

# AfriVox: Probing Multilingual and Accent Robustness of Speech LLMs

Anonymous ACL submission

## Abstract

Recent advances in multimodal and speech-native large language models (LLMs) have delivered impressive speech recognition, translation, understanding, and question-answering capabilities for high-resource languages. However, African languages and non-native French or English accents remain dramatically under-represented in benchmarks limiting the understanding and applicability of leading LLMs for millions of francophone and anglophone users in low-resource settings. We presents AfriVox, an open-source benchmark (including novel domain-specific and unscripted datasets) across 20 African languages, African-accented French, Arabic, and 100+ African English accents, contrasting leading multimodal speech LLMs with traditional unimodal automatic speech transcription (ASR) and translation (AST) models. Our analysis reveals significant language coverage variation, surprising LLM translation performance gains (e.g. Gemini), robustness concerns with unscripted speech, and substantial performance disparities for "supported" African languages. We profile the strengths, limitations, and language support of each model, and conduct the first targeted fine-tuning of a modern speech LLM (Qwen2.5-Omni) for three Nigerian languages, exceeding SOTA, and achieving up to 54% relative WER reduction and significant BLEU gains, offering practical guidance for implementers seeking to serve local language users.

## 1 Introduction

The transformative impact of LLMs in global technology—especially speech-enabled and multimodal LLMs—has opened new frontiers for human-computer interaction (AlSaad et al., 2024). Major recent breakthroughs, such as OpenAI’s GPT-4o (Hurst et al., 2024), Google Gemini (Google DeepMind, 2024), and Meta’s SeamlessM4T (Barrault et al., 2023), have enabled voice-based applications that promise to make informa-

tion and services more accessible, especially in regions where text literacy and high-resource language proficiency may be limiting factors (Peng et al., 2025).

Across Africa, LLM-powered systems are already being deployed in sectors like health, agriculture, and financial inclusion, operating in large languages via text interfaces (Olatunji et al., 2023; Nazi and Peng, 2024; Al-Garadi et al., 2025). However, as voice-native and multilingual LLMs have rapidly improved (Bai et al., 2024; Google DeepMind, 2024), technology implementers across Africa are eager to shift towards more natural, relatable, and intuitive speech-driven interfaces that truly reflect users’ language preferences and linguistic diversity (Sanni et al., 2025).

Despite this demand, no comprehensive benchmark exists that systematically evaluates modern speech LLMs on African languages and accents (Adelani et al., 2025; Ojo et al., 2025). Existing benchmarks such as MLS, mSTEB, NaijaVoices, and ML-SUPERB 2.0 include very limited African language coverage and lack recent domain-specific, real-world unscripted speech, especially for emerging LLM architectures (Pratap et al., 2020a; Beyene et al., 2025; Emezue et al., 2025; Shi et al., 2024). Most performance claims are based on high-resource languages, providing little actionable guidance to African technology teams deciding whether to trust LLMs for local deployment (Reid et al., 2021).

To bridge this gap, we introduce AfriVox, a unified benchmark suite aggregating and extending multiple African speech datasets and releasing two novel datasets covering parliamentary speech from 4 countries and health-focused conversations in 20 African languages. We use AfriVox to answer two critical questions for implementers: (1) Which speech LLMs reliably support certain African languages and (2) How do leading multimodal LLMs compare with traditional leading ASR/AST models

on understanding realistic African speech? Should implementers switch from unimodal ASR models to LLMs? Our contributions are as follows:

- We curate and open-source **AfriVox**, the most comprehensive benchmark to date for African language ASR and AST, with detailed language support analysis.
- We open-source 2 novel datasets under a CC-BY-NC-SA license: **Afrispeech-Parliamentary**, transcribed accented-English parliamentary proceedings from 4 African countries; and **Afrivox-Medical**, a health-focused read-speech translation and transcription dataset in 19 African languages.
- We conduct the first systematic, reproducible evaluation of state-of-the-art speech LLMs and unimodal models across 20 languages and 100+ English, French, and Arabic accents.
- We provide detailed error analysis and practical guidance, including fine-tuning experiments with Qwen2.5-Omni on major Nigerian languages using only moderate data.

By clarifying the current capabilities and limitations of speech LLMs on African languages, we aim to empower both researchers and implementers to build more equitable language technology.

## 2 Related works

Recent years have seen remarkable progress in speech and multimodal large language models (LLMs) (Yu et al.), driven by advances in self-supervised learning, scaling laws, and reinforcement learning techniques (Ghosh et al., 2024). However, these improvements have disproportionately benefited high-resource languages, with African languages still underrepresented in both model training and evaluation (Adelani et al., 2025; Ojo et al., 2025).

**Multilingual Speech Benchmarks:** Benchmarks such as MLS (Multilingual LibriSpeech) (Pratap et al., 2020a), mSTEB (Beyene et al., 2025), and ML-SUPERB 2.0 (Shi et al., 2024) have provided valuable evaluation resources, but offer limited coverage of African languages, and their data is primarily read speech or synthetic in nature. ML-SUPERB 2.0 and mSTEB in particular have improved multilingual evaluation rigor, yet it covers only a handful of African

languages and lacks representation of diverse accents and real-world conversational domains (Pratap et al., 2020a). Our benchmark, AfriVox, addresses these gaps by including (a) a broader and more granular set of African languages and accents, (b) domain-specific, real-world audio (e.g., parliamentary sessions, healthcare dialogues), and (c) explicit evaluation of both unimodal and state-of-the-art multimodal LLMs.

**Speech and Multimodal LLMs:** Large-scale unimodal models such as Whisper (Radford et al., 2023a), MMS (Denisov and Vu, 2024), and Parakeet (Galvez et al., 2024) have demonstrated robust speech recognition performance in high-resource settings, but their reliability in African language tasks remains largely anecdotal (Ojo et al., 2025). Recent multimodal models—including Google AudioPaLM (Rubenstein et al., 2023), Meta SeamlessM4T (Barrault et al., 2023), Qwen-Audio (Chu et al., 2024), and Gemini (Google DeepMind, 2024)—promise to unify speech, text, and translation tasks, but have yet to be systematically benchmarked on African data (Adelani et al., 2025).

**Parameter-Efficient Fine-Tuning (PEFT):** Scaling LLMs for downstream tasks in low-resource settings can be prohibitively expensive. PEFT approaches such as LoRA (Karimi Mahabadi et al., 2021), Adapters (Han et al., 2024), and QLoRA (Dettmers et al., 2023) enable practical model adaptation by training only a small subset of parameters. However, most prior studies have focused on high-resource or Asian languages (Bai et al., 2024); little is known about their impact on speech LLMs for African contexts (Emezue et al., 2025).

**African Speech Datasets:** Public African speech corpora—including NCHLT (Barnard et al., 2014), CommonVoice (Ardila et al., 2020), and FLEURS (Conneau et al., 2023)—have played a vital role, but coverage, accent diversity, and domain relevance remain limited. Recent datasets such as AfriSpeech (Olatunji et al., 2025) and NaijaVoices (Emezue et al., 2025) have begun to address these challenges. Our work builds on and expands these efforts, contributing new datasets and a unified benchmark for comprehensive, reproducible evaluation.

**Distinctive Contributions:** To our knowledge, this work is the first to 1) Aggregate and compare both unimodal and multimodal speech LLMs

across 20+ African languages and 100+ English accents, 2) Include new, diverse African speech test sets, 3) Provide practical, data-driven guidance for implementers on the suitability of LLMs vs. traditional ASR for local deployment, and 4) Systematically analyze model performance, error types, and the impact of PEFT for low-resource African speech recognition and translation.

