and compared their performance across various tasks.

C.2 Supervised Finetuning

We perform supervised fine-tuning to tailor LLMs for translation tasks. To train sentence-level MT systems, we use all parallel sentences from AFRIDOC-MT to construct the training set, enabling the LLMs to translate across multiple directions and domains. Following Zhu et al. (2024b), we augment the parallel data with translation instructions, which are randomly sampled from a pre-defined set of 31 MT instructions for each training example.¹³ To train document-level MT systems, we follow the same process, but train on longer segments formed by concatenating multiple sentences. When fine-tuning, we use a learning rate of $5e^{-6}$ and an effective batch size of 64. Models are trained for only one epoch, as further training does not result in improvements and may even lead to performance degradation.

Similarly, we fine-tuned the 1.3B version of NLLB-200 for sentence and pseudo-document (with 10 sentences) translation using the Fairseq (Ott et al., 2019) codebase. We used all

¹³We use the same instruction set as described in (Zhu et al., 2024b).

Setting	X → eng	eng → X
Sentence		
sentence	512	512
Document		
5	4096	4096
10	4096	4096
25	1024	8192 (11264)
Full	2048	16384 (32768)

Table 12: The maximum number of tokens set for decoder-only LLMs when translating between English and African languages, and vice versa. Special cases for Amharic are indicated in brackets.

the training examples from 30 language directions across both domains. The model was fine-tuned for 50k steps using a learning rate of $5e^{-5}$, token batch size of 2048 and a gradient accumulation of 2. The checkpoint with the lowest validation loss was selected as the best model for evaluation.

C.3 Evaluation setup

The models were evaluated using different tools. For example, both the NLLB-200 and M2M-100 models were evaluated with the Fairseq codebase, while Toucan and MADLAD-400 were evaluated using the Hugging Face (HF) codebase. All other LLMs, including LLama3.1 (both instruction-tuned and SFT models), Gemma, and Aya-101, were evaluated using EleutherAI LM Evaluation Harness (lm-eval) tool (Biderman et al., 2024). In all cases, greedy decoding was used.

The models evaluated have different context lengths. For encoder-decoder models, M2M-100 and NLLB have a maximum sequence length of 1024 and 512 respectively. Aya-101 and MADALAD, based on the T5 architecture, do not have a pre-specified maximum sequence length, so we fixed their maximum sequence length to 1024 for all experiments involving encoder-decoder models. However, for decoder-only models, Gemma and LLaMAX3 (based on LLama3) have a maximum sequence length of 8192, while LLama3.1 has a maximum sequence length of 128K. Since all the decoder-only models were evaluated using LM Evaluation Harness, we used a similar setup for them, selecting the maximum length based on the specific needs of each model.

Table 12 shows the maximum number of generation tokens we set when translating between English and African languages. These numbers were chosen based on the statistics from Table 11. However, for Amharic, when translating pseudo-

documents with 25 sentences and full documents, there were instances exceeding the 95th percentile derived from the training statistics. Therefore, we increased the token limit specifically for Amharic.

C.4 Evaluation prompts

While the translation models we evaluated do not require prompts, MADLAD-400, requires a prefix of the form <2xx> token, which is prepended to the source sentence. Here, xx indicates the target language using its language code (e.g., “sw” for Swahili). Similarly, Toucan uses just the target language ISO-693 code as prefix, which is prepended to the source sentence (e.g., “swa” for Swahili). For other models, including Aya-101, we used three different prompts for sentence-level translation and document translation experiments. The main difference between the prompts for these tasks is the explicit mention of “text” or “document” within the prompt, as shown in Table 23. For the base models Gemma2, Llama3.1, LLaMAX3, and Aya-101, we prompted them directly using the respective prompts. However, for the instruction-tuned versions of Gemma2 and Llama3.1, we used their respective chat templates. For all Alpaca-based models, including our SFT models, we used the Alpaca template.

C.5 Evaluation metrics

We evaluate translation quality with BLEU (Papineni et al., 2002) and ChrF (Popović, 2015) using SacreBLEU¹⁴ (Post, 2018). We run significance tests using bootstrap resampling and report the 95% confidence interval for the scores, based on a sample size of 1000. We also use AfriCOMET¹⁵ (Wang et al., 2024a) to evaluate the quality of the translation outputs. We report the chrF scores of the best prompt for each model and language direction in the main paper, with all additional results provided in the Appendix D. For document-level experiments, we evaluated the LLMs using the same three prompts as in the sentence-level experiment. For evaluation, we used BLEU and chrF scores but excluded AfriCOMET due to its backbone model, AfroXLM-R-L (Alabi et al., 2022; Adelani et al., 2024a), having a context length of 512 tokens. This made it impractical to compute COMET scores for document-level outputs.

¹⁴case:mixed|eff:no| tok:13a|smooth:exp|v:2.3.1,
¹⁵<https://huggingface.co/masakhane/africomet-stl-1.1>

C.6 GPT-4o as an evaluator for machine translation

We use GPT-4o to assess the quality of translation output, as demonstrated by Sun et al. (2025), which shows a correlation with human judgment. Due to the cost of this task, we limited our evaluation to a few selected models, including Aya-101, GPT-3.5, GPT-4o, and LLaMAX3 fine-tuned on AFRIDOC-MT sentences and pseudo-documents of 10 sentences. We compared translations performed at the sentence level and pseudo-document level in terms of fluency, content errors, and cohesion errors—specifically lexical (LE) and grammatical (GE) errors—using the same definitions as Sun et al. (2025).

Below are the prompts used to evaluate documents using GPT-4o for fluency, content errors, and cohesion errors—specifically lexical (LE) and grammatical (GE) errors.

