

LLM-Assisted Auditing of Human Translations for Low-resource Languages

Anonymous ACL submission

Abstract

Large-scale translation projects for low-resource languages mostly rely on human translators to ensure cultural and linguistic fidelity. However, even professionally produced translations often contain subtle translation errors that are difficult to detect. Manual quality control at scale becomes prohibitively expensive, creating a major bottleneck in the development of high-quality Natural Language Processing (NLP) resources. Recent advances in multilingual large language models (LLMs) offer promising support for annotation workflows with human-in-the-loop settings. In this work, we investigate the use of LLMs to assist in auditing translation quality, enabling more efficient quality control pipelines for low-resource African languages. We audit translations in 11 African languages using the MAFAND-MT dataset, combining LLM-as-a-judge, native-speaker human review, and automated metrics. Our quality-audited version of MAFAND-MT test set yields performance gains across all languages, with BLEU scores ranging from 0.4 to 9.27 and chrF scores ranging from 0.3 to 8.69. Our findings further indicate that state-of-the-art LLMs, such as GPT-5.1, can assist in auditing translation quality and suggesting candidate corrections for low-resource languages. However, they remain far from being a stand-alone solution for the automatic correction of human translations in African languages.

1 Introduction

Machine Translation (MT) is a fundamental and prominent task in natural language processing (NLP), essential for global communication and information access (Anastasopoulos et al., 2020). For many low-resource languages, particularly those in Africa, a common method for developing benchmark datasets is through human translation of existing resources from higher-resource languages such as English and French (Adelani et al., 2025a). Therefore, the quality of these translated datasets

is very crucial, as it directly impacts the evaluation and development of MT systems, ultimately determining their reliability for real-world use. High-quality human translations should satisfy at least three key criteria: fluency in the target language, adequacy in preserving the semantic content of the source text, and the target language’s cultural context (Freitag, Markus and Foster, George and Grangier, David and Ratnakar, Viresh and Tan, Qijun and Macherey, Wolfgang, 2021).

However, human translation, while indispensable for cultural and nuanced understanding, is not immune to error (Han et al., 2021; Lin et al., 2022). Translators may introduce typos, grammatical mistakes, fluency issues, and bilingual (code mixing) errors (Lin et al., 2022). These errors can stem from various factors, including the use of imperfect auxiliary translation tools, errors by native translators, the translator’s proficiency in the target language, and the ambiguity of the source content to be translated (Han et al., 2021; Lin et al., 2022). Table 1 shows examples from the MAFAND-MT test set where human-translated text in Amharic, Hausa, Igbo, Swahili, and Twi languages contains such errors, creating a "garbage in, garbage out" risk for MT models and evaluations (Adelani et al., 2022). Furthermore, errors that propagate into benchmark datasets systematically bias evaluation and hinder the development of robust MT models (Koehn and Knowles, 2017).

Despite their importance, ensuring the quality of human translations at scale remains a major challenge. Exhaustive manual review by professional translators is financially unsustainable (Sambasivan et al., 2021). Consequently, many projects face a difficult trade-off between scale, cost, and quality, potentially allowing errors to propagate into valuable resources.

Recent advances in multilingual Large Language Models (LLMs) offer a promising path toward multidisciplinary problem-solving capabilities (Treviso

MAFAND translated dataset	Corrected translation
"eng": "Date: Thursday, July 31, 2014", "amh": "የዛን ዘጠኝ ጥግረዎች ፍትህ ይገባቸዋል"	"eng": "Date: Thursday, July 31, 2014", "amh": ቀን፡ ሐሙስ፣ ሐምሌ 31፣ 2014
"eng": "Ian Simbota represented the Association of Persons with Albinism at the court.", "kin": "Igihe umucamanza yasomaga, icyampangayikishije ni uko igice cya Padiri Muhosha [n'abandi] bahamwe n'icyaha ..."	"eng": "Ian Simbota represented the Association of Persons with Albinism at the court.", "kin": "Ian Simbota yari ahagarariye ishyirahamwe ry'abantu bafite ubumuga bw'uruhu mu rukiko."
"eng": "Every soul shall have a taste of death.", "hau": "Ko shakka babu akwai wani lokaci da rayuwa za ta zo karshe Kowane rai mai dandanar mutuwa ne(Suratu Al Imrana 3:185)."	"eng": "Every soul shall have a taste of death.", "hau": "Ubangiji, muna d~aukin jiran wannan damar"
"eng": "policemen has claimed ownership of Dino melaye", "ibo": "Ndị uweojii egbochiela ụlọ Dino Melaye"	"eng": "policemen has claimed ownership of Dino melaye", "ibo": "Ndị uweojii ekwuola na ha nwe Dino Melaye"
"eng": "Ian Simbota represented the Association of Persons with Albinism at the court.", "swa": "Igihe umucamanza yasomaga, icyampangayikishije ni uko igice cya Padiri Muhosha [n'abandi] bahamwe n'icyaha"	"eng": "Ian Simbota represented the Association of Persons with Albinism at the court.", "swa": "Ian Simbota yari ahagarariye ishyirahamwe ry'abantu bafite ubumuga bw'uruhu mu rukiko."
"eng": "A local councilor, Jabu Zondo, visited the area yesterday to reprimand the incident.", "twi": "Kurow no mu kaunsila, eabu Zondo koo beaee ho nnera kohwee dee esii no"	"eng": "A local councilor, Jabu Zondo, visited the area yesterday to reprimand the incident.", "twi": "Kurow no mu gyanatufoo panyin, Jabu Zondo, koo beaee no nnera ko kaa won anim."

Table 1: **Examples of translation errors from the MAFAND-MT test dataset.** The Table shows example cases in Amharic (amh) Hausa (hau), Igbo (ibo), Swahili (swa), and Twi (twi) translations (Adelani et al., 2022). The red marked text is the incorrect translation of the given English-sourced (eng) text, and the blue is the correct translation using native speakers of the target language.

et al., 2024; Feng et al., 2025a). In this work, we investigate whether LLMs can act as assist *first-pass filters*, automatically identifying translation errors with a higher likelihood of containing errors and thus proceeding for expert review. Specifically, we explore the following three research questions: **RQ1** Can large language models (LLMs) assist in detecting and correcting human translation errors in low-resource African languages (how are they good enough to judge the translation quality of low-resource languages)? **RQ2** What types of translation errors are commonly found in machine translation (MT) resources for African languages? and **RQ3** How does translation quality review improve the performance of machine translation systems in low-resource languages?

Our contributions are threefold: (1) We introduce a pipeline for LLM-assisted quality assurance of translated resources; (2) We provide a detailed error analysis of a subset of human-translated MAFAND-MT datasets (11 languages), (3) We explore the different kinds errors exist in African MT resources; and (4) We offer insights on using LLMs as a cost-effective alternative in the reviewing of a high-quality human translation dataset for low-resource languages. Our findings demonstrate that

an LLM-human workflow can help develop reliable and accurate MT datasets and systems.

2 Related Work

2.1 Auditing Training Corpora Quality

Recent efforts to improve NLP for African languages have increasingly emphasized both the scale and quality of training corpora. Early work focused on constructing large-scale web-crawled multilingual pretraining datasets (Xue et al., 2021; Vegi et al., 2022b; Tonja et al., 2024), demonstrating the feasibility of incorporating a broader set of African languages into foundation models. However, the reliance on web-crawled data introduced substantial noise, including mistranslations, misalignments, and non-parallel content, which introduces greater degradation in data quality for low-resource languages.

To mitigate these quality issues, subsequent studies have focused on improving the quality of translation datasets through filtering, cleaning strategies, and manual or semi-automatic audits (Zhang et al., 2020; Kreutzer et al., 2022). This evolution reflects a growing recognition that data quality, rather than volume, is a critical bottleneck for machine translation performance in low-resource settings.