### 3 Methodology

Dataset	Hours	Speakers	Accents
NCHLT	2.24	8	1
AfriSpeech-200	18.68	750	108
CV-17 En-Afr	0.11	46	9
Afrispeech-Parl	42.17	~1651	4
<b>Total</b>	<b>63.20</b>	<b>~2455</b>	<b>108</b>

Table 1: Summary of African-accented English speech datasets.

Language	Region	Language Family	# Speakers
Afrikaans	South	IndoWest (Germanic)	7.2M
Akan	West	Niger-Congo (Kwa)	24M
Amharic	East	Afro-Asiatic (Semitic)	35M
Egyptian Arabic	North	Afro-Asiatic (Semitic)	78M
French	West	Indo-European (Romance)	320M
Fula	West	Niger-Congo (Atlantic)	36.8M
Gaa	West	Niger-Congo (Kwa)	0.7M
Hausa	West	Afro-Asiatic (Chadic)	54M
Ibo	West	Niger-Congo (Volta-Niger)	31M
Kinyarwanda	East	Niger-Congo (Bantu)	15M
Luganda	East	Niger-Congo (Bantu)	5.6M
Northern Sotho	South	Niger-Congo (Bantu)	4.6M
Shona	South	Niger-Congo (Bantu)	8.4M
Southern Sotho	South	Niger-Congo (Bantu)	5.6M
Swahili	East	Niger-Congo (Bantu)	87M
Tswana	South	Niger-Congo (Bantu)	8.2M
Twi	West	Niger-Congo (Kwa)	4.4M
Xhosa	South	Niger-Congo (Bantu)	8M
Yoruba	West	Niger-Congo (Yoruboid)	45M
Zulu	South	Niger-Congo (Bantu)	13.6M

Table 2: Language, region, family, and number of speakers.

Dataset	# Langs	Hours	Speakers
NCHLT	6	12.75	36
CV-17	10	16.89	670
FLEURS	13	14.44	1595
OpenSLR	3	0.31	372
Bible TTS	3	0.47	3
NaijaVoices <sup>1</sup>	3	1800	5000
FISD <sup>2</sup>	3	0.05	23
AfriVox-Medical <sup>3</sup>	19	36.63	1179
<b>Total Hours</b>		<b>1878.52</b>	

Table 3: Summary of multilingual speech datasets.

### 3.1 Benchmark Design and Datasets

We design the AfriVox benchmark to evaluate speech LLMs and ASR/AST models on realistic African language and accent use-cases. This benchmark unifies and expands existing corpora, incorporating both new and public datasets to maximize coverage and relevance.

#### 3.1.1 African-Accented English Speech (AES)

**Sources:** NCHLT (Barnard et al., 2014), AfriSpeech-200 (Olatunji et al., 2025), Common Voice 17 (Ardila et al., 2020) (filtered for African accents), and a newly-curated AfriSpeech-Parl dataset<sup>1</sup> with transcribed Parliamentary Proceedings from 4 African countries (Ghana, Kenya, Nigeria, and South Africa).

**Coverage:** Over 63 hours, 2,000+ speakers, 12 countries, and 108 distinct African English accents (Table 1).

**Curation:** Common Voice was filtered using speaker metadata and manual accent validation. Parliamentary recordings were transcribed and quality-controlled by native speakers; only utterances with >80% reviewer ratings were included.

#### 3.1.2 Multilingual African Speech (MLS)

**Sources:** Existing open source transcription datasets—NCHLT, Common Voice 17 (filtered for African languages), FLEURS (Conneau et al., 2023), OpenSLR, BibleTTS (Meyer et al., 2022), NaijaVoices (Emezue et al., 2025)—and newly-curated AfriVox-Medical<sup>3</sup>, a health-related read-speech multilingual translation and transcription dataset of simulated text conversations across 20 languages).

For translation, we include FLEURS, CoVoST (Wang et al., 2020), NaijaVoices, IWSLT-LRST (Cettolo et al., 2017), and AfriVox-Medical<sup>3</sup>.

**Coverage:** 20 languages across 7 datasets, 8,000+ speakers, >1,800 hours of audio (Tables 2 and 3). Coverage includes both high-population and low-resource languages, and features diverse linguistic families (Niger-Congo, Afro-Asiatic, etc.).

<sup>1</sup><https://huggingface.co/datasets/naijavoices/naijavoices-dataset>

<sup>2</sup><https://github.com/Ashesi-Org/Financial-Inclusion-Speech-Dataset>

<sup>3</sup>URL to be added after anonymity period

<sup>4</sup>URL to be added after anonymity period

**Ethics and Quality:** All audio files are mono-channel WAV at 16kHz. All data is either open-source or collected with explicit consent. New data is transcribed and quality-checked by native speakers. All contributors and reviewers were fairly compensated via a crowdsourcing platform<sup>5</sup>.

## 3.2 Models Evaluated

We benchmarked a mix of unimodal and multimodal models:

**Unimodal ASR:** Canary (Puvvada et al., 2024), Parakeet (NeMo and Suno.ai, 2023), Whisper (Medium/Large) (Radford et al., 2023b), MMS (with/without language adapters) (Pratap et al., 2024)

**Unimodal AST [X->En]:** Whisper, MMS

**Multimodal LLMs:** Meta SeamlessM4T (Aharoni et al., 2019), Google Gemini-2.0-Flash (Google DeepMind, 2024), OpenAI GPT-4o (OpenAI et al., 2024), Alibaba Qwen2.5-Omni (Chu et al., 2024)

Languages supported for each language are presented in Table 6. Models were chosen for their reported state-of-the-art performance, public availability, language coverage (supporting one or more African languages), or relevance to real-world deployment in Africa. All were used in their pre-trained, off-the-shelf forms unless otherwise specified.

## 3.3 Fine-Tuning

**Data:** We fine-tuned on the NaijaVoices dataset—1,800 hours, 5,000+ speakers, balanced by gender and age, spanning 3 Nigerian languages—Hausa, Igbo, and Yoruba.

**Model:** Qwen2.5-Omni, a 10B multimodal and multilingual LLM was selected for PEFT due to its open-source availability, multilingual support, and relatively small size (compute limitations).

**Fine-tuning:** We fine-tuned on four NVIDIA 3090 GPUs with approximately 280 hours of speech per language, using LoRA (rank 8, alpha 32) applied to all linear layers while freezing the vision encoder. We trained for three epochs using a learning rate of 1e-4 and a warmup ratio of 0.05 with bfloat16 precision, with a batch size of 256. Prompt formatting details are included in the Appendix A

<sup>1</sup>URL to be added after anonymity period

## 3.4 Evaluation

### 3.4.1 Tasks

**Automatic Speech Recognition (ASR):** Transcribe audio into native script.

**Automatic Speech Translation (AST):** Translate audio into English text.

**Prompting:** All models were tested with consistent, standardized prompts (zero-shot and few-shot) for fairness and reproducibility (see Appendix A for details).

**Post-processing:** Outputs were normalized for punctuation, casing, and diacritics to ensure comparability.

**Reproducibility:** All code, model configurations, and new data will be open-sourced; results are reported for single runs.