- **Fluency:** GPT-4o is prompted to rate the fluency of a document on a scale from 1 to 5, where 5 indicates high fluency and 1 represents low fluency. This evaluation is conducted without providing any reference document. For the final fluency score, we report the average rating across all documents. Below we provide the prompt used.

```
Please evaluate the fluency of the
following text in <<target>>.

-----

### **Instructions:**

- **Task**: Evaluate the fluency of
the text.

- Scoring: Provide a score from 1 to
5, where:

- **5**: The text is **highly
fluent**, with no grammatical
errors, unnatural wording, or
stiff syntax.
- **4**: The text is **mostly
fluent**, with minor errors
that do not impede
understanding.
- **3**: The text is **moderately
fluent**, with noticeable
errors that may slightly
affect comprehension.
- **2**: The text has **low
fluency**, with frequent
errors that hinder
understanding.
- **1**: The text is **not fluent
**, with severe errors that
```

Model	Setup	$eng \rightarrow X$					$X \rightarrow eng$				
		d-CHRF \uparrow	Fluency \uparrow	CE \downarrow	LE \downarrow	GE \downarrow	d-CHRF \uparrow	Fluency \uparrow	CE \downarrow	LE \downarrow	GE \downarrow
Aya-101	Sent	45.8 _{12.6}	2.5 _{0.8}	10.1 _{1.3}	4.3 _{0.8}	3.4 _{0.4}	64.9 _{4.0}	3.1 _{0.2}	19.0 _{2.5}	12.5 _{1.3}	5.7 _{0.5}
	Doc10	42.8 _{15.5}	3.0 _{0.6}	9.4 _{1.8}	2.8 _{0.5}	2.1 _{0.2}	64.2 _{5.4}	3.4 _{0.2}	14.4 _{1.0}	9.4 _{1.4}	4.4 _{0.3}
GPT-3.5	Sent	44.2 _{21.4}	2.3 _{1.4}	11.8 _{6.6}	5.5 _{3.0}	4.1 _{2.3}	57.9 _{8.9}	2.8 _{0.4}	13.1 _{3.2}	7.6 _{2.8}	4.7 _{1.7}
	Doc10	30.2 _{5.5}	2.4 _{0.4}	7.5 _{0.3}	2.7 _{0.4}	2.1 _{0.3}	28.5 _{3.0}	3.0 _{0.2}	8.4 _{0.3}	2.9 _{0.3}	2.1 _{0.2}
LLaMAX3-SFT ₁	Sent	58.8 _{10.1}	3.6 _{0.4}	11.1 _{1.2}	4.2 _{0.9}	3.0 _{0.7}	61.6 _{6.1}	3.3 _{0.4}	11.5 _{1.8}	5.8 _{1.3}	3.2 _{0.2}
	Doc10	31.8 _{2.8}	2.6 _{0.5}	8.9 _{0.6}	2.9 _{0.6}	2.2 _{0.3}	28.4 _{2.1}	3.0 _{0.3}	8.8 _{0.2}	3.2 _{0.2}	2.0 _{0.1}
LLaMAX3-SFT ₁₀	Doc10	55.0 _{11.8}	3.8 _{0.6}	10.0 _{1.4}	2.7 _{0.9}	1.9 _{0.5}	69.3 _{3.1}	4.3 _{0.2}	9.5 _{1.0}	5.2 _{0.8}	2.5 _{0.5}

Table 13: Document-level evaluation in the *tech* domain, judged by GPT-4o. Compares sentence- vs. document-level outputs on Fluency (1–5 scale), Content Errors (CE), Lexical (LE), and Grammatical Cohesion Errors (GE).

```

1858         make it difficult to
1859         understand.
1860     - **Explanation**: Support your
1861       score with specific examples to
1862       justify your evaluation.
1863
1864     -----
1865
1866     ### **Output Format**:**
1867
1868     Provide your evaluation in the
1869     following JSON format:
1870
1871     ```
1872     {
1873       "Fluency": {
1874         "Score": "<the score>",
1875         "Explanation": "<your
1876           explanation on how you made
1877           the decision>"
1878       }
1879     }
1880     ```
1881
1882     -----
1883
1884     **Text to Evaluate**:**
1885
1886     <<hypothesis>>
1887
1888     Answer:
  
```

1890 • **Accuracy:** GPT-4 is prompted to identify and
 1891 list the mistakes, such as incorrect translations,
 1892 omissions, additions, and any other errors, by
 1893 comparing the model’s output to the reference
 1894 translation. After identifying these errors, we
 1895 count all of them and compute the average
 1896 across all documents, reporting that as the
 1897 content error (CE). Below is the prompt used.

```

1898     Please evaluate the accuracy of the
1899     following translated text in <<
1900     target>> by comparing it to the
1901     provided reference text.
1902
1903     -----
1904
1905     ### **Instructions**:**
1906
1907     - **Task**: Compare the text to the
1908       reference text.
  
```

```

- Identify Mistakes: List all
  mistakes related to accuracy.
- Mistake Types:
  - **Wrong Translation**:
    Incorrect meaning or
    misinterpretation leading to
    wrong information.
  - **Omission**: Missing words,
    phrases, or information
    present in the reference
    text.
  - **Addition**: Extra words,
    phrases, or information not
    present in the reference
    text.
  - **Others**: Mistakes that are
    hard to define or categorize
    .
- **Note**: If the text expresses
  the same information as the
  reference text but uses
  different words or phrasing, it
  is not considered a mistake.
- **Provide a List**: Summarize all
  mistakes without repeating the
  exact sentences. Provide an
  empty list if there are no
  mistakes.

-----

### **Output Format**:**

Provide your evaluation in the
following JSON format:

```
{
 "Accuracy": {
 "Mistakes": [
 "<list of all mistakes in the
 text with format 'Mistake
 Types: summarize the
 mistake', provide an empty
 list if there are no
 mistakes>"
]
 }
}
```

-----
  
```

1968
1969
1970
1971
1972
1973
1974

1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034

```
**Reference Text:**  
  
<<reference>>  
  
**Text to Evaluate:**  
  
<<hypothesis>>
```