136 Despite these advances, existing auditing ap- 186
137 proaches remain largely human-intensive, limiting 187
138 their scalability across languages and domains. In 188
139 this work, we position LLMs as a complementary 189
140 tool for MT data auditing, examining their ability to 190
141 judge if the given translation is whether correct or 191
142 not, identify translation errors and suggests correct 192
143 translations. 193

144 2.2 Evaluations of Test Dataset Quality 194

145 Language models are commonly evaluated on 196
146 downstream tasks, with Machine Translation (MT) 197
147 serving as a central benchmark for assessing 198
148 cross-lingual capabilities. For African languages, 199
149 MAFAND-MT (Adelani et al., 2022) is one widely 200
150 used evaluation dataset that has supported numer- 201
151 ous studies on MT training and evaluations (Ojo 202
152 et al., 2025; Abdulmumin et al., 2022; Vegi et al., 203
153 2022a; Nzeyimana, 2024; Tang et al., 2024; Ji et al., 204
154 2025; Singh et al., 2025). The validity of con- 205
155 clusions drawn from such benchmarks critically 206
156 depends on the quality of their underlying transla- 207
157 tions. 208

158 The work by Abdulmumin et al. (2024) demon- 209
159 strated that even human-translated evaluation 210
160 datasets are susceptible to translation errors and 211
161 identified and corrected issues for some African 212
162 languages (Hausa, Sepedi, Xitsonga, isiZulu) in 213
163 the FLORES dataset. These findings, together with 214
164 prior analyses (Freitag, Markus and Foster, George 215
165 and Grangier, David and Ratnakar, Viresh and Tan, 216
166 Qijun and Macherey, Wolfgang, 2021), highlight 217
167 the need for systematic auditing and validation of 218
168 MT evaluation datasets to ensure reliable bench- 219
169 marking. 220

170 However, existing approaches primarily rely on 221
171 additional rounds of human translation and expert 222
172 review, which are costly and difficult to scale across 223
173 languages and datasets. In contrast, the use of 224
174 LLMs for auditing evaluation data quality remains 225
175 underexplored. In this work, we investigate the ex- 226
176 tent to which LLMs can assist in auditing MT eval- 227
177 uation datasets by identifying translation errors and 228
178 inconsistencies, and we analyze their agreement 229
179 with human judgments to better understand when 230
180 LLM-based auditing can reduce cost and when hu- 231
181 man oversight remains essential. 232

182 2.3 LLM-as-a-Judge Translation Review 232

183 LLMs have demonstrated strong performance as 233
184 an evaluator in various tasks, including automated 234
185 data quality control (Gu et al., 2025), dataset an-

notation assistance (Tan et al., 2024; Belay et al., 186
2025), identifying error types for machine transla- 187
tion dataset (Feng et al., 2025b; Kim, 2025), and 188
research paper summarization (Liu et al., 2024). 189
Within machine translation research, LLM-assisted 190
translation error detection and correction have been 191
explored primarily in English as an automatic post- 192
editing (APE) (Berger et al., 2024; Freitag, Markus 193
and Foster, George and Grangier, David and Rat- 194
nakar, Viresh and Tan, Qijun and Macherey, Wolf- 195
gang, 2021; Lu et al., 2024), translation quality 196
evaluation (Qian et al., 2024), and automatic cor- 197
rection of human translations (Lin et al., 2022). 198

199 Despite these advances, the use of LLMs as trans- 200
lation reviewers, capable of identifying, categor- 201
izing, and correcting translation errors, remains 202
unexplored, mainly for low-resource and African 203
languages. Recent evaluations of MAFAND-MT 204
(Adelani et al., 2022) reveal the presence of trans- 205
lation errors and varying degrees of semantic mis- 206
alignment, as evidenced by low automatic quality 207
scores reported in prior work using metrics such 208
as COMET (Falcão et al., 2024). These under- 209
score the need for scalable, systematic translation- 210
review methods that extend beyond manual in- 211
spection. In this work, we investigate the role 212
of LLMs as translation reviewers for African- 213
language datasets. Specifically, we assess their 214
ability to judge whether the translation is correct, 215
identify common types of translation errors, pro- 216
pose correction candidates, and support a human- 217
in-the-loop auditing pipeline. By analysing agree- 218
ment between LLM-based judgments and human 219
verification, we aim to clarify both the potential 220
and the limitations of LLMs as assistants for transla- 221
tion data quality auditing. 222

222 3 Translated African Languages Dataset 222

223 A growing number of human-translated datasets are 224
225 available for African languages. These datasets are 226
227 mostly translated from English and French source 228
229 texts and span diverse domains. For machine trans- 230
231 lation NLP tasks, prominent datasets that include 232
233 African languages are FLORES (Guzmán et al., 2019), NLLB (Team et al., 2022), HornMT¹, and MAFAND-MT (Adelani et al., 2022). Recently, domain-specific datasets have also been created, such as AFRIDOC-MT (Alabi et al., 2025) and AfriMed-QA (Nimo et al., 2025) for health-related

¹A multi-way parallel news corpus for languages in the Horn of Africa; <https://github.com/asmeLashteka/HornMT>

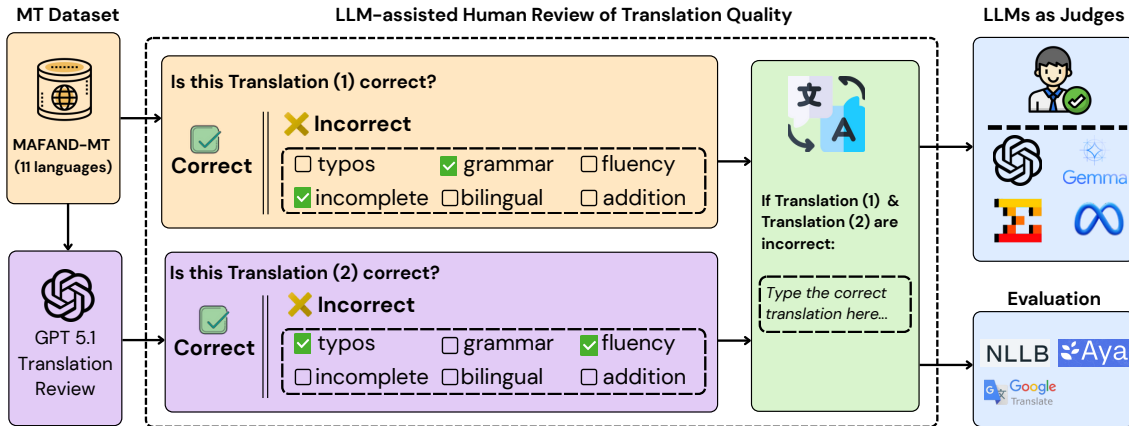


Figure 1: **LLM-assisted pipeline for MT dataset quality assessment and correction.** The figure illustrates the workflow for judging whether the translation is correct or not, identifying translation errors, and generating correction suggestions.

234 data, AfriGSM (Adelani et al., 2025b) for math
 235 word problems, and AfriMMLU and AfriXNLI
 236 (Adelani et al., 2025b) for general knowledge and
 237 reasoning.

238 **The MAFAND-MT Translation Dataset** The
 239 Masakhane Anglo and Franco Africa News Dataset
 240 for Machine Translation (MAFAND-MT) is one
 241 of the frequently used MT evaluation datasets for
 242 African languages (Adelani et al., 2022). This
 243 dataset covers 20 African languages: 15 are trans-
 244 lated from English into target languages, and the
 245 remaining 5 are translated from French sources.
 246 MAFAND-MT data is professionally translated
 247 by native speakers of the target languages with
 248 a compensation (Adelani et al., 2022). However,
 249 we observed a range of common MT error types
 250 presented in Table 1.

Category	Error description
Typos	Misspellings or character-level mistakes in the translation
Grammar	Grammatical errors (e.g., agreement, tense, syntax)
Fluency	Unnatural or awkward phrasing; non-native flow
Bilingual	Interference or overly literal translation from English
Incomplete	Translation omits some(all) part(s) of the source meaning
Addition	Adds information not present in the source
Omission	Removes information present in the source

Table 2: **Error categories reported during manual analysis of translation quality.** The table summarizes recurrent error types observed with their description.

4 Translation Review Pipeline

251 We used a two-stage review pipeline to assess and
 252 review translation quality. In the first stage, LLM
 253 automatically evaluates each translation pair and
 254 flags potential errors. In the second stage, native-

256 speaker verify the flagged cases and provide cor-
 257 rections where necessary, shown in Figure 1.