### 3.4.2 Metrics and Human Evaluation

**ASR:** Word Error Rate (WER)

**AST:** BLEU (Papineni et al., 2002), chrF (Popović, 2015), and AfriCOMET-STL (Wang et al., 2023).

**Human Evaluation:** Conducted for translation quality validation and metric selection; see Appendix Table 17.

### 3.4.3 Addressing Benchmark Contamination

We note and analyze the potential for older public datasets to appear in model pretraining, and explicitly distinguish between “old” (NCHLT, Common-Voice) and “new” (AfriSpeech-200, Afrispeech-Parl) data in analysis to measure true generalization.

## 4 Results and Analysis

Tables 4 and 5 present the transcription results on the African-Accented English Speech and Multilingual African Speech datasets. Results presented are for single runs. The results indicate that, in most cases, unimodal models outperformed the multimodal models. While Table 6 shows multimodal models edges over unimodal models on the speech translation task. Additionally, Table 8 shows the comparison between the results of the base and fine-tuned Qwen 2.5 Omni model. A detailed breakdown of results by individual languages is provided in Appendix A. We provide the following analysis based on the findings from our experimental results:



Model	Old			New	
	Lib	NC	CV	Af	Parl
Canary	1.48	<b>10.05</b>	<b>8.41</b>	38.03	27.38
Parakeet	<b>1.40</b>	11.33	9.48	34.96	21.89
Whisper M	3.02	10.17	12.39	30.81	28.53
Whisper L	2.01	10.10	12.54	<b>26.49</b>	<b>19.29</b>
MMS	12.63	32.11	23.09	61.19	107.41
M4T	2.89	32.96	10.40	49.75	54.68
Gemini	3.03	14.19	13.76	28.12	21.63
GPT-Aud.	5.26	86.52	26.76	36.54	41.88
Qwen2	1.60	25.14	11.16	49.61	57.43

Table 4: Word Error Rates (WER) across African-accented English speech data sources and Librispeech test-clean [Lib](Panayotov et al., 2015). Af: Afrispeech, NC: NCHLT, CV: Common Voice, Parl: Parliamentary Proceedings. Models in the top section are unimodal ASRs while those below are multimodal LLMs.

#### 4.1 Widespread Variation in African Language and Accent Performance and Support

Table 4 and 5 reveal that, despite recent advances and better coverage of African languages, both unimodal and multimodal speech models exhibit substantial performance gaps on African languages and non-native English accents when compared with large languages and native accents (Multilingual Librispeech). Wide variation within models exist, most evident with Seamless and Whisper for supported languages. Consistent with multilingual claims in its documentation, Gemini outperforms GPT-4o by a wide margin sometimes with 2-4x better WER.

**Unusually High Error Rates for Supported Languages:** On African-accented English, state-of-the-art unimodal ASR models (e.g., Whisper Large-v3) display a 10–15x increase in Word Error Rate (WER) compared to standard benchmarks—for example, WER rises from 2.0% (LibriSpeech) to 26–38% (Afrispeech, NCHLT). For African languages, WERs routinely exceed 50% and, for some languages (e.g., Yoruba, Hausa, Swahili), surpass 100%, despite self-reported "support" for these languages, indicating nearly unintelligible output. These results suggest that simply including African data in pretraining does not provide performance guarantees.

**Multimodal Model Language Coverage:** Multimodal LLMs (e.g., Gemini, SeamlessM4T, GPT-4o) support more African languages than most unimodal ASR/AST models and can be prompted with-

out explicit language labels, but their accuracy often lags unimodal models for transcription. SeamlessM4Tv2, for example, shows particularly strong results for Southern and Eastern African languages, providing clues about the language distribution in its training data.

#### 4.2 Transcription vs. Translation: Unimodal and Multimodal Model Trends

**Transcription (ASR):** As shown in Table 5, Unimodal models, especially MMS with language adapters, outperform multimodal LLMs for exact transcription in most African languages. Gemini stands out, outperforming MMS across multiple supported languages, indicating progress towards more inclusive multimodal LLMs. However, with WERs still over 20% for several languages and accented speech, top ASR models and LLMs still struggle with accent/language diversity and noisy or spontaneous speech.

**Translation (AST):** Multimodal models (Table 6), especially Gemini and SeamlessM4T, significantly outperform unimodal baselines on low-resource African language audio-to-English translation. They achieve higher BLEU and AfriCOMET-STL scores, and provide more semantically faithful translations, particularly on longer, context-rich utterances. Appendix Table 17 shows AfriCOMET-STL’s correlation with human evaluation.

#### 4.3 Robustness to Real-World Speech

**Realistic Noisy Conditions:** As shown in Table 4 All models perform worst on the parliamentary proceedings dataset, which contains high ambient noise, overlapping speakers, and real-world spontaneous speech. Here, WERs for even the best models double relative to clean, read speech, demonstrating that accented English speech transcription is still an unsolved problem. MMS is most notable in this regard, with a 5x collapse in WER, likely demonstrating an over-reliance on clean/read speech during training.

**Accent and Dialect Variability:** Table 5 reveals a consistent trend with accented French. Besides GPT-4o, inclusion of accent-diverse datasets exposes weaknesses in all models, with WER dropping by roughly 2x. Performance is notably worse on underrepresented accents and dialects—even for languages like French with larger training resources.

Language	Canary-1b	Whisper medium	Whisper large-v3	MMS-1b all	Qwen2.5	Seamless-M4T Large-v2	Gpt-4o audio-preview	Gemini-2.0 flash
English (M. Lib)	<b>3.03</b>	6.80	3.53	17.63	16.32	4.68	9.63	6.63
French (M. Lib)	<b>4.06</b>	8.90	5.38	19.30	10.43	6.82	22.71	5.23
Spanish (M. Lib)	-	-	-	17.35	-	6.76	21.25	<b>3.22</b>
Afrikaans	-	68.87	45.43	48.73	-	18.41	84.36	<b>18.02</b>
Akan	-	-	-	<b>62.92</b>	-	-	104.02	67.04
Amharic	-	447.26	165.83	67.52	-	<b>44.05</b>	245.4	55.88
Arabic	-	39.49	29.72	44.94	-	51.26	31.88	<b>14.44</b>
French	9.67	13.95	9.31	33.93	24.14	15.90	22.29	<b>9.12</b>
Fulani	-	-	-	<b>56.78</b>	-	86.85	157.03	66.11
Ga	-	-	-	-	-	-	172.73	<b>87.27</b>
Hausa	-	180.29	95.11	40.47	-	-	118.60	<b>38.48</b>
Igbo	-	-	-	<b>50.33</b>	-	70.03	112.23	66.68
Kinyarwanda	-	-	-	<b>36.73</b>	-	-	135.75	58.44
Luganda	-	-	-	28.85	-	<b>16.39</b>	131.19	59.89
Pedi	-	-	-	<b>41.43</b>	-	-	119.29	70.69
Sesotho	-	-	-	-	-	-	158.21	<b>59.30</b>
Shona	-	193.21	110.35	<b>30.7</b>	-	76.05	90.51	38.84
Swahili	-	117.7	62.75	28.37	-	<b>16.25</b>	73.96	25.88
Tswana	-	-	-	-	-	-	133.46	<b>54.85</b>
Twi	-	-	-	<b>51.09</b>	-	-	98.86	67.13
Xhosa	-	-	-	42.24	-	-	130.79	<b>39.32</b>
Yoruba	-	213.88	93.77	39.59	-	<b>37.43</b>	101.14	43.42
Zulu	-	-	-	43.19	-	52.53	135.84	<b>30.02</b>

Table 5: **WER (%) by model and language on the Multilingual African Speech transcription dataset.** Bold values mark the lowest (best) WER for each language. "-" indicates the language is not supported by the model. The first section of the table shows baseline performance on Multilingual LibriSpeech (Pratap et al., 2020b)

#### 4.4 Fine-Tuning Unlocks Substantial Gains

**Parameter-Efficient Fine-Tuning (PEFT):** Table 7 zooms in on model performance for the 3 languages selected for fine-tuning. Although all 3 languages were unsupported by Qwen2.5-Omni, Table 8 shows that fine-tuning on just 280 hours per language from NaijaVoices yields a 54% reduction in WER and up to 21-point gains in BLEU for Igbo, Hausa, and Yoruba, exceeding SOTA (MMS) on Igbo. AfriCOMET-STL (translation performance) more than doubles for all three languages, exceeding SOTA (Gemini) on Igbo.