- **Cohesion:** GPT-4 is prompted to rate cohesion-related mistakes, including lexical and grammatical errors, in the model’s output, comparing it to the reference translation. We count each error individually, compute the average across the documents, and report them as lexical errors (LE) and grammatical errors (GE). Below is the prompt template we used.

```
Please evaluate the cohesion of the following translated text in <<target>> by comparing it to the provided reference text.  
  
-----  
  
### **Instructions:**  
  
- **Task:** Evaluate the cohesion of the text.  
  
- **Definition:** Cohesion refers to how different parts of a text are connected using language structures like grammar and vocabulary. It ensures that sentences flow smoothly and the text makes sense as a whole.  
  
- Identify Mistakes: List all mistakes related to cohesion.  
  
- Separate the mistakes into:  
  
  - **Lexical Cohesion Mistakes:** Issues with vocabulary usage, incorrect or missing synonyms, or overuse of certain words that disrupt the flow.  
  
  - **Grammatical Cohesion Mistakes:** Problems with pronouns, conjunctions, or grammatical structures that link sentences and clauses.  
  
- **Provide Lists:** Provide separate lists for lexical cohesion mistakes and grammatical cohesion mistakes. Provide empty lists if there are no mistakes.  
  
-----  
  
### **Output Format:**  
  
Provide your evaluation in the following JSON format:
```

```
““  
{  
  "Cohesion": {  
    "Lexical Cohesion Mistakes": [  
      "<list of all mistakes in the text one by one, provide an empty list if there are no mistakes>"  
    ],  
    "Grammatical Cohesion Mistakes": [  
      "<list of all mistakes in the text one by one, provide an empty list if there are no mistakes>"  
    ]  
  }  
}  
””  
  
-----  
  
**Reference Text:**  
  
<<reference>>  
  
**Text to Evaluate:**  
  
<<hypothesis>>
```

2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064

Fluency can only have values between 1 and 5. However, the other metrics, including CE, GE, and LE, do not have a specific range and can take on any value because they are counts. Refer to (Sun et al., 2025) for more details about these metrics.

2066
2067
2068
2069
2070

C.7 Human Evaluation Setup

2071

Beyond using GPT-4o as a judge, we also conduct human evaluation on a subset of outputs from GPT-3.5, LLaMAX3-SFT₁, and LLaMAX3-SFT₁₀ for two domains, focusing specifically on translation into five African languages due to cost constraints. Translation into English was excluded, as existing automatic metrics, including GPT-based evaluations, are already reliable for this direction.

2072
2073
2074
2075
2076
2077
2078
2079

For the human evaluation, three native speakers of the African languages—primarily translators involved in the dataset creation—were recruited. Each annotator was assigned 80 documents to evaluate, tasked with marking as many error spans as possible and rating the overall quality on a scale from 0 to 100. This annotation followed the error span annotation (ESA) (Kocmi et al., 2024) protocol as implemented within the Appraise Evaluation Framework (Federmann, 2018). To assess consistency and inter-annotator agreement, 30 of the 80 documents were shared among all three annotators. Table 14 shows statistics for 80 documents

2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092

Model	Setup	Full		Shared	
		health	tech	health	tech
GPT-3.5	Sent.	5	5	-	5
	Pseudo.	5	5	-	5
LLaMAX3-SFT ₁	Sent.	5	5	5	-
	Pseudo.	5	5	5	-
LLaMAX3-SFT ₁₀	Pseudo	5	5	5	5
Total		25	25	15	15

Table 14: The number of documents annotated by each annotator for human direct assessment.

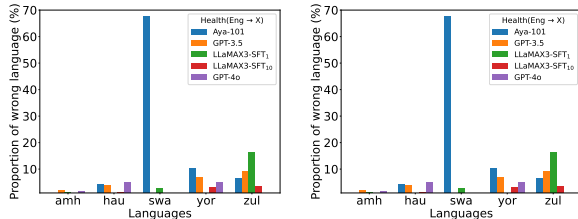


Figure 5: Rate of off-target translation ($k = 10$).

sampled from the models in both domains for each annotator. Each annotator was remunerated with \$55¹⁶.

C.8 Qualitative Analysis

Alongside the human direct assessment of the translation outputs, we shared a subset of the outputs with one author per language, each a native speaker. They were tasked with analyzing the outputs to answer two key questions: (1) What common errors or flaws do the models exhibit across different setups? and (2) How fluent are the translation outputs produced by the models across the various settings?

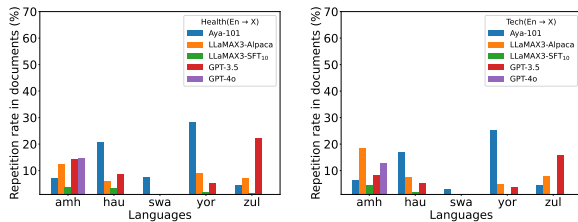


Figure 6: Word repetition rate in the pseudo-document translation ($k = 10$).

D More experimental results

D.1 Sentence-level evaluation

Given that AFRIDOC-MT is a document-level translation dataset, and due to the limited context length of most translation models and LLMs, which makes it impossible to translate a full document at once, we opted to translate the sentences within the documents and then merge them back to form the complete document. This also serves as a baseline for document-level translation. In the main

¹⁶Annotation protocol.

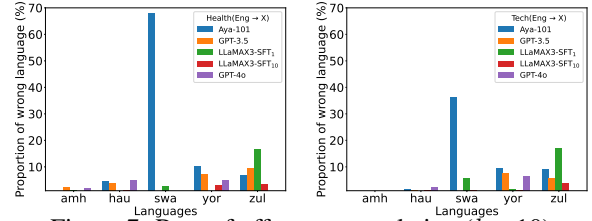


Figure 7: Rate of off-target translation ($k = 10$).

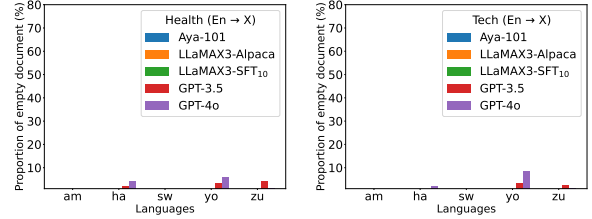


Figure 8: Proportion of empty outputs for pseudo-documents.

paper, we present the results for the best prompt for each language pair and model using d-chrF. In this section, we also provide the full results on the merged documents using d-chrF and d-BLEU in Tables 16 and 17. Furthermore, we present results for evaluating just the sentences (without merging them back into documents) using s-BLEU, s-chrF, and s-COMET in Tables 18 and 19. In Figures 18 to 21, we provide plots that summarize some of the results in the table for a few models. Although the main findings are summarized in the main draft, below are some other points we identify.