258 **LLM-assisted Translation Review** We used
 259 GPT-5.1² as a judge to review the translation. We
 260 probed GPT-5.1 for each parallel text pair to assess
 261 translation quality and classify the translation as
 262 correct or incorrect. We further instructed this
 263 LLM to suggest the types of translation error(s)
 264 presented in Table 2 and the correct translation
 265 versions if the reply was incorrect at first.

266 **Human Translation Correction** We assign a
 267 minimum of two native-speaker volunteers per lan-
 268 guage for a total of 11 languages. Human trans-
 269 lation reviews translation errors flagged by LLMs as
 270 incorrect and verifies the corrected translations
 271 proposed by the LLMs. To facilitate native speaker
 272 review of LLM suggested corrections, we design
 273 an interactive annotation interface that displays the
 274 source English text, the original MAFAND-MT hu-
 275 man translation, and the LLM proposed alternative.
 276 If either translation is flagged as incorrect, the tool
 277 highlights common error categories. Annotators
 278 are also provided with an option to supply a new
 279 translation if both existing options are incorrect.
 280 Details of the translation review guidelines and the
 281 annotation interface are provided in Appendix A.

5 Human vs LLM Audit Agreement

282 We analyze the agreement between humans and
 283 GPT-5.1 in translation quality review and assess
 284 whether the LLM’s suggested corrections are use-
 285 ful.
 286

²<https://openai.com/index/gpt-5-1/>, Dec 2025

Language	Translation direction	# Test set	LLM predict "Correct"	LLM predict "Incorrect"	Human vs LLM Trans. Agree. %	Human vs LLM Errors. Agree.
Amharic	eng-amh	1,037	271	766	0.85	0.23
Hausa	eng-hau	1,500	680	820	0.23	0.25
Igbo	eng-ibo	1,500	884	964	0.78	0.18
Kinyarwanda	eng-kin	1,006	390	616	0.55	0.00
Luo (Dholuo)	eng-luo	1,500	398	1,102	0.96	0.09
Nigerian Pidgin	eng-pcm	1,564	1009	555	0.43	0.20
Shona	eng-sna	1,005	368	637	0.24	0.15
Swahili (Kiswahili)	eng-swa	1,835	748	787	0.30	0.05
Tswana (Setswana)	eng-tsn	1,500	844	656	0.75	0.18
Twi (Akan-Twi)	eng-twi	1,500	265	1285	0.91	0.11
Yoruba	eng-yor	1,558	585	973	0.45	0.16

Table 3: **The MAFAND-MT test set dataset details with human and LLM agreement analysis.** LLM predict **Correct** and LLM predict **Incorrect** columns are the number of translation responses from LLMs. **Human vs LLM Translation Agreement** is the percentage agreement between humans and LLM to say the original translation is correct or incorrect. **Human vs LLM Translation Agreement** is Cohen’s Kappa translation error label agreement between Human and LLM. We target only the test set data, and each language has its own Source (English) and target translations.

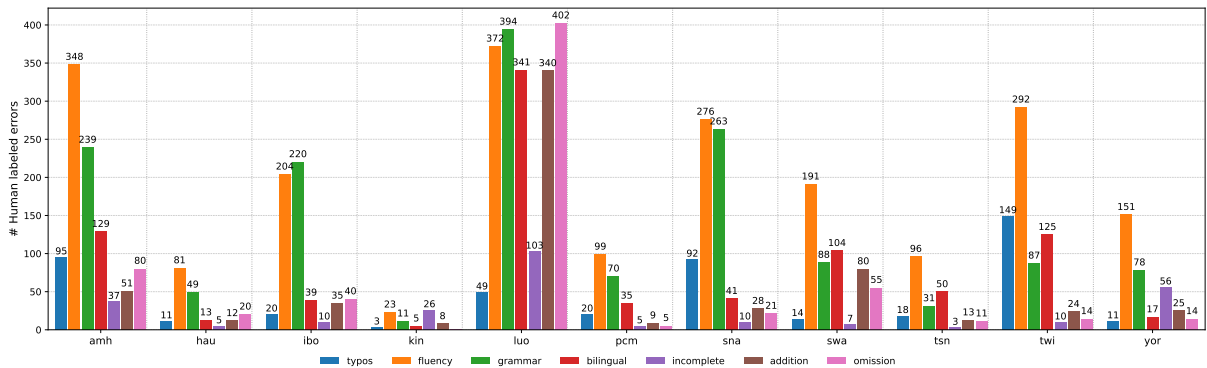


Figure 2: The statistics of Human-labeled translation errors for the MAFAND-MT test dataset. The bar graph illustrates the number of types of translation errors where Amharic (amh), Luo-Dholuo (luo), Igbo (ibo), Shona (sna), Swahili (swa), and Twi (twi) languages have more error types statistically.

Can LLMs assist in reviewing translation quality for low-resource languages?

Based on AfroBench, a benchmark for evaluating LLMs on African languages, proprietary LLMs such as the GPT family consistently outperform widely used open-source models for the machine translation task (Ojo et al., 2025). Motivated by this observation, we used the latest GPT-5.1 as a translation quality reviewer. Native speakers of each target language are then presented with two options: the original MAFAND-MT translation and the revised translation produced by the GPT-5.1 (reviewer model).

The level of agreement between GPT-5.1 and humans in assessing translation quality is reported in Table 3. The results suggest that GPT-5.1 can support translation quality verification within a human-in-the-loop framework. Notably, based on GPT-

generated revisions, a substantial number of translations were judged by native speakers to be of higher quality than the original human-produced translations, for example, amh (118), hau (628), yor (124), swa (548), and pcm (320).

Regarding agreement on translation error labels between the LLM (GPT-5.1) and human native speakers (shown in Table 3, *Human vs. LLM Errors Agreement* column), the Cohen’s Kappa scores are consistently low, ranging from 0 (minimum) to 0.25 (maximum). This low agreement can be attributed to several factors: (1) translation errors are annotated in a multi-label setting, (2) GPT-5.1 tends to over-predict multiple error types for a single translation pair, and (3) there is substantial disagreement in cases where native speakers labeled most LLM translations as incorrect.

State-of-the-art LLMs, such as GPT-5.1, can as-

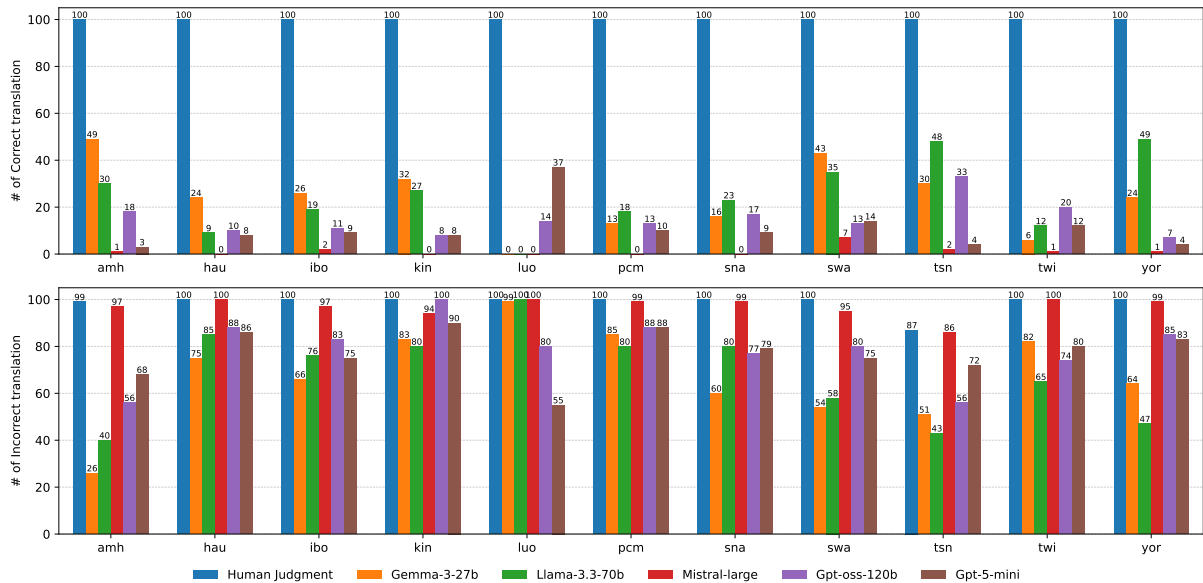


Figure 3: Human judgments versus LLM-as-a-judge on 200 randomly selected samples: 100 translations labeled as correct and 100 as incorrect by human evaluators.

sist with machine translation quality auditing and provide translation suggestions for low-resource languages. A considerable number of translations that were generated by LLM were judged by native speakers to be of acceptable quality. Moreover, a majority of annotators (64.7%) reported that qualifying each translation took 1–3 minutes, and 29.4% reported 30 seconds to 1 minute, indicating that LLM pre-auditing can provide a substantial time-saving benefit for human annotators. Regarding the helpfulness of adding LLM-generated translation as an option during quality audits of machine translation data, native speaker responds 42.2% yes, it was helpful, 42.1% partially useful, and 17.6% not helpful. However, the selected LLM reviewer (GPT-5.1) remains far from a stand-alone solution for the automatic correction of human translations in African languages. High-quality, corrected MAFAND-MT test set data will therefore be valuable for further machine translation experiments and evaluations.