**Low-Resource Potential:** These results demonstrate that, even with moderate in-domain data, open-source speech LLMs can be rapidly adapted for African languages using PEFT, offering a viable path for local teams.

#### 4.5 Error Analysis

**Verbatim vs. Paraphrase:** Multimodal models frequently paraphrase or summarize rather than provide exact transcriptions (Figure 1), which is unsuitable for many ASR use cases. In contrast, unimodal ASR models are more likely to attempt verbatim output, albeit with higher rates of inser-

tion and substitution errors on low-resource languages.

<p><b>Example 1 [Af]: Paraphrasing and Audio Description</b>  <b>Reference:</b> Adana spoke with doctor  <b>Qwen2-Audio:</b> A woman is saying Adana spoke with doctor</p>
<p><b>Example 2 [Parl.]: Content Description</b>  <b>Reference:</b> We had legislation in front of this house to push down funds to the lowest levels of service delivery in the counties, namely the wards. What we have discussed this morning is that a lot of areas are against.  <b>GPT Audio:</b> The audio content discusses legislation aimed to allocate funds to the lowest levels of service delivery in counties, specifically the wards. It indicates that there is some disagreement or istance to this approach in various areas.</p>

Figure 1: Examples of paraphrasing and audio description.

**Hallucinations:** Both Whisper and Canary sometimes hallucinate content—repeating text or filling silent segments with unrelated words as shown in Figure 2. Multimodal models are prone to “helpful” completions (Figure 2), such as generating plausible answers to questions not present in the audio.

**Contextual Mistranslations:** In AST tasks, multimodal models occasionally substitute synonyms

Language	Canary 1b	Whisper medium	Whisper large-v3	Qwen2.5	SeamlessM4T Large-v2	Gpt-4o audio-preview	Gemini-2.0 flash
Afrikaans	-	0.57	0.65	-	0.73	0.71	<b>0.80</b>
Akan	-	-	-	-	-	0.34	<b>0.38</b>
Amharic	-	0.23	0.27	-	0.64	0.42	<b>0.79</b>
Arabic	-	0.65	0.70	-	0.80	0.81	<b>0.85</b>
French	0.65	0.70	0.73	<b>0.80</b>	0.79	0.78	<b>0.80</b>
Fulani	-	-	-	-	0.19	0.30	<b>0.35</b>
Ga	-	-	-	-	-	0.24	<b>0.29</b>
Hausa	-	0.16	0.19	-	0.17	0.37	<b>0.65</b>
Igbo	-	-	-	-	0.25	0.29	<b>0.37</b>
Kinyarwanda	-	-	-	-	-	0.29	<b>0.54</b>
Luganda	-	-	-	-	0.57	0.47	<b>0.59</b>
Pedi	-	-	-	-	-	0.31	<b>0.39</b>
Sesotho	-	-	-	-	0.23	0.35	<b>0.50</b>
Shona	-	0.18	0.21	-	<b>0.73</b>	0.47	0.61
Swahili	-	0.32	0.42	-	-	0.76	<b>0.81</b>
Tswana	-	-	-	-	<b>0.56</b>	0.32	0.46
Twi	-	-	-	-	<b>0.41</b>	0.33	0.32
Xhosa	-	-	-	-	-	0.35	<b>0.66</b>
Yoruba	-	0.18	0.20	-	-	0.36	<b>0.49</b>
Zulu	-	-	-	-	-	0.40	<b>0.71</b>

Table 6: **AfriComet-STL scores** across the languages for each model. "-" means the language is not supported by the model. The highlighted scores are the best score per language

Language	Whisper medium	Whisper large-v3	MMS-1b all	Seamless-M4T-v2 Large	Gpt-4o audio-preview	Gemini-2.0 flash	Qwen2.5
Hausa	186.23	96.99	<b>39.37</b>	-	119.74	52.16	126.81
Igbo	-	-	<b>48.81</b>	66.27	117.84	87.32	198.68
Yoruba	213.41	97.51	<b>44.05</b>	44.62	107.25	78.46	120.84

Table 7: **Transcription WER % for each model–language pair on the NaijaVoices subset of the Multilingual African Speech transcription dataset**; the lowest (best) WER per language is shown in bold. "-" indicates the language is not supported by the model.

Language	ASR (WER)		AST (STL)	
	Base	Finetuned	Base	Finetuned
Hausa	126.81	<b>50.54</b>	0.19	<b>0.39</b>
Igbo	198.68	<b>42.41</b>	0.18	<b>0.54</b>
Yoruba	120.84	<b>71.29</b>	0.20	<b>0.29</b>

Table 8: Qwen-Omni2 ASR (WER score) and AST (AfriComet-STL) Performance Before and After Fine-Tuning

or miss important words (Figure 3), producing contextually plausible but non-literal translations—highlighted by AfriCOMET-STL (Figure 6, which better captures adequacy than BLEU alone.

**Noise Sensitivity:** All models suffer from degraded output under overlapping speech and real-world noise, with frequent failures to segment speakers or filter background sounds, indicating model’s failure to adequately generalize to real-

world spontaneous speech.

#### 4.6 Implications for Inclusive Voice Technology

Our findings have clear implications for implementers, researchers, and product teams:

**Model Selection:** For applications requiring exact transcription—such as legal or medical records—unimodal ASR models remain preferable where they support the target language. However, for conversational interfaces or translation tasks, recent multimodal LLMs (e.g. Gemini) offer broader language coverage and better semantic translation, even in low-resource settings.

**Fine-Tuning Value:** The dramatic improvements achieved with PEFT fine-tuning on Qwen2.5-Omni (Figure 8) highlight a promising pathway for African NLP practitioners. Moderate, domain-specific datasets can unlock substantial gains, mak-





datasets (e.g., CommonVoice, NCHLT) may overlap with pretraining data for popular models, possibly inflating apparent model performance relative to unseen, truly out-of-domain audio. Our results on newly-curated datasets are more reliable but still limited by size and scope.

**Evaluation Scope:** Most evaluations focus on transcription and direct audio-to-English translation. We do not benchmark the full range of speech LLM multimodal abilities (e.g., dialog, spoken question answering), nor do we exhaustively test different prompting strategies or task configurations due to compute constraints.

**Fine-Tuning Experiments:** Our parameter-efficient fine-tuning is limited to three Nigerian languages, using moderate (not minimal) amounts of labeled data. Results may not generalize to ultra-low-resource languages or domains with dramatically less data available.