M2M-100 is not competitive Neither version of M2M-100, which was once a state-of-the-art translation model, is competitive with other translation models such as Toucan, NLLB-200, and MADLAD-400, even when compared to models of similar sizes, across all metrics at both the sentence and document levels.

Base LLMs are not translators for African languages. Base LLMs without instruction tuning and supervised fine-tuning, such as Gemma2 and LLaMAX3, do not show competitive translation performance either. This can be explained by the fact that they are just language models with limited coverage of African languages. However, LLaMAX3, which was trained on more than 100 languages, including African languages, through continued pre-training, shows improved performance, surpassing LLama3.1-IT.

Amharic and Yorùbá are the worst performing language directions. When translating from English into African languages, our results show that both Amharic and Yoruba perform the least

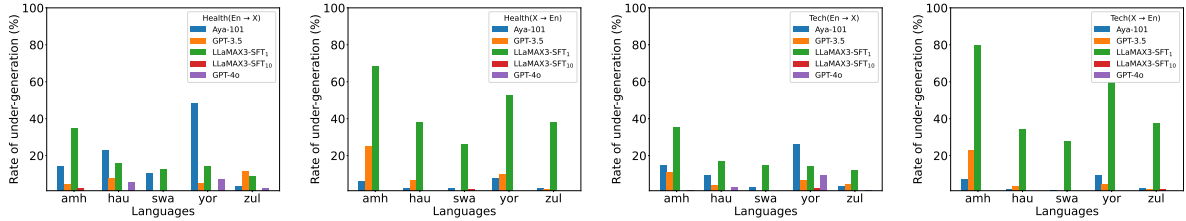


Figure 9: Rate of under-generation in pseudo-document translation ($k = 10$)

effectively. This may be attributed to specific properties of these languages, such as the use of non-Latin script in Amharic and the use of diacritics in Yoruba, which in turn increase the tokenization rate of these languages by the different model tokenizers.

D.2 Document-level evaluation

For document-level evaluation, we split the documents into chunks of 10 sentences and translate these chunks using the different models. In Tables 20 and 21 we provide the full results on the merged pseudo-documents using d-chrF and d-BLEU. And below are some other relevant points from the results. It is important to note that we also trained and evaluated NLLB-200 for pseudo-document translation. However, due to its 512-token maximum sequence length, it is not competitive. Nevertheless, the results still show the influence of fine-tuning. Below are other findings.

Gemma2-IT shows better translation performance. Compared to the sentence-level setup, where Gemma2-IT and LLaMAX3-Alpaca achieved similar performance on average, in the pseudo-document setup, Gemma2-IT not only outperforms LLaMAX3-Alpaca but also surpasses GPT-3.5. Although we cannot provide an exact explanation for this performance, we hypothesize that its pre-training setup might be a contributing factor.

Fine-tuning data has an impact on translation quality. Our results show that both LLaMA3.1 and LLaMAX3 models, when fine-tuned on sentences, performed significantly worse on pseudo-document evaluations compared to the same models fine-tuned on pseudo-documents for both domains. All these models were trained using a similar setup, with the primary difference being the data used for fine-tuning.

Language-specific performance trends Overall, no clear trend is observed in MT performance across language family classes. However, Amharic

Model	Setup	amh	hau	swh	yor
GPT-3.5	Sent	18.3	42.1	63.8	12.6
	Doc10	4.8	32.3	64.2	6.9
LLaMAX3-SFT ₁	Sent	58.1	86.0	61.0	66.1
	Doc10	19.0	54.5	36.8	40.1
LLaMAX3-SFT ₁₀	Sent	54.3	83.7	61.8	62.3
	Doc10	54.3	83.7	61.8	62.3

Table 15: Average DA score (scale 0–100) from three human evaluators per language in the *tech* domain.

(a non-Latin script language) and Yorùbá (a heavily diacriticized language) result in the lowest chrF scores, while Swahili—the most widely spoken indigenous African language—performs best.

D.3 Findings from GPT-4o as a judge

In Tables 8 and 13 we present the average GPT-4o evaluation results for four models. When translating into African languages, there is no clear pattern: for example, GPT-3.5, despite having the lowest fluency score, also had the fewest content, lexical, and grammatical errors, which is counterintuitive. In contrast, when translating into English, the pattern is clear and consistent: translations of pseudo-documents show better fluency and fewer errors overall. These results suggest that using GPT-4o as a translation judge is not yet well-suited for low-resource languages.

D.4 Findings from human evaluation

We were able to obtain DA scores from three annotators for all languages except Zulu. For each language, we calculated inter-annotator agreement using Krippendorff’s alpha α over 30 document instances. We obtained α scores of 0.46, 0.57, 0.48, and 0.81 for Amharic, Hausa, Swahili, and Yorùbá, respectively. These are relatively low scores, except for Yorùbá. We present the average DA scores in Tables 9 and 15 for the *health* and *tech* domains, respectively. The results show that annotators rate documents translated at the sentence level as higher quality than those translated at the pseudo-document level. Additionally, GPT-3.5 received the lowest ratings among the three models. LLaMAX3-SFT₁, a model trained on sentence-level data, was rated the best across all languages when evaluated

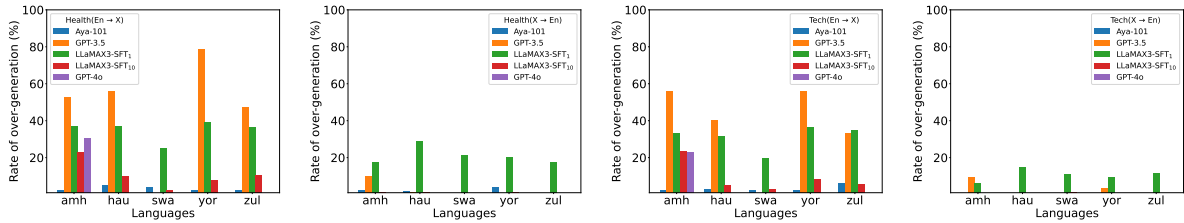


Figure 10: Rate of over-generation in pseudo-document translation ($k = 10$)