Can LLMs serve as judges of translation quality for low-resource languages? We select the following popular open-source LLMs for the LLM-as-a-judge evaluation: Gemma-3-27B (Team et al., 2025), LLaMA-3.3-70B (Grattafiori et al., 2024), GPT-oss-120B (OpenAI et al., 2025), and Mistral-123B (OpenAI et al., 2024). In addition, we include the closed-source model GPT-5-mini (Hurst et al., 2024). We evaluate these models on 200 randomly sampled, human-labeled instances (100 correct and

100 incorrect translations), where incorrect translations are further annotated with fine-grained error labels.

As shown in Figure 3, the agreement between human judgments and LLM-based judges is substantially low. In particular, for translations labeled as correct by human annotators, most LLM judges incorrectly classify them as incorrect. Models such as Mistral-123B, GPT-oss-120B, and GPT-5-mini fail to identify correct translations reliably. A similar trend is observed for incorrect translations: the LLM judges tend to label nearly all translation pairs as incorrect and frequently overpredict multiple error categories for a single translation. The agreement on error labeling between humans and LLM-as-judge models is reported in Appendix 6.

How is the COMET score a reliable quality estimation for low-resource languages? We apply SSA-COMET-QE (Sub-Saharan African Crosslingual Optimized Metric for Evaluation of Translation) - an improved version of AfriCOMET (Wang et al., 2024), a robust and automatic metric for machine translation quality estimation (Rei et al., 2020). It receives a pair (source sentence, translation in target language), and returns a score ranging from 0 (semantically unrelated) to 1 (high quality) that reflects the quality of the translation (Li et al., 2025). The quality estimation score (SSA-COMET-QE) for the MAFAND-MT dataset before and after quality check is presented in Figure 4. Based on the SSA-COMET-QE score, each language has

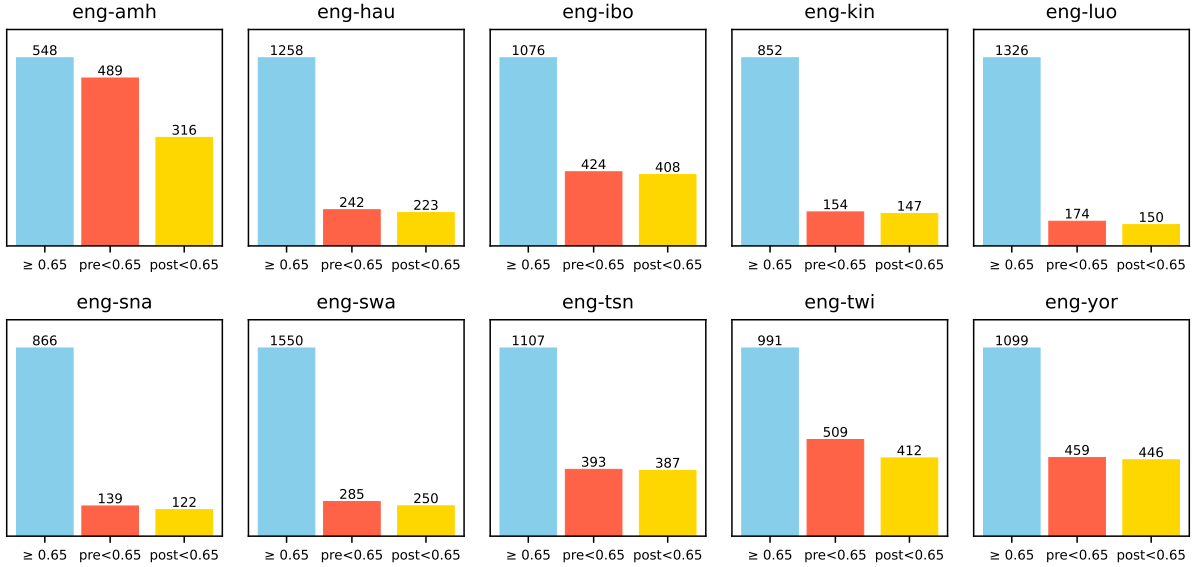


Figure 4: **SSA-COMET-QE translation scores across language pairs before (pre) and after (post) applying quality audit.** Scores range from 0 to 1, with higher values indicating better translation quality. The en-pcm direction has lower COMET scores; only 21% of the data have >0.5 SSA-COMET-QE because the language is not included in the specifically fine-tuned model.

translations ranging from 70 to 400 parallel texts that have a quality estimation score of less than 0.6. As illustrated in Figure 4, an empty source and output or a single character or word translation can still receive a low score. Across all languages, the maximum score observed is 0.84; even perfect translations do not approach a score of 1.0, and most translation outputs fall within a narrow range between 0.65 and 0.70. Notably, some outputs in an incorrect language receive higher scores than null or random outputs produced in the correct language. These findings indicate a divergence between SSA-COMET-QE scores and human judgments in low-resource language settings, warranting further investigation into the reliability and behaviour of such evaluation metrics for low-resource languages. While we improved the statistics of low-scoring translations, the improvement remains modest due to the quality of the SSA-COMET-QE metrics, as shown in Table 4.

6 Benchmarking Improved Dataset

Automated scores provide a cost-effective and rapid approximation of quality, which is essential for machine translation system performance and for quick feedback on evolving models (Kocmi et al., 2024). While human judgment remains the gold standard, we evaluate our approach in three ways: human judgments, LLM-as-a-judge assessments,

Language	Source text	Translation	SSA-COMET score
amh	NULL	NULL	0.38
amh	l	l	0.44
amh	" "	" "	0.48
ibo	Naira	Naira	0.54
ibo	Ihiala	Ihiala	0.44
swa	ICANNWiki	ICANNWiki	0.49
swa	(CC BY 2.0)	(CC BY 2.0)	0.55
pcm	"GOD DEY,"	"GOD DEY,"	0.33

Table 4: Translation pair examples with SSA-COMET-QE score. Short translation examples from the dataset: empty (NULL), single-character, and a word with perfect translation but a low SSA-COMET-QE score.

and automatic evaluation metrics such as BLEU and chrF for MT output. To make benchmark results, we select popular open-source machine translation models such as multilingual Aya-101 (Üstün et al., 2024), NLLB 600M and NLLB 3B (Team et al., 2022), and Google Translate³.