**Noise and Real-World Testing:** While AfriVox includes challenging real-world audio, our robustness analysis is not exhaustive. Further work should explore adversarial noise, code-switching, and multi-speaker dialog in more depth.

Despite these constraints, AfriVox establishes a practical and extensible blueprint for ongoing evaluation and improvement of speech and text LLMs in Africa. We hope this work will catalyze further open data sharing, community-driven evaluation, and development of voice AI systems that genuinely serve Africa’s linguistic diversity.

## References

- 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 Chukwuneke, Happy Buzaaba, Blessing Sibanda, Godson Kalipe, Jonathan Mukiibi, Salomon Kabongo, Foutse Yuehgoh, Mmasibidi Setaka, Lolwethu Ndolela, and 8 others. 2025. [IrokoBench: A New Benchmark for African Languages in the Age of Large Language Models](#). *arXiv preprint*. ArXiv:2406.03368 [cs].
- Roece Aharoni, Melvin Johnson, and Orhan Firat. 2019. [Massively multilingual neural machine translation](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3874–3884, Minneapolis, Minnesota. Association for Computational Linguistics.
- Mohammed Al-Garadi, Tushar Mungle, Abdulaziz Ahmed, Abeed Sarker, Zhuqi Miao, and Michael E. Matheny. 2025. [Large Language Models in Healthcare](#). *arXiv preprint*. ArXiv:2503.04748 [cs].
- Rawan AlSaad, Alaa Abd-alrazaq, Sabri Boughorbel, Arfan Ahmed, Max-Antoine Renault, Rafat Damseh, and Javaid Sheikh. 2024. [Multimodal Large Language Models in Health Care: Applications, Challenges, and Future Outlook](#). *Journal of Medical Internet Research*, 26(1):e59505. Company: Journal of Medical Internet Research Distributor: Journal of Medical Internet Research Institution: Journal of Medical Internet Research Label: Journal of Medical Internet Research Publisher: JMIR Publications Inc., Toronto, Canada.
- Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. 2020. [Common voice: A massively-multilingual speech corpus](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4218–4222, Marseille, France. European Language Resources Association.
- Ye Bai, Jingping Chen, Jitong Chen, Wei Chen, Zhuo Chen, Chuang Ding, Linhao Dong, Qianqian Dong, Yujiao Du, Kepan Gao, and 1 others. 2024. [Seed-asr: Understanding diverse speech and contexts with llm-based speech recognition](#). *arXiv preprint arXiv:2407.04675*.
- Etienne Barnard, Marelle H. Davel, Charl van Heerden, Febe de Wet, and Jaco Badendorst. 2014. The nchlt speech corpus of the south african languages. In *4th Workshop on Spoken Language Technologies for Under-Resourced Languages (SLTU 2014)*, pages 194–200.
- Loïc Barrault, Yu-An Chung, Mariano Cora Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, and 1 others. 2023. [Seamlessm4t-massively multilingual & multimodal machine translation](#). *arXiv preprint arXiv:2308.11596*.
- Luel Hagos Beyene, Vivek Verma, Min Ma, Jesujoba O Alabi, Fabian David Schmidt, Joyce Nakatumba-Nabende, and David Ifeoluwa Adelani. 2025. [msteb: Massively multilingual evaluation of llms on speech and text tasks](#). *arXiv preprint arXiv:2506.08400*.
- Mauro Cettolo, Marcello Federico, Luisa Bentivogli, Jan Niehues, Sebastian Stüker, Katsuhito Sudoh, Koichiro Yoshino, and Christian Federmann. 2017. [Overview of the IWSLT 2017 evaluation campaign](#). In *Proceedings of the 14th International Conference on Spoken Language Translation*, pages 2–14, Tokyo, Japan. International Workshop on Spoken Language Translation.
- Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, Chang Zhou, and Jingren Zhou. 2024. [Qwen2-audio technical report](#). *Preprint*, arXiv:2407.10759.