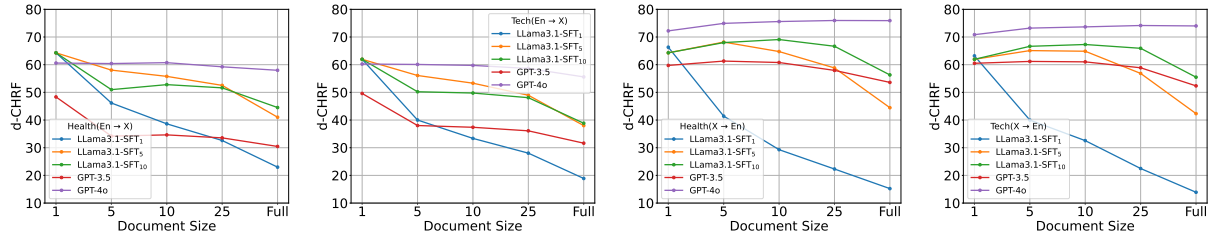


Figure 11: Average chrF score across languages for documents of different sizes.

on sentences. However, when evaluated on pseudo-documents, its performance was rated lower than that of LLaMAX3-SFT₁₀. These findings are consistent with the d-chrF scores for the models, but they do not align with the evaluations from GPT-4o as a judge.

E More discussion and analysis

Do LLMs generate translations in the correct target languages? We evaluate whether these models understand the task by generating outputs in the target languages using the OpenLID (Burchell et al., 2023) language identification model. Our results show that these models rarely generate outputs in the wrong language when translating to English. However, when translating to African languages, there is a higher likelihood of incorrect language translations, particularly with open models (Figure 7). These off-target languages often include English, and other languages including other African languages.

What is the effect of document length on translation quality? We compare the average d-chrF scores obtained by selected models, including GPT-3.5/4 and LLaMA3.1-SFT_k where $k = \{1, 5, 10\}$. The evaluation was conducted across all pseudo-document lengths: 1, 5, 10, 25, and the full length. Figure 3 shows that for translations into African languages, d-chrF scores decrease as document length increases. A similar trend is observed for the reverse translation, except for GPT-4o, which shows an increasing trend.

What language benefits more from supervised finetuning? We focus on the sentence-level task

and translated across all 30 directions for which the model was trained, evaluating both NLLB-200 (1.3B) and its fine-tuned version using d-chrF. Figures 15 and 16 show performance improvements after supervised fine-tuning of NLLB-200 for both domains. The results shows that translating into Yorùbá, which is the direction with the lowest d-chrF score from English among all the languages, benefited the most. One major factor contributing to this is the presence of diacritics. Furthermore, looking at their actual performances and not just the differences, our results show that translations into Swahili and English—both relatively high-resource languages—yield higher BLEU and CHRf scores (see Figures 13 and 14), even after supervised fine-tuning. Hence, there is much to be done to improve translation performance between low-resource language pairs.

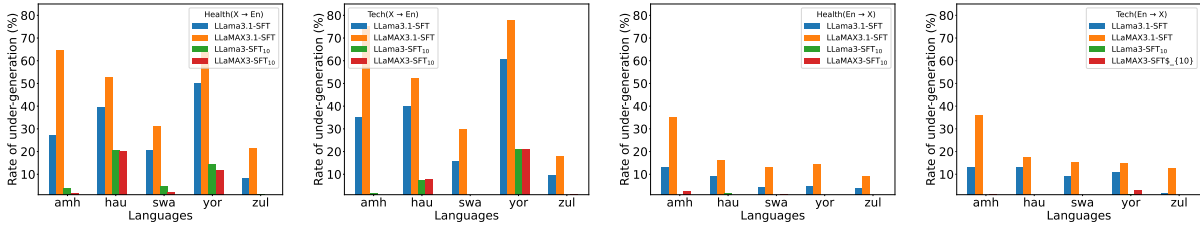


Figure 12: Rate of under-generation in our SFT models.

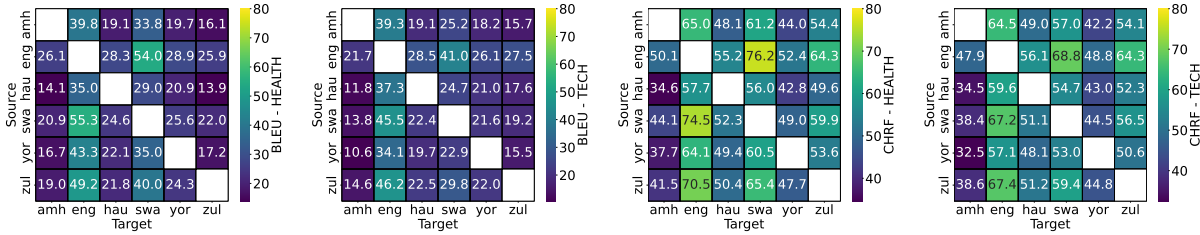


Figure 13: s-BLEU and s-chrF pair-wise comparison of supervised finetuning of NLLB-1.3B on AFRIDOC-MT

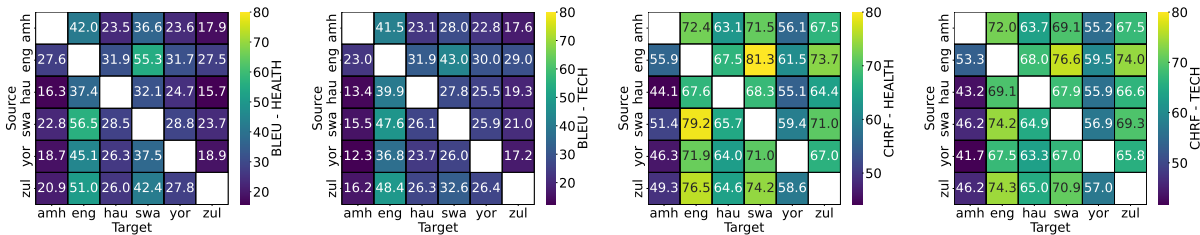


Figure 14: d-BLEU and d-chrF pair-wise comparison of supervised finetuning of NLLB-1.3B on AFRIDOC-MT