Results Analysis The benchmark results are presented in Table 5. As the results show, the quality-audited version (corrected MAFAND-MT) outperforms the original evaluation across all settings, indicating that the dataset has been improved. Google Translate outperforms other open-source models, with NLLB-3B being the strongest open alternative in terms of parameter size, followed by NLLB-600M. However, the result remains close to the

³<https://cloud.google.com/translate>, Dec 2025

Models	Metrics	amh	hau	ibo	kin	luo	pcm	sna	swa	tsn	twi	yor	Avg.
NLLB 600M	BLEU (MAFAND-MT)	4.92	7.68	17.00	23.11	11.02	7.83	8.72	27.18	25.23	8.10	8.57	13.58
	BLEU (Corrected)	10.04	10.04	18.26	23.13	12.58	7.87	10.41	30.00	25.36	10.10	8.84	15.15
	chrF (MAFAND-MT)	25.49	36.83	47.22	55.70	39.92	27.89	42.55	56.18	56.04	36.87	29.78	41.32
	chrF (Corrected)	34.06	38.81	48.10	55.72	41.61	27.93	44.56	58.32	56.09	38.43	30.22	43.08
Aya 101	BLEU (MAFAND-MT)	3.40	7.12	9.66	09.64	2.84	13.52	5.85	19.96	3.96	2.73	4.22	7.54
	BLEU (Corrected)	6.10	8.98	10.27	9.61	2.92	13.50	6.32	21.83	4.02	3.25	4.21	8.27
	chrF (MAFAND-MT)	20.00	34.88	37.35	37.85	13.94	45.71	27.98	47.77	23.06	23.96	17.16	29.97
	chrF (Corrected)	25.65	36.46	37.90	37.88	14.05	45.74	28.81	49.29	23.08	24.21	17.21	30.93
NLLB 3B	BLEU (MAFAND-MT)	5.62	8.44	20.24	26.60	12.43	4.59	9.39	28.81	28.00	8.69	10.88	14.88
	BLEU (Corrected)	11.23	10.59	21.60	26.61	14.25	04.51	10.92	32.01	28.17	10.71	11.26	16.53
	chrF (MAFAND-MT)	26.51	37.87	50.13	59.28	41.80	11.63	43.18	57.54	57.79	38.92	32.36	41.55
	chrF (Corrected)	35.20	39.81	51.01	59.31	43.66	11.56	45.02	59.91	57.85	40.71	32.83	43.35
Google Trans.	BLEU (MAFAND-MT)	6.35	8.71	15.60	25.33	8.49	00.00	10.69	30.81	31.98	8.87	14.26	14.64
	BLEU (Corrected)	15.62	11.41	16.77	25.20	13.16	00.00	12.58	34.73	32.05	10.85	14.95	17.03
	chrF (MAFAND-MT)	27.99	38.85	48.91	64.71	38.02	00.00	45.44	59.52	61.38	40.37	37.33	42.05
	chrF (Corrected)	39.65	41.12	49.86	64.67	41.21	00.00	47.55	62.31	61.40	42.30	37.99	44.37

Table 5: **Zero-shot evaluation benchmark results.** The result compares the original (MAFAND-MT) with the **corrected** translation version of the MAFAND-MT test set. All translation directions are from English to target languages. Nigerian Pidgin (pcm) is not supported by Google Translate.

original. This might be due to two main reasons: 1) the target languages are low-resource languages - the evaluated models do not well represent the languages, 2) the evaluation metrics problem, such as COMET, as discussed in Table 4.

7 Native-Speakers Feedback

Following completion of the translation quality audit, we conducted a survey to gather qualitative insights from native speaker annotators regarding the source text quality, the use of LLM-generated suggestions, and the time they spent on the task. The primary feedbacks are summarized below:

Quality Issues in English Source Text: Annotator feedback revealed significant challenges stemming from the quality and composition of the source (English) text. In particular, source-side noise was observed in segments derived from social media (X/Twitter); as the entries frequently contained platform-specific metadata, such as user handles (@usernames) and hashtags(#), and suffered from syntactic fragmentation due to character limits. Additionally, annotators identified instances of language leakage, where the source text was labeled as English but included content in other languages, which the annotators were unable to interpret. Such issues negatively impacted the annotation process and introduced ambiguity in translation and sentiment interpretation.

Literal Translation: Annotators observed that GPT-5.1 often defaulted to overly literal translations, struggling to balance literal and conceptual

meaning, especially for metaphors. This issue was exacerbated by archaic or unnatural terms in human references, which conflicted with modern usage. As a result, current benchmarks may over-reward word-level matching while overlooking native fluency, and, in some cases, GPT-5.1 produced non-existent words in the target language.

8 Conclusions

In this work, we evaluated a subset of a widely used machine translation evaluation dataset for African languages (MAFAND-MT), covering 11 languages and all test set splits, with support from native speakers. The evaluation process involved judging whether each translation was correct or incorrect, labeling the type(s) of translation error(s) for incorrect translations, and producing corrected translations when necessary. Our analysis revealed that the original translations contained various types of errors relative to the source text. We corrected the MAFAND-MT test set using native speakers and LLMs as assistants at different stages. We show that attention should be given to translated evaluation sets, and that relying solely on automatic evaluation metrics for MT quality evaluation may not align with human assessments. A combination of human evaluation, using LLMs as judges, and automatic metrics is recommended. The improved MAFAND-MT test set and the accompanying quality-audit annotation tool, provide valuable resources for researchers conducting further machine translation quality analysis and evaluation.

490 Limitations

491 Our work is not without limitations.

492 First, it focuses on a single MT dataset because
493 recruiting volunteer native speakers for each target
494 language is difficult. However, our pipeline is re-
495 producible and this work can be extended to other
496 African languages' translated dataset such as 1)
497 machine translation dataset: FLORES 101 (Goyal
498 et al., 2022) and FLORES+ (Gordeev et al., 2024)
499 and 2) health (e.g., AFRIDOC-MT (Alabi et al.,
500 2025) and AfriMed-QA (Nimo et al., 2025)), 3)
501 mathematics word problem (e.g., AfriGSM (Ade-
502 lani et al., 2025b)), and 4) general knowledge
503 and reasoning (e.g., AfriMMLU and AfriXNLI
504 (Adelani et al., 2025b)) and MAFAND-MT dataset
505 (Adelani et al., 2022).

506 Second, we focused only on the quality audit
507 of the test set, as it is urgent, and research work
508 reports are based on test set results. The same way
509 can be extended for other split sets, such as training
510 and validation sets.

511 References

512 Idris Abdulmumin, Michael Beukman, Jesujoba O. Al-
513 abi, Chris Emezue, Everlyn Asiko, Tosin Adewumi,
514 Shamsuddeen Hassan Muhammad, Mofetoluwa
515 Adeyemi, Oreen Yousuf, Sahib Singh, and Tajud-
516 deen Rabi Gwadabe. 2022. [Separating Grains from
517 the Chaff: Using Data Filtering to Improve Multi-
518 lingual Translation for Low-Resourced African Lan-
519 guages](#). In *Proceedings of the Seventh Conference on
520 Machine Translation (WMT)*, pages 1001–1014, Abu
521 Dhabi, United Arab Emirates (Hybrid). Association
522 for Computational Linguistics.

523 Idris Abdulmumin, Sthembiso Mkhwanazi, Mahlatse
524 Mbooi, Shamsuddeen Hassan Muhammad,
525 Ibrahim Said Ahmad, Neo Putini, Miehleketo Math-
526 ebula, Matimba Shingange, Tajuddeen Gwadabe,
527 and Vukosi Marivate. 2024. [Correcting FLORES
528 Evaluation Dataset for Four African Languages](#). In
529 *Proceedings of the Ninth Conference on Machine
530 Translation*, pages 570–578, Miami, Florida, USA.
531 Association for Computational Linguistics.

532 David Ifeoluwa Adelani, Jesujoba Oluwadara Alabi,
533 Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel
534 Reid, Dana Ruitter, Dietrich Klakow, Peter Nabende,
535 Ernie Chang, Tajuddeen Gwadabe, Freshia Sackey,
536 Bonaventure F. P. Dossou, Chris Emezue, Colin
537 Leong, Michael Beukman, Shamsuddeen H. Muham-
538 mad, Guyo D. Jarso, Oreen Yousuf, and 26 others.
539 2022. [A Few Thousand Translations Go a Long Way!
540 Leveraging Pre-trained Models for African News
541 Translation](#). In *Proceedings of the 2022 Conference
542 of the North American Chapter of the Association for*

*Computational Linguistics: Human Language Tech-
543 nologies*, pages 3053–3070, Seattle, United States.
544 Association for Computational Linguistics. 545

546 David Ifeoluwa Adelani, Jessica Ojo, Israel Abebe
547 Azime, Jian Yun Zhuang, Jesujoba Alabi, Xu-
548 anli He, Millicent Ochieng, Sara Hooker, Andiswa
549 Bukula, En-Shiun Annie Lee, and 1 others. 2025a.
550 [Irokobench: A new benchmark for african languages
551 in the age of large language models](#). In *Proceedings
552 of the 2025 Conference of the Nations of the Amer-
553 icas Chapter of the Association for Computational
554 Linguistics: Human Language Technologies (Volume
555 1: Long Papers)*, pages 2732–2757.