643	Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang,	Zabir Al Nazi and Wei Peng. 2024. <a href="#">Large language</a>	697
644	Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara	<a href="#">models in healthcare and medical domain: A review.</a>	698
645	Rivera, and Ankur Bapna. 2023. Fleurs: Few-shot	In <i>Informatics</i> , volume 11, page 57. MDPI. Issue: 3.	699
646	learning evaluation of universal representations of		
647	speech. In <i>2022 IEEE Spoken Language Technology</i>	NVIDIA NeMo and Suno.ai. 2023. <a href="#">Parakeet tdt 1.1b:</a>	700
648	<i>Workshop (SLT)</i> , pages 798–805. IEEE.	<a href="#">An asr model with fastconformer and tdt decoder.</a>	701
649	Pavel Denisov and Ngoc Thang Vu. 2024. Teaching	Jessica Ojo, Odunayo Ogundepo, Akintunde Oladipo,	702
650	a multilingual large language model to understand	Kelechi Ogueji, Jimmy Lin, Pontus Stenertorp, and	703
651	multilingual speech via multi-instructional training.	David Ifeoluwa Adelani. 2025. <a href="#">AfroBench: How</a>	704
652	<i>arXiv preprint arXiv:2404.10922</i> .	<a href="#">Good are Large Language Models on African Lan-</a>	705
653	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and	<a href="#">guages?</a> <i>arXiv preprint</i> . ArXiv:2311.07978 [cs].	706
654	Luke Zettlemoyer. 2023. Qlora: Efficient finetuning	Tobi Olatunji, Tejumade Afonja, Aditya Yadavalli,	707
655	of quantized llms. <i>Advances in neural information</i>	Chris Chinenye Emezue, Sahib Singh, Bonaventure	708
656	<i>processing systems</i> , 36:10088–10115.	F. P. Dossou, Joanne Osuchukwu, Salomey Osei,	709
657	Chris Emezue, NaijaVoices Community, Busayo	Atnafu Lambebo Tonja, Naome Etori, and Clinton	710
658	Awobade, Abraham Owodunni, Sewade Ogun, Han-	Mbataku. 2023. <a href="#">AfriSpeech-200: Pan-African ac-</a>	711
659	del Emezue, Gloria Monica Tobechukwu Emezue,	<a href="#">cented speech dataset for clinical and general domain</a>	712
660	Nefertiti Nneoma Emezue, Bunmi Akinremi, David	<a href="#">ASR</a> . <i>Transactions of the Association for Computa-</i>	713
661	Adelani, and Chris Pal. 2025. <a href="#">The naijavoices dataset:</a>	<i>tional Linguistics</i> , 11:1669–1685.	714
662	<a href="#">Cultivating large-scale, high-quality, culturally-rich</a>	Tobi Olatunji, Charles Nimo, Abraham Owodunni, Tas-	715
663	<a href="#">speech data for african languages</a> .	sallah Abdullahi, Emmanuel Ayodele, Mardhiyah	716
664	Daniel Galvez, Vladimir Bataev, Hainan Xu, and Tim	Sanni, Chinemelu Aka, Folafunmi Omofoye, Foutse	717
665	Kaldewey. 2024. Speed of light exact greedy de-	Yuehgoh, Timothy Faniran, Bonaventure F. P. Dos-	718
666	coding for rnn-t speech recognition models on gpu.	sou, Moshood Yekini, Jonas Kemp, Katherine Heller,	719
667	<i>arXiv preprint arXiv:2406.03791</i> .	Jude Chidubem Omeke, Chidi Asuzu MD, Naome A.	720
668	Sreyan Ghosh, Sonal Kumar, Ashish Seth, Chandra Ki-	Etori, Aimérou Ndiaye, Ifeoma Okoh, and 7 others.	721
669	ran Reddy Evuru, Utkarsh Tyagi, S. Sakshi, Oriol	2025. <a href="#">AfriMed-QA: A Pan-African, Multi-Specialty,</a>	722
670	Nieto, Ramani Duraiswami, and Dinesh Manocha.	<a href="#">Medical Question-Answering Benchmark Dataset.</a>	723
671	2024. <a href="#">GAMA: A Large Audio-Language Model with</a>	<i>arXiv preprint</i> . ArXiv:2411.15640 [cs].	724
672	<a href="#">Advanced Audio Understanding and Complex Rea-</a>	OpenAI, Aaron Hurst, Adam Lerer, Adam P. Goucher,	725
673	<a href="#">soning Abilities</a> . <i>arXiv preprint</i> . ArXiv:2406.11768	Adam Perelman, Aditya Ramesh, Aidan Clark, A. J.	726
674	[cs].	Ostrow, Akila Welihinda, Alan Hayes, Alec Radford,	727
675	Google DeepMind. 2024. <a href="#">Gemini 2.0 flash: Built for</a>	Aleksander Mądry, Alex Baker-Whitcomb, Alex Beu-	728
676	<a href="#">the agentic era</a> .	tel, Alex Borzunov, Alex Carney, Alex Chow, Alex	729
677	Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and	Kirillov, Alex Nichol, and 400 others. 2024. <a href="#">GPT-4o</a>	730
678	Sai Qian Zhang. 2024. Parameter-efficient fine-	<a href="#">System Card</a> . <i>arXiv preprint</i> . ArXiv:2410.21276	731
679	tuning for large models: A comprehensive survey.	[cs].	732
680	<i>arXiv preprint arXiv:2403.14608</i> .	Vassil Panayotov, Guoguo Chen, Daniel Povey, and San-	733
681	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam	jeev Khudanpur. 2015. <a href="#">Librispeech: An asr corpus</a>	734
682	Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow,	<a href="#">based on public domain audio books</a> . In <i>2015 IEEE</i>	735
683	Akila Welihinda, Alan Hayes, Alec Radford, and 1	<i>International Conference on Acoustics, Speech and</i>	736
684	others. 2024. Gpt-4o system card. <i>arXiv preprint</i>	<i>Signal Processing (ICASSP)</i> , pages 5206–5210.	737
685	<i>arXiv:2410.21276</i> .	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	738
686	Rabeeh Karimi Mahabadi, James Henderson, and Se-	Jing Zhu. 2002. Bleu: a method for automatic evalu-	739
687	bastian Ruder. 2021. Compacter: Efficient low-rank	ation of machine translation. In <i>Proceedings of the</i>	740
688	hypercomplex adapter layers. <i>Advances in Neural</i>	<i>40th annual meeting of the Association for Computa-</i>	741
689	<i>Information Processing Systems</i> , 34:1022–1035.	<i>tional Linguistics</i> , pages 311–318.	742
690	Josh Meyer, David Ifeoluwa Adelani, Edresson	Jing Peng, Yucheng Wang, Yu Xi, Xu Li, Xizhuo Zhang,	743
691	Casanova, Alp Öktem, Daniel Whitenack Julian We-	and Kai Yu. 2025. <a href="#">A Survey on Speech Large Lan-</a>	744
692	ber, Salomon Kabongo, Elizabeth Salesky, Iroko	<a href="#">guage Models</a> . <i>arXiv preprint</i> . ArXiv:2410.18908	745
693	Orife, Colin Leong, Perez Ogayo, and 1 others.	[eess] version: 3.	746
694	2022. Biblelets: a large, high-fidelity, multilingual,	Maja Popović. 2015. chrF: character n-gram f-score for	747
695	and uniquely african speech corpus. <i>arXiv preprint</i>	automatic mt evaluation. In <i>Proceedings of the tenth</i>	748
696	<i>arXiv:2207.03546</i> .	<i>workshop on statistical machine translation</i> , pages	749
		392–395.	750
		Vineel Pratap, Andros Tjandra, Bowen Shi, Paden	751
		Tomasello, Arun Babu, Sayani Kundu, Ali Elkahky,	752