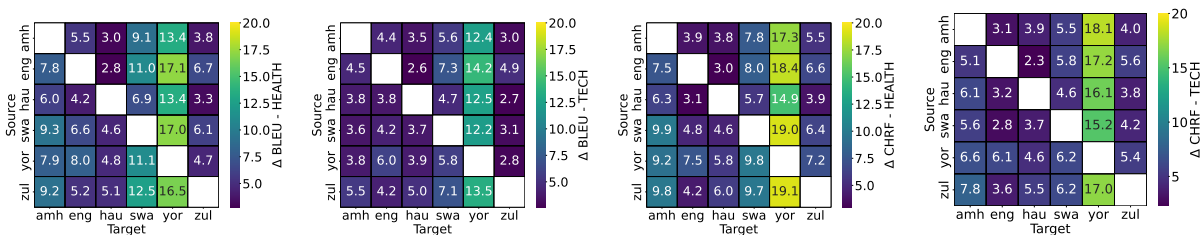


Figure 15: Change (Δ) in s-BLEU and s-chrF for sentence evaluation comparing NLLB1.3B before and after supervised finetuning on AFRIDOC-MT

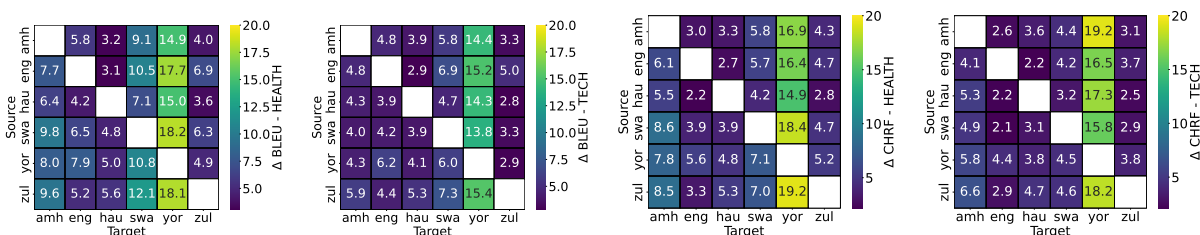


Figure 16: Change (Δ) in d-BLEU and d-chrF for sentence evaluation comparing NLLB1.3B before and after supervised finetuning on AFRIDOC-MT

Model	Setup	$eng \rightarrow X$					$X \rightarrow eng$				
		d-CHRF \uparrow	Fluency \uparrow	CE \downarrow	LE \downarrow	GE \downarrow	d-CHRF \uparrow	Fluency \uparrow	CE \downarrow	LE \downarrow	GE \downarrow
Amharic											
Aya-101	Sent	36.6	2.8	9.0	3.8	2.5	64.6	3.0	15.8	9.6	5.1
	Doc10	28.7	2.5	8.9	2.9	2.1	61.6	3.6	12.6	8.5	3.8
GPT-3.5	Sent	20.4	1.1	7.2	10.8	3.1	48.3	3.0	9.5	3.7	2.4
	Doc10	11.6	1.0	3.3	1.8	1.6	41.6	4.2	6.6	2.0	1.3
LLaMAX3-SFT ₁	Sent	46.8	3.5	10.0	2.9	2.0	66.6	3.5	11.5	6.0	2.9
	Doc10	24.1	1.7	10.4	1.9	1.8	22.6	3.0	9.0	2.5	1.6
LLaMAX3-SFT ₁₀	Doc10	37.8	2.6	7.3	1.8	1.6	64.0	4.1	8.7	4.6	2.7
Hausa											
Aya-101	Sent	56.4	2.9	9.1	4.1	3.7	61.5	2.9	17.1	10.2	5.6
	Doc10	48.5	2.9	8.3	2.3	1.9	62.3	3.2	14.6	8.4	4.5
GPT-3.5	Sent	44.3	1.4	11.3	5.2	5.8	52.4	2.5	13.2	6.8	3.7
	Doc10	23.1	1.1	7.4	2.9	2.8	52.7	4.1	9.2	4.1	2.3
LLaMAX3-SFT ₁	Sent	62.5	3.2	9.5	4.0	3.5	58.9	2.9	9.6	4.1	3.0
	Doc10	29.0	2.3	8.8	2.3	1.9	22.9	2.6	9.0	3.0	1.9
LLaMAX3-SFT ₁₀	Doc10	51.9	3.9	14.6	2.2	1.9	66.7	4.2	9.6	4.5	2.6
Swahili											
Aya-101	Sent	44.7	1.3	9.5	5.0	3.4	70.8	3.3	17.5	9.9	4.8
	Doc10	34.7	1.8	8.2	5.0	3.0	71.2	3.8	14.1	9.8	4.0
GPT-3.5	Sent	76.7	4.9	4.2	1.8	1.2	75.0	3.6	12.4	7.6	3.7
	Doc10	76.1	4.9	3.9	0.8	0.6	77.7	4.7	7.1	4.5	2.0
LLaMAX3-SFT ₁	Sent	73.1	3.6	12.3	4.8	3.4	73.1	3.9	11.8	8.4	2.7
	Doc10	42.2	3.2	9.5	3.6	2.5	33.1	3.1	8.9	3.2	1.9
LLaMAX3-SFT ₁₀	Doc10	74.4	4.4	10.5	3.6	2.2	77.8	4.6	8.9	6.2	2.7
Yoruba											
Aya-101	Sent	31.2	1.2	8.2	5.4	3.1	57.9	2.3	23.8	13.6	6.8
	Doc10	18.7	1.4	6.7	2.7	2.2	56.1	2.9	18.7	10.1	5.1
GPT-3.5	Sent	21.3	1.1	7.8	5.5	3.9	52.1	2.6	12.9	5.0	3.3
	Doc10	10.1	1.0	4.1	2.3	2.1	51.7	3.9	9.0	3.6	2.1
LLaMAX3-SFT ₁	Sent	57.5	4.0	12.3	3.8	2.7	64.7	3.2	12.2	6.1	3.1
	Doc10	33.8	2.7	8.1	2.4	1.8	27.2	3.2	9.1	3.4	2.0
LLaMAX3-SFT ₁₀	Doc10	52.2	4.2	10.3	2.4	1.7	71.8	4.4	10.2	5.9	2.4
Zulu											
Aya-101	Sent	58.6	2.7	13.7	6.0	3.6	67.4	2.9	19.1	12.2	8.6
	Doc10	54.9	3.1	16.1	3.9	2.7	69.0	3.3	15.7	10.7	4.5
GPT-3.5	Sent	51.1	1.5	20.9	10.0	6.1	59.5	2.6	15.9	8.8	5.7
	Doc10	29.2	1.3	10.1	4.0	3.2	61.1	4.0	10.8	5.3	2.9
LLaMAX3-SFT ₁	Sent	67.5	3.3	12.4	5.5	3.6	70.5	3.6	12.5	7.2	3.8
	Doc10	33.7	2.4	8.4	2.7	2.2	31.5	3.1	8.6	3.4	2.1
LLaMAX3-SFT ₁₀	Doc10	55.0	3.6	12.1	3.0	2.0	74.1	4.4	9.4	5.5	2.5