556 David Ifeoluwa Adelani, Jessica Ojo, Israel Abebe Az-
557 ime, Jian Yun Zhuang, Jesujoba Oluwadara Alabi,
558 Xuanli He, Millicent Ochieng, Sara Hooker, Andiswa
559 Bukula, En-Shiun Annie Lee, Chiamaka Ijeoma
560 Chukwunke, Happy Buzaaba, Blessing Kudzaishe
561 Sibanda, Godson Koffi Kalipe, Jonathan Mukiibi,
562 Salomon Kabongo Kabenamualu, Foutse Yuehgoh,
563 Mmasibidi Setaka, Lolwethu Ndolela, and 8 others.
564 2025b. [IrokoBench: A New Benchmark for African
565 Languages in the Age of Large Language Models](#).
566 In *Proceedings of the 2025 Conference of the Na-
567 tions of the Americas Chapter of the Association for
568 Computational Linguistics: Human Language Tech-
569 nologies (Volume 1: Long Papers)*, pages 2732–2757,
570 Albuquerque, New Mexico. Association for Compu-
571 tational Linguistics.

572 Jesujoba Oluwadara Alabi, Israel Abebe Azime,
573 Míaoran Zhang, Cristina España-Bonet, Rachel
574 Bawden, Dawei Zhu, David Ifeoluwa Adelani,
575 Clement Oyeleke Odoje, Idris Akinade, Iffat Maab,
576 Davis David, Shamsuddeen Hassan Muhammad, Neo
577 Putini, David O. Ademuyiwa, Andrew Caines, and
578 Dietrich Klakow. 2025. [AFRIDOC-MT: Document-
579 level MT Corpus for African Languages](#). In *Proceed-
580 ings of the 2025 Conference on Empirical Methods in
581 Natural Language Processing*, pages 27758–27794,
582 Suzhou, China. Association for Computational Lin-
583 guistics.

584 Antonios Anastasopoulos, Alessandro Cattelan, Zi-
585 Yi Dou, Marcello Federico, Christian Federmann,
586 Dmitriy Genzel, Francisco Guzmán, Junjie Hu, Mac-
587 duff Hughes, Philipp Koehn, Rosie Lazar, Will Lewis,
588 Graham Neubig, Mengmeng Niu, Alp Öktem, Eric
589 Paquin, Grace Tang, and Sylwia Tur. 2020. [TICO-19:
590 the Translation Initiative for COvid-19](#). In *Proceed-
591 ings of the 1st Workshop on NLP for COVID-19 (Part
592 2) at EMNLP 2020*, Online. Association for Compu-
593 tational Linguistics.

594 Tadesse Destaw Belay, Kedir Yassin Hussien,
595 Sukairaj Hafiz Imam, Ibrahim Said Ahmad, Isa
596 Inuwa-Dutse, Abrham Belete Haile, Grigori Sidorov,
597 Iqra Ameer, Idris Abdulmumin, Tajuddeen Gwad-
598 abe, Vukosi Marivate, Seid Muhie Yimam, and
599 Shamsuddeen Hassan Muhammad. 2025. [The
600 Rise of AfricaNLP: Contributions, Contributors,
601 and Community Impact \(2005-2025\)](#). *Preprint*,
602 arXiv:2509.25477.

719	sera Tapo, Nishant Subramani, Artem Sokolov, Clay-	<i>Findings of the Association for Computational Lin-</i>	777
720	tone Sikasote, Monang Setyawan, Supheakmungkol	<i>guistics: NAACL 2024</i> , pages 182–195, Mexico City,	778
721	Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, An-	Mexico. Association for Computational Linguistics.	779
722	nette Rios, Isabel Papadimitriou, Salomey Osei, Pe-		
723	dro Ortiz Suarez, and 33 others. 2022. Quality at	Jessica Ojo, Odunayo Ogundepo, Akintunde Oladipo,	780
724	a Glance: An Audit of Web-Crawled Multilingual	Kelechi Ogueji, Jimmy Lin, Pontus Stenertorp, and	781
725	Datasets . <i>Transactions of the Association for Com-</i>	David Ifeoluwa Adelani. 2025. AfroBench: How	782
726	<i>putational Linguistics</i> , 10:50–72.	Good are Large Language Models on African Lan-	783
		guages? In <i>Findings of the Association for Computa-</i>	784
727	Senyu Li, Jiayi Wang, Felermino D. M. A. Ali, Colin	<i>tional Linguistics: ACL 2025</i> , pages 19048–19095,	785
728	Cherry, Daniel Deutsch, Eleftheria Briakou, Rui	Vienna, Austria. Association for Computational Lin-	786
729	Sousa-Silva, Henrique Lopes Cardoso, Pontus Stene-	guistics.	787
730	torp, and David Ifeoluwa Adelani. 2025. SSA-		
731	COMET: Do LLMs Outperform Learned Metrics	OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason	788
732	in Evaluating MT for Under-Resourced African Lan-	Ai, Sam Altman, Andy Applebaum, Edwin Arbus,	789
733	guages? In <i>Proceedings of the 2025 Conference on</i>	Rahul K. Arora, Yu Bai, Bowen Baker, Haiming Bao,	790
734	<i>Empirical Methods in Natural Language Processing</i> ,	Boaz Barak, Ally Bennett, Tyler Bertao, Nivedita	791
735	pages 12990–13009, Suzhou, China. Association for	Brett, Eugene Brevdo, Greg Brockman, Sebastien	792
736	Computational Linguistics.	Bubeck, and 108 others. 2025. gpt-oss-120b & gpt-	793
		oss-20b Model Card . <i>Preprint</i> , arXiv:2508.10925.	794
737	Jessy Lin, Geza Kovacs, Aditya Shastry, Joern Wue-	OpenAI, Albert : Jiang, Alexandre Sablayrolles, Alexis	795
738	bker, and John DeNero. 2022. Automatic Correc-	Tacnet, Alok Kothari, Antoine Roux, Arthur Mensch,	796
739	tion of Human Translations . In <i>Proceedings of the</i>	Audrey Herblin-Stoop, Augustin Garreau, Austin	797
740	<i>2022 Conference of the North American Chapter of</i>	Birky, Bam4d, Baptiste Bout, Baudouin de Moni-	798
741	<i>the Association for Computational Linguistics: Hu-</i>	cault, Blanche Savary, Carole Rambaud, Caroline	799
742	<i>man Language Technologies</i> , pages 494–507, Seattle,	Feldman, Devendra Singh Chaplot, Diego de las	800
743	United States. Association for Computational Lin-	Casas, Diogo Costa, and 51 others. 2024. Mistral	801
744	guistics.	Large 2 . https://huggingface.co/mistralai/	802
		Mistral-Large-Instruct-2407 . Large-scale	803
745	Yixin Liu, Kejian Shi, Katherine He, Longtian Ye,	instruct-tuned language model released by Mistral	804
746	Alexander Fabbri, Pengfei Liu, Dragomir Radev, and	AI.	805
747	Arman Cohan. 2024. On Learning to Summarize		
748	with Large Language Models as References . In <i>Pro-</i>	Shenbin Qian, Archchana Sindhujan, Minnie Kabra,	806
749	<i>ceedings of the 2024 Conference of the North Amer-</i>	Diptesh Kanojia, Constantin Orasan, Tharindu Ranas-	807
750	<i>ican Chapter of the Association for Computational</i>	inghe, and Fred Blain. 2024. What do Large Lan-	808
751	<i>Linguistics: Human Language Technologies (Volume</i>	guage Models Need for Machine Translation Eval-	809
752	<i>1: Long Papers)</i> , pages 8647–8664, Mexico City,	uation? In <i>Proceedings of the 2024 Conference on</i>	810
753	Mexico. Association for Computational Linguistics.	<i>Empirical Methods in Natural Language Processing</i> ,	811
		pages 3660–3674, Miami, Florida, USA. Association	812
754	Qingyu Lu, Baopu Qiu, Liang Ding, Kanjian Zhang,	for Computational Linguistics.	813
755	Tom Kocmi, and Dacheng Tao. 2024. Error Analysis		
756	Prompting Enables Human-Like Translation Eval-	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	814
757	uation in Large Language Models . In <i>Findings of</i>	Lavie. 2020. COMET: A Neural Framework for MT	815
758	<i>the Association for Computational Linguistics: ACL</i>	Evaluation . In <i>Proceedings of the 2020 Conference</i>	816
759	<i>2024</i> , pages 8801–8816, Bangkok, Thailand. Associ-	<i>on Empirical Methods in Natural Language Process-</i>	817
760	ation for Computational Linguistics.	<i>ing (EMNLP)</i> , pages 2685–2702, Online. Association	818
		for Computational Linguistics.	819
761	Charles Nimo, Tobi Olatunji, Abraham Toluwase	Nithya Sambasivan, Shivani Kapania, Hannah Highfill,	820
762	Owodunni, Tassallah Abdullahi, Emmanuel Ayo-	Diana Akrong, Praveen Paritosh, and Lora M Aroyo.	821
763	dele, Mardhiyah Sanni, Ezinwanne C. Aka, Fola-	2021. “Everyone wants to do the model work, not	822
764	funmi Omofoye, Foutse Yuehgoh, Timothy Faniran,	the data work” : Data Cascades in High-Stakes AI . In	823
765	Bonaventure F. P. Dossou, Moshood O. Yekini, Jonas	<i>Proceedings of the 2021 CHI Conference on Human</i>	824
766	Kemp, Katherine A Heller, Jude Chidubem Omeke,	<i>Factors in Computing Systems, CHI ’21</i> , New York,	825
767	Chidi Asuzu Md, Naome A Etori, Aïméroù Ndiaye,	NY, USA. Association for Computing Machinery.	826
768	Ifeoma Okoh, and 7 others. 2025. AfriMed-QA:		
769	A Pan-African, Multi-Specialty, Medical Question-	Pratik Rakesh Singh, Kritarth Prasad, Mohammadi Zaki,	827
770	Answering Benchmark Dataset . In <i>Proceedings</i>	and Pankaj Wasnik. 2025. In-Domain African Lan-	828
771	<i>of the 63rd Annual Meeting of the Association for</i>	guages Translation Using LLMs and Multi-armed	829
772	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	Bandits . In <i>Proceedings of the Sixth Workshop on</i>	830
773	pages 1948–1973, Vienna, Austria. Association for	<i>African Natural Language Processing (AfricaNLP</i>	831
774	Computational Linguistics.	<i>2025)</i> , pages 167–175, Vienna, Austria. Association	832
		for Computational Linguistics.	833
775	Antoine Nzeyimana. 2024. Low-resource neural ma-		
776	chine translation with morphological modeling . In		