753	Zhaoheng Ni, Apoorv Vyas, Maryam Fazel-Zarandi,	Changhan Wang, Juan Pino, Anne Wu, and Jiatao Gu.	810
754	Alexei Baevski, Yossi Adi, Xiaohui Zhang, Wei-Ning	2020. Covost: A diverse multilingual speech-to-text	811
755	Hsu, Alexis Conneau, and Michael Auli. 2024. <a href="#">Scal-</a>	translation corpus. <i>arXiv preprint arXiv:2002.01320</i> .	812
756	<a href="#">ing speech technology to 1,000+ languages</a> . <i>Journal</i>		
757	<i>of Machine Learning Research</i> , 25(97):1–52.		
758	Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel	Jiayi Wang, David Ifeoluwa Adelani, Sweta Agrawal,	813
759	Synnaeve, and Ronan Collobert. 2020a. <a href="#">MLS: A</a>	Marek Masiak, Ricardo Rei, Eleftheria Briakou, Ma-	814
760	<a href="#">Large-Scale Multilingual Dataset for Speech Re-</a>	rine Carpuat, Xuanli He, Sofia Bourhim, Andiswa	815
761	<a href="#">search</a> . ArXiv:2012.03411 [eess].	Bukula, and 1 others. 2023. Afrimte and africommet:	816
762		Enhancing comet to embrace under-resourced african	817
763	Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel	languages. <i>arXiv preprint arXiv:2311.09828</i> .	818
764	Synnaeve, and Ronan Collobert. 2020b. <a href="#">Mls: A</a>		
765	<a href="#">large-scale multilingual dataset for speech research</a> .	Wenqi Yu, Siyin Wang, Xiaoyu Yang, Xianzhao Chen,	819
766	In <i>Interspeech 2020</i> , pages 2757–2761.	Xiaohai Tian, Jun Zhang, Guangzhi Sun, Lu Lu, Yux-	820
767		uan Wang, and Chao Zhang. Salmonn-omni: A	821
768	Krishna C. Puvvada, Piotr Żelasko, He Huang, Olek-	speech understanding and generation llm in a codec-	822
769	sii Hrinchuk, Nithin Rao Koluguri, Kunal Dhawan,	free full-duplex framework.	823
770	Somshubra Majumdar, Elena Rastorgueva, Zhehuai		
771	Chen, Vitaly Lavrukhin, Jagadeesh Balam, and Boris	<b>A Appendix</b>	824
772	Ginsburg. 2024. <a href="#">Less is more: Accurate speech</a>	<b>A.1 Automatic Speech Recognition</b>	825
773	<a href="#">recognition &amp; translation without web-scale data</a> . In	<b>A.1.1 ASR Prompts</b>	826
774	<i>Interspeech 2024</i> , pages 3964–3968.	For automatic speech recognition (ASR), we eval-	827
775	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	uate three prompting strategies. The first employs	828
776	man, Christine McLeavey, and Ilya Sutskever. 2023a.	a simple instruction: “Transcribe this audio.” The	829
777	Robust speech recognition via large-scale weak su-	second includes language specificity: “Transcribe	830
778	per- In <i>International conference on machine</i>	the entire audio in {source_language}.” The third is	831
779	learning, pages 28492–28518. PMLR.	a few-shot variant of the second prompt, which pro-	832
780		vides two audio-transcription exemplars as demon-	833
781	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	strations to guide the model’s output.	834
782	man, Christine McLeavey, and Ilya Sutskever. 2023b.		
783	Robust speech recognition via large-scale weak super-	<b>A.2 Automatic Speech Translation</b>	835
784	vision. In <i>Proceedings of the 40th International Con-</i>	<b>A.2.1 AST Prompting Strategies</b>	836
785	<i>ference on Machine Learning</i> , ICML’23. JMLR.org.	We evaluate three AST prompting strategies:	837
786	Machel Reid, Junjie Hu, Graham Neubig, and Yu-		
787	taka Matsuo. 2021. <a href="#">AfroMT: Pretraining Strate-</a>	1. <b>Zero-shot translation:</b>	838
788	<a href="#">gies and Reproducible Benchmarks for Transla-</a>	“Given audio in {source_language}, trans-	839
789	<a href="#">tion of 8 African Languages</a> . <i>arXiv preprint</i> .	late to English.”	840
790	ArXiv:2109.04715 [cs].		
791	Paul K Rubenstein, Chulayuth Asawaroengchai,	2. <b>Zero-shot transcriptiontranslation:</b>	841
792	Duc Dung Nguyen, Ankur Bapna, Zalán Borsos,	“Given audio in {source_language}, first	842
793	Félix de Chaumont Quirry, Peter Chen, Dalia El	transcribe the speech, then translate the tran-	843
794	Badawy, Wei Han, Eugene Kharitonov, and 1 others.	script into English.”	844
795	2023. Audiopalm: A large language model that can		
796	speak and listen. <i>arXiv preprint arXiv:2306.12925</i> .	3. <b>Few-shot variants:</b>	845
797		For each of the above prompts, we prepend	846
798	Mardhiyah Sanni, Tassallah Abdullahi, Devendra D.	two example audio–translation pairs to pro-	847
799	Kayande, Emmanuel Ayodele, Naome A. Etori,	vide in-context demonstrations of the desired	848
800	Michael S. Mollel, Moshood Yekini, Chibuzor	behavior.	849
801	Okocha, Lukman E. Ismaila, Folafunmi Omofoye,		
802	Boluwatife A. Adewale, and Tobi Olatunji. 2025.	We found the Zero-shot transcriptiontranslation	850
803	<a href="#">Afrispeech-Dialog: A Benchmark Dataset for Spon-</a>	gives the best result as it encourages the model	851
804	<a href="#">taneous English Conversations in Healthcare and Be-</a>	to understand the audio by first transcribing, before	852
805	<a href="#">yond</a> . <i>arXiv preprint</i> . ArXiv:2502.03945 [cs].	attempting to translate.	853
806			
807	Jiatong Shi, Shih-Heng Wang, William Chen, Mar-	<b>A.2.2 Performance Across Sources</b>	854
808	tijn Bartelds, Vanya Bannihatti Kumar, Jinchuan		
809	Tian, Xuankai Chang, Dan Jurafsky, Karen Livescu,		
	Hung-yi Lee, and Shinji Watanabe. 2024. <a href="#">ML-</a>		
	<a href="#">SUPERB 2.0: Benchmarking Multilingual Speech</a>		
	<a href="#">Models Across Modeling Constraints, Languages,</a>		
	<a href="#">and Datasets</a> . <i>arXiv preprint</i> . ArXiv:2406.08641		
	[cs].		

Language	Whisper medium	Whisper large-v3	MMS-1b all	SeamlessM4T-v2 Large	Gpt-4o audio-preview	Gemini-2.0 flash
Afrikaans	44.49	30.93	26.48	18.64	32.20	<b>13.77</b>
Amharic	441.81	205.81	34.71	86.45	118.45	<b>19.10</b>
Arabic	–	11.06	36.28	9.29	6.64	<b>4.42</b>
Fulani	–	–	<b>56.78</b>	–	157.03	74.62
Hausa	158.21	86.13	<b>31.39</b>	–	100.85	34.92
Igbo	–	–	<b>44.60</b>	102.95	110.63	66.07
Luganda	–	–	45.77	<b>37.62</b>	89.34	52.98
Pedi	–	–	<b>31.29</b>	–	110.12	90.11
Shona	222.30	116.51	<b>29.60</b>	76.46	97.43	54.45
Swahili	99.04	41.51	22.22	<b>11.98</b>	29.92	12.37
Xhosa	–	–	<b>44.58</b>	–	124.79	56.94
Yoruba	204.21	87.18	34.29	<b>31.03</b>	82.98	42.04
Zulu	–	–	40.30	50.56	110.88	<b>32.03</b>

Table 9: **WER % for each model–language pair on the FLEURS subset of the Multilingual African Speech transcription dataset**; the lowest (best) WER per language is shown in bold. "-" indicates the language is not supported by the model.

Language	Canary 1b	Whisper medium	Whisper large-v3	MMS-1b all	Qwen2.5	SeamlessM4T-v2 Large	GPT-4o audio-preview	Gemini-2.0 flash
Afrikaans	–	52.54	32.7	36.99	–	<b>14.69</b>	47.48	17.64
Akan	–	–	–	<b>62.90</b>	–	–	103.97	76.53
Arabic	–	45.74	33.10	75.29	–	–	32.76	<b>23.67</b>
French	13.14	16.32	10.65	41.74	24.00	16.80	12.11	<b>8.02</b>
Hausa	–	129.55	93.68	43.22	–	–	125.96	<b>39.55</b>
Igbo	–	–	–	<b>53.61</b>	–	68.97	104.18	77.30
Kinyarwanda	–	–	–	<b>46.65</b>	–	–	134.26	65.19
Pedi	–	–	–	<b>46.67</b>	–	–	124.27	76.72
Sesotho	–	–	–	–	–	–	172.76	<b>77.59</b>
Shona	–	150.31	101.27	<b>32.33</b>	–	75.46	80.30	45.04
Swahili	–	112.09	48.11	34.17	–	18.87	42.74	<b>16.30</b>
Tswana	–	–	–	–	–	–	135.98	<b>72.81</b>
Twi	–	–	–	<b>50.55</b>	–	–	102.58	80.66
Xhosa	–	–	–	<b>43.62</b>	–	–	122.86	46.54
Yoruba	–	157.12	88.98	43.05	–	<b>30.44</b>	134.79	54.02
Zulu	–	–	–	48.41	–	52.49	129.38	<b>35.19</b>

Table 10: **WER % for each model–language pair on the Intron-AfriVox subset of the Multilingual African Speech transcription dataset**; the lowest (best) WER per language is shown in bold. "-" indicates the language is not supported by the model.

Language	Whisper medium	Whisper large-v3	MMS-1b all	Seamless-M4T-v2 Large	Gpt-4o audio-preview	Gemini-2.0 flash
Amharic	427.57	155.51	76.16	<b>23.94</b>	280.17	280.17
Swahili	132.67	73.47	40.56	<b>26.39</b>	93.58	93.58

Table 11: **WER % for each model–language pair on the ALFFA subset of the Multilingual African Speech transcription dataset**; the lowest (best) WER per language is shown in bold. "-" indicates the language is not supported by the model.