Table 22: Document-level evaluation in the Health domain, judged by GPT-4o. Compares sentence- vs. document-level outputs on Fluency (1–5 scale), Content Errors (CE), Lexical (LE), and Grammatical Cohesion Errors (GE). Best scores in bold.

Prompt 1

{system_prompt}

Translate the following {source_language} text to {target_language}:

Provide only the translation.

{source_language} text: {{source_sentence}}

{target_language} text:

Prompt 2

{system_prompt}

Translate the following {domain} text from {source_language} to {target_language}:

Provide only the translation.

{source_language} document: {{source_document}}

{target_language} document:

Prompt 3

{system_prompt}

Please provide the {target_language} translation for the following {source_language} text:{{source_document}}

Provide only the translation.

Prompt 1

{system_prompt}

Translate the following {source_language} document to {target_language}:

Provide only the translation.

{source_language} document: {{source_document}}

{target_language} document:

Prompt 2

{system_prompt}

Translate the following {domain} document from {source_language} to {target_language}:

Provide only the translation.

{source_language} document: {{source_document}}

{target_language} document:

Prompt 3

{system_prompt}

Please provide the {target_language} translation for the following {source_language} document:{{source_document}}

Provide only the translation.

Table 23: The task prompts used for evaluating LLMs are applied to both sentence-level and document-level translation tasks.

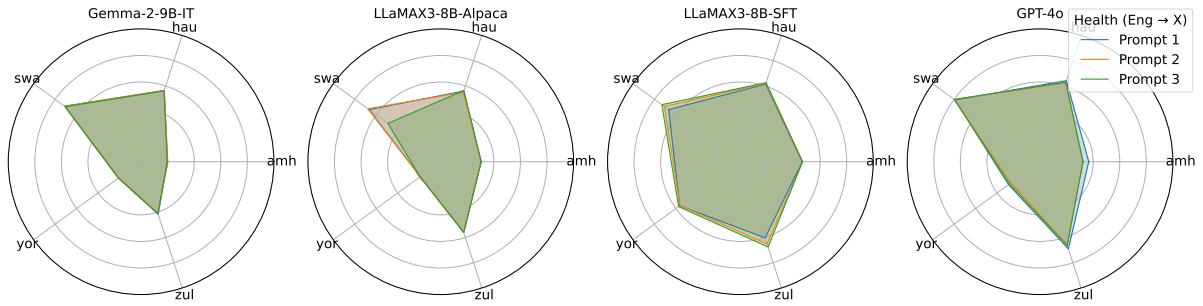


Figure 17: d-chrF scores for some LLMs for sentence-level translation using different prompts when translating into African languages

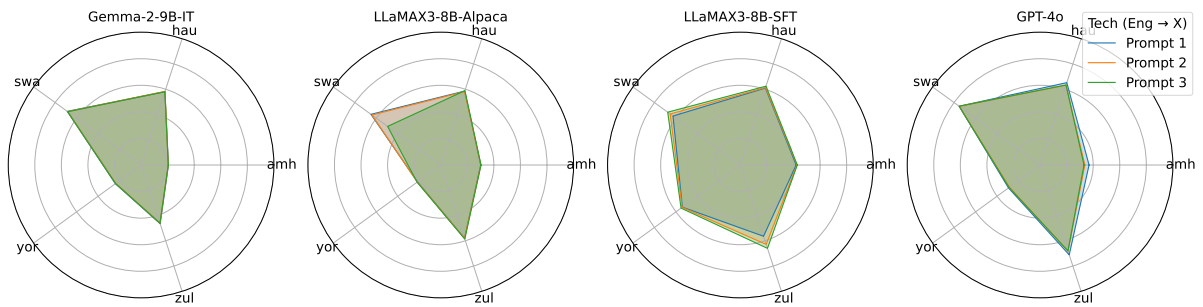


Figure 18: d-chrF scores for some LLMs for sentence-level translation using different prompts when translating into African languages

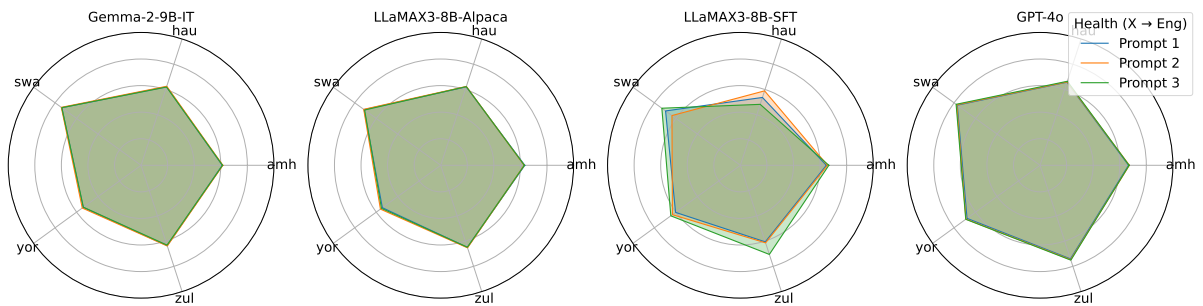


Figure 19: d-chrF scores for some LLMs for sentence-level translation using different prompts when translating into English