834	Zhen Tan, Dawei Li, Song Wang, Alimohammad	Pavanpankaj Vegi, Sivabhavani J, Biswajit Paul,	894
835	Beigi, Bohan Jiang, Amrita Bhattacharjee, Man-	Prasanna K R, and Chitra Viswanathan. 2022a.	895
836	sooreh Karami, Jundong Li, Lu Cheng, and Huan Liu.	ANVITA-African: A Multilingual Neural Machine	896
837	2024. Large Language Models for Data Annotation	Translation System for African Languages . In <i>Pro-</i>	897
838	and Synthesis: A Survey . In <i>Proceedings of the 2024</i>	<i>ceedings of the Seventh Conference on Machine</i>	898
839	<i>Conference on Empirical Methods in Natural Lan-</i>	<i>Translation (WMT)</i> , pages 1090–1097, Abu Dhabi,	899
840	<i>guage Processing</i> , pages 930–957, Miami, Florida,	United Arab Emirates (Hybrid). Association for Com-	900
841	USA. Association for Computational Linguistics.	putational Linguistics.	901
842	Gongbo Tang, Oreen Yousuf, and Zeying Jin. 2024. Im-	Pavanpankaj Vegi, Sivabhavani J, Biswajit Paul, Abhi-	902
843	proving BERTScore for Machine Translation Evalu-	nav Mishra, Prashant Banjare, Prasanna K R, and	903
844	ation Through Contrastive Learning . <i>IEEE Access</i> ,	Chitra Viswanathan. 2022b. WebCrawl African : A	904
845	12:77739–77749.	Multilingual Parallel Corpora for African Languages .	905
846	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya	In <i>Proceedings of the Seventh Conference on Ma-</i>	906
847	Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin,	<i>chine Translation (WMT)</i> , pages 1076–1089, Abu	907
848	Tatiana Matejovicova, Alexandre Ramé, Morgane	Dhabi, United Arab Emirates (Hybrid). Association	908
849	Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey	for Computational Linguistics.	909
850	Cideron, Jean bastien Grill, Sabela Ramos, Edouard	Jiayi Wang, David Ifeoluwa Adelani, Sweta Agrawal,	910
851	Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev,	Marek Masiak, Ricardo Rei, Eleftheria Briakou,	911
852	and 197 others. 2025. Gemma 3 Technical Report .	Marine Carpuat, Xuanli He, Sofia Bourhim, An-	912
853	Preprint , arXiv:2503.19786.	diswa Bukula, Muhidin Mohamed, Temitayo Ola-	913
854	NLLB Team, Marta R. Costa-jussà, James Cross, Onur	toye, Tosin Adewumi, Hamam Mokayed, Christine	914
855	Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hef-	Mwase, Wangui Kimotho, Foutse Yuehgho, An-	915
856	fernan, Elahe Kalbassi, Janice Lam, Daniel Licht,	nuoluwapo Aremu, Jessica Ojo, and 39 others. 2024.	916
857	Jean Maillard, Anna Sun, Skyler Wang, Guillaume	AfriMTE and AfriCOMET: Enhancing COMET to	917
858	Wenzek, Al Youngblood, Bapi Akula, Loic Barrault,	Embrace Under-resourced African Languages . In	918
859	Gabriel Mejia Gonzalez, Prangthip Hansanti, and	<i>Proceedings of the 2024 Conference of the North</i>	919
860	20 others. 2022. No Language Left Behind: Scal-	<i>American Chapter of the Association for Computa-</i>	920
861	ing Human-Centered Machine Translation . <i>Preprint</i> ,	<i>tional Linguistics: Human Language Technologies</i>	921
862	arXiv:2207.04672.	<i>(Volume 1: Long Papers)</i> , pages 5997–6023, Mexico	922
863	Atnafu Lambebo Tonja, Israel Abebe Azime,	City, Mexico. Association for Computational Lin-	923
864	Tadesse Destaw Belay, Mesay Gemedo Yigezu,	guistics.	924
865	Moges Ahmed Ah Mehamed, Abinew Ali Ayele,	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale,	925
866	Ebrahim Chekol Jibril, Michael Melese Woldeyohan-	Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and	926
867	nis, Olga Kolesnikova, Philipp Slusallek, Dietrich	Colin Raffel. 2021. mT5: A Massively Multilingual	927
868	Klakow, and Seid Muhie Yimam. 2024. EthioLLM:	Pre-trained Text-to-Text Transformer . In <i>Proceed-</i>	928
869	Multilingual Large Language Models for Ethiopian	<i>ings of the 2021 Conference of the North American</i>	929
870	Languages with Task Evaluation . In <i>Proceedings</i>	<i>Chapter of the Association for Computational Lin-</i>	930
871	<i>of the 2024 Joint International Conference on</i>	<i>guistics: Human Language Technologies</i> , pages 483–	931
872	<i>Computational Linguistics, Language Resources</i>	498, Online. Association for Computational Linguis-	932
873	<i>and Evaluation (LREC-COLING 2024)</i> , pages	tics.	933
874	6341–6352, Torino, Italia. ELRA and ICCL.	Biao Zhang, Philip Williams, Ivan Titov, and Rico Sen-	934
875	Marcos V Treviso, Nuno M Guerreiro, Sweta Agrawal,	nrich. 2020. Improving Massively Multilingual Neu-	935
876	Ricardo Rei, José Pombal, Tania Vaz, Helena Wu,	ral Machine Translation and Zero-Shot Translation .	936
877	Beatriz Silva, Daan Van Stigt, and Andre Martins.	In <i>Proceedings of the 58th Annual Meeting of the As-</i>	937
878	2024. xTower: A Multilingual LLM for Explaining	<i>sociation for Computational Linguistics</i> , pages 1628–	938
879	and Correcting Translation Errors . In <i>Findings of the</i>	1639, Online. Association for Computational Linguis-	939
880	<i>Association for Computational Linguistics: EMNLP</i>	tics.	940
881	<i>2024</i> , pages 15222–15239, Miami, Florida, USA.	Ahmet Üstün, Viraat Aryabumi, Zheng Yong, Wei-Yin	
882	Association for Computational Linguistics.	Ko, Daniel D’souza, Gbemileke Onilude, Neel Bhan-	
883	Ahmet Üstün, Viraat Aryabumi, Zheng Yong, Wei-Yin	dari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Fred-	
884	Ko, Daniel D’souza, Gbemileke Onilude, Neel Bhan-	die Vargus, Phil Blunsom, Shayne Longpre, Niklas	
885	dari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Fred-	Muennighoff, Marzieh Fadaee, Julia Kreutzer, and	
886	die Vargus, Phil Blunsom, Shayne Longpre, Niklas	Sara Hooker. 2024. Aya Model: An Instruction Fine-	
887	Muennighoff, Marzieh Fadaee, Julia Kreutzer, and	tuned Open-Access Multilingual Language Model .	
888	Sara Hooker. 2024. Aya Model: An Instruction Fine-	In <i>Proceedings of the 62nd Annual Meeting of the</i>	
889	tuned Open-Access Multilingual Language Model .	<i>Association for Computational Linguistics (Volume 1:</i>	
890	In <i>Proceedings of the 62nd Annual Meeting of the</i>	<i>Long Papers)</i> , pages 15894–15939, Bangkok, Thai-	
891	<i>Association for Computational Linguistics (Volume 1:</i>	land. Association for Computational Linguistics.	
892	<i>Long Papers)</i> , pages 15894–15939, Bangkok, Thai-		
893	land. Association for Computational Linguistics.		