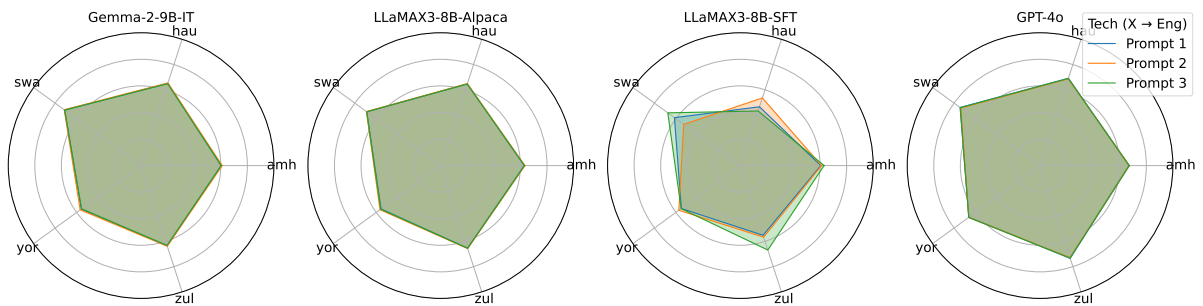


Figure 20: d-chrF scores for some LLMs for sentence-level translation using different prompts when translating into English

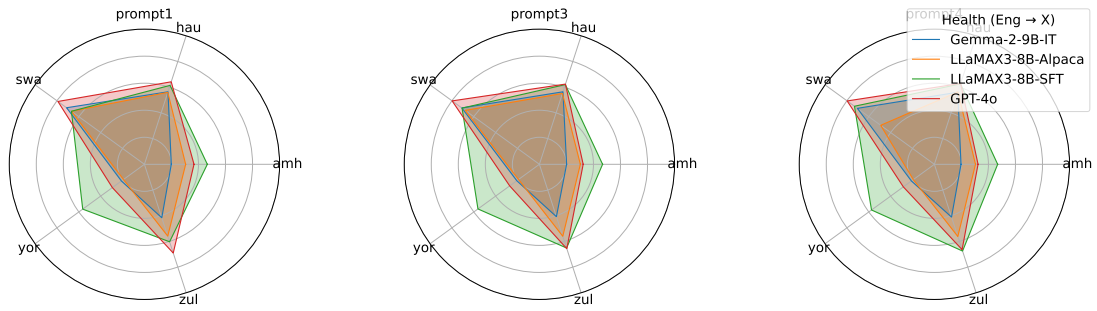


Figure 21: d-chrF scores for some LLMs for sentence-level translation using different prompts when translating into African languages

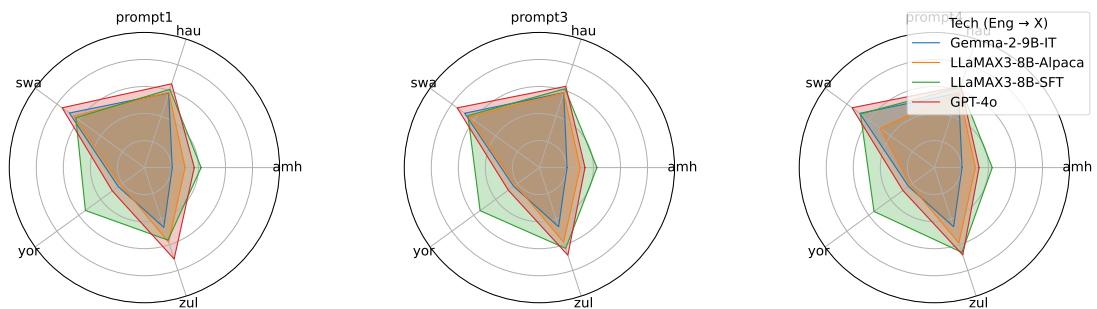


Figure 22: d-chrF scores for some LLMs for sentence-level translation using different prompts when translating into African languages

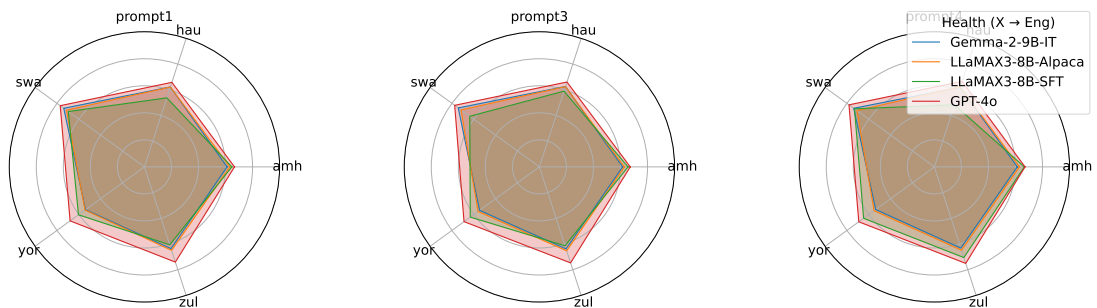


Figure 23: d-chrF scores for some LLMs for sentence-level translation using different prompts when translating into English

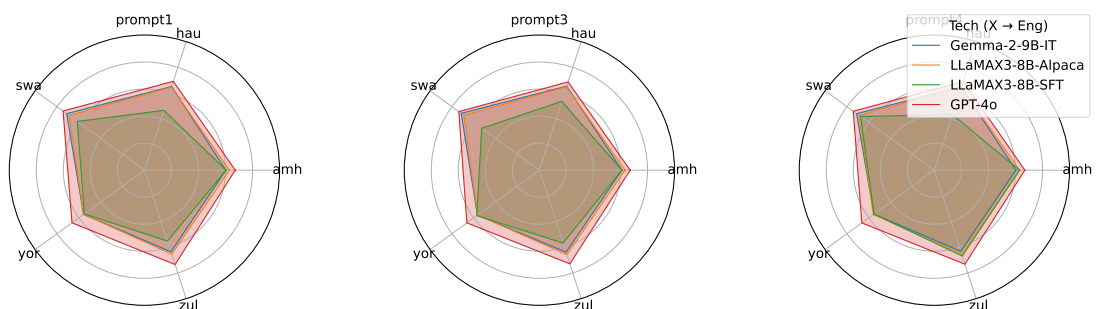


Figure 24: d-chrF scores for some LLMs for sentence-level translation using different prompts when translating into English