Appendix

941

A Translation Annotation Guideline

942

The native speakers do not have information about the two translation option sources, which are the original human translation and LLM translation. The native speakers were given the following guidelines.

943

Reviewing translation errors: At this stage, the native speakers read and evaluate the translation quality in parallel with the given source English text.

944

945

946

- Read the source English sentence.
- Read the translated in Translation 1 and Translation 2 (One is the original translation, and the other is the LLM translation; we randomly shuffle the content positions of the two translations when displaying for the annotators).
- For both Translation 1 and Translation 2, choose Correct or Incorrect.
- If any of the translations are Incorrect, select one or more error types (multi-label error selection approach) that best describe the error by ticking from the list of error types, shown in Table 2.

947

948

949

950

951

952

953

Correcting the translations: If both Translation 1 and Translation 2 are marked Incorrect with the corresponding error types, provide a new correct translation in the given text box; the UI is shown in Appendix B.

954

955

956

B Annotation Tool UI

957

Figure B shows a screenshot of our machine translation quality audit annotation tool UI.

958

Field	Text	Annotation
English	The Commander of the Faithful (a.s.) teaches us a lesson.	
Corrected translation	Type the correct translation here... Please provide the corrected translation only if both translations above are incorrect.	
Translation 1	Amirul Muminin (a.s) yana koyarwa da mu, yana fadin cewa: 'Wanda ya nemi gaskiya amma ya kuskure, bai yi daidai da wanda ya nemi karya kuma ya same ta ba'.	<input type="checkbox"/> Typos <input type="checkbox"/> Grammar <input type="checkbox"/> Fluency <input type="checkbox"/> Bilingual <input checked="" type="checkbox"/> Incomplete <input checked="" type="checkbox"/> Addition <input type="checkbox"/> Omission
Translation 2	Amirul Muminin (a.s.) yana koyar da mu darasi.	<input type="checkbox"/> Incorrect <input type="checkbox"/> Choose <input checked="" type="checkbox"/> Correct <input type="checkbox"/> Incorrect

Figure 5: **Annotation tool UI for Hausa language.** The tool will be publicly released upon acceptance of the work for further machine translation and other NLP dataset quality audits with additional features.

C LLM-as-a-Judge Prompts

Prompt: LLM-as-a-judge for translation quality analysis

You are an expert translation quality analyst with deep knowledge of machine translation evaluation from English to African languages. Analyze the following English → {target_lang_name} translation..

Possible translation error types (choose one or more when incorrect):

- Typos : misspellings or character mistakes in the translation
- Grammar : grammatical errors (agreement, tense, syntax)
- Fluency : unnatural or awkward phrasing / non-native flow
- Bilingual : interference or literal translation from English
- Incomplete : translation omits part(s) of the source meaning
- Addition : adds information not present in the source
- Omission : removes information present in the source

SOURCE (English):

{eng_text}

TRANSLATION ({target_language_name}):

{tgt_text}

Your task is to follow the below rules exactly:

- 1) Decide if the translation is correct or incorrect contextually. If correct, respond with status "correct". Only mark as 'incorrect' when the meaning changes. Do NOT mark minor differences as errors.
- 2) If incorrect, set status "incorrect", pick one or more error types from the taxonomy, and give a short explanation for each type of error.
- 3) IF incorrect, ALSO PROVIDE a corrected translation in the target language in the field "correct_translation" (a fluent, natural translation that preserves source meaning correctly).
- 4) If correct, set "correct_translation" to null.
- 5) Return ONLY valid JSON (no extra commentary). Use this exact structure:

```
{{
  "status": "correct" | "incorrect",
  "errors": [
    {"type": "<one_of_taxonomy>", "description": "<short explanation>"}
  ],
  "correct_translation": "<correct text or null>"
}}
```

Taxonomy reference:

{taxonomy} Return only the JSON.

D Translation Error Labeling Agreement Between Human vs LLMs

Figure 6 the overlap between human and LLM-as-a-judge for translation error labeling. The statistics show that LLMs are overpredicting error types relative to humans in the targeted low-resource languages.

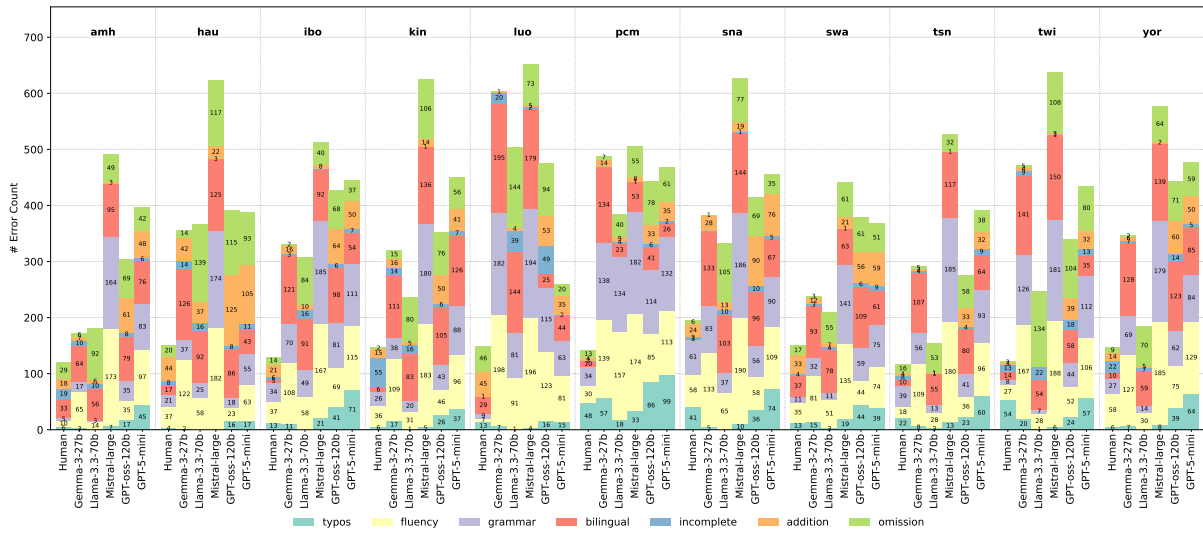


Figure 6: Translation error labeling overlap between Human and LLM-as-a-judge.