Language	Canary 1b	Whisper medium	Whisper large-v3	MMS-1b all	Qwen	Seamless-M4T-v2 Large	Gpt-4o audio-preview	Gemini-2.0 flash
French	<b>5.49</b>	7.69	11.10	24.53	24.00	14.82	34.55	12.67

Table 12: **WER % for each model–language pair on the OpenSLR subset of the Multilingual African Speech transcription dataset**; the lowest (best) WER per language is shown in bold. "-" indicates the language is not supported by the model.

Language	MMS-1b all	Gpt-4o audio-preview	Gemini-2.0 flash
Akan	<b>77.78</b>	133.33	94.44
Ga	–	172.73	<b>114.55</b>
Twi	<b>75.00</b>	184.38	150.00

Table 13: **WER % for each model–language pair on the Ashesi Financial Inclusion subset of the Multilingual African Speech transcription dataset**; the lowest (best) WER per language is shown in bold. "-" indicates the language is not supported by the model.

Language	Whisper medium	Whisper large-v3	MMS-1b all	Seamless-M4T-v2 Large	Gpt-4o- audio-preview	Gemini-2.0 flash
Afrikaans	52.30	37.65	27.09	<b>13.80</b>	57.07	17.55
Amharic	513.92	183.28	<b>52.69</b>	92.51	183.54	130.17
Arabic	36.24	18.33	27.66	68.27	31.73	<b>11.94</b>
Hausa	270.36	91.49	<b>27.20</b>	–	109.09	40.53
Igbo	–	–	60.71	<b>42.86</b>	246.43	82.14
Kinyarwanda	–	–	<b>32.75</b>	–	136.35	84.26
Luganda	–	–	28.51	<b>15.97</b>	132.04	80.73
Swahili	120.74	71.30	24.50	<b>14.11</b>	92.47	26.33
Twi	–	–	<b>57.53</b>	–	123.29	93.15
Yoruba	294.01	99.43	<b>38.63</b>	39.91	96.48	103.57

Table 14: **WER % for each model–language pair on the Common Voice subset of the Multilingual African Speech transcription dataset**; the lowest (best) WER per language is shown in bold. "-" indicates the language is not supported by the model.

Language	Whisper- medium	Whisper large-v3	MMS-1b all	Seamless-M4T-v2 Large	Gpt-4o audio-preview	Gemini-2.0 flash
Afrikaans	99.00	68.31	71.32	<b>25.01</b>	151.50	48.94
Pedi	–	–	<b>42.03</b>	–	119.29	90.75
Sesotho	–	–	–	–	133.43	<b>104.33</b>
Tswana	–	–	–	–	127.82	<b>85.19</b>
Xhosa	–	–	<b>31.93</b>	–	171.43	56.70
Zulu	–	–	<b>28.10</b>	56.43	208.26	44.64

Table 15: **WER % for each model–language pair on the NCHLT subset of the Multilingual African Speech transcription dataset**; the lowest (best) WER per language is shown in bold. "-" indicates the language is not supported by the model.

Language	Whisper medium	Whisper large-v3	MMS-1b all	Seamless-M4T-v2 Large	Gpt-4o audio-preview	Gemini-2.0 flash
Hausa	112.01	102.16	<b>39.37</b>	–	110.46	104.58
Twī	–	–	–	<b>51.53</b>	89.81	78.04
Yoruba	118.50	106.66	<b>24.63</b>	27.23	84.70	44.94

Table 16: **WER % for each model–language pair on the BibleTTS subset of the Multilingual African Speech transcription dataset**; the lowest (best) WER per language is shown in bold. "-" indicates the language is not supported by the model.

Language	Metric	Fluency $r$	Adequacy $r$
<b>Akan</b>	BLEU	−0.09	0.58
	ChrF	−0.24	<b>0.68</b>
	AfriComet-STL	<b>0.07</b>	0.61
<b>Igbo</b>	BLEU	<b>0.10</b>	0.63
	ChrF	−0.11	0.69
	AfriComet-STL	−0.04	<b>0.93</b>
<b>Pedi</b>	BLEU	0.05	<b>0.78</b>
	ChrF	0.26	0.68
	AfriComet-STL	<b>0.38</b>	0.61
<b>Shona</b>	BLEU	0.38	0.44
	ChrF	0.48	0.73
	AfriComet-STL	<b>0.67</b>	<b>0.86</b>
<b>Swahili</b>	BLEU	0.43	0.47
	ChrF	0.56	0.70
	AfriComet-STL	<b>0.67</b>	<b>0.76</b>
<b>Twī</b>	BLEU	0.43	0.34
	ChrF	0.44	0.36
	AfriComet-STL	<b>0.52</b>	<b>0.60</b>
<b>Yoruba</b>	BLEU	0.30	0.61
	ChrF	0.40	<b>0.76</b>
	AfriComet-STL	<b>0.47</b>	0.70
<b>Average</b>	BLEU	0.23	0.52
	ChrF	0.40	0.66
	AfriComet-STL	<b>0.48</b>	<b>0.70</b>

Table 17: Pearson correlations ( $r$ ) between automatic metrics and human evaluations of fluency and adequacy for automatic speech translation.

Language	Canary 1b	Whisper medium	Whisper large-v3	Qwen2.5	SeamlessM4T Large-v2	Gpt-4o audio-preview	Gemini-2.0 flash
Afrikaans	–	19.39	23.2	–	27.62	31.59	<b>38.76</b>
Akan	–	–	–	–	–	2.44	<b>5.15</b>
Amharic	–	0.8	0.71	–	15.61	4.2	<b>24.88</b>
Arabic	–	17.97	20.34	–	27.69	31.06	<b>34.68</b>
French	24.46	27.39	28.92	41.40	33.38	41.27	<b>43.57</b>
Fulani	–	–	–	–	0.58	1.05	<b>2.41</b>
Ga	–	–	–	–	–	0.49	<b>1.06</b>
Hausa	–	0.71	0.71	–	0.31	6.23	<b>21.06</b>
Igbo	–	–	–	–	1.92	2.97	<b>5.82</b>
Kinyarwanda	–	–	–	–	–	1.99	<b>10.91</b>
Luganda	–	–	–	–	<b>15.97</b>	7.77	13.79
Pedi	–	–	–	–	–	3.19	<b>6.34</b>
Sesotho	–	–	–	–	–	4.11	<b>11.23</b>
Shona	–	0.4	0.52	–	2.11	6.78	<b>12.56</b>
Swahili	–	2.84	5.47	–	23.27	26.78	<b>32.62</b>
Tswana	–	–	–	–	–	3.72	<b>9.59</b>
Twi	–	–	–	–	–	<b>2.83</b>	2.48
Xhosa	–	–	–	–	–	4.71	<b>19.9</b>
Yoruba	–	0.24	0.37	–	<b>14.39</b>	4.89	11.77
Zulu	–	–	–	–	8.17	6.57	<b>22.9</b>

Table 18: **BLEU scores for each model–language pair on the Multilingual African Speech translation dataset;** the highest (best) BLEU score per language is shown in bold. "-" indicates the language is not supported by the model.

Language	Gemini-2.0 flash	GPT-4o audio-preview	SeamlessM4T-v2 Large	Whisper Large	Whisper Medium	Canary-1b	Qwen2.5
Afrikaans	<b>64.33</b>	56.39	56.13	50.33	45.58	–	–
Akan	<b>29.86</b>	25.01	–	–	–	–	–
Amharic	<b>56.62</b>	29.62	43.48	17.06	13.57	–	–
Arabic	<b>63.10</b>	59.26	55.53	47.85	44.38	–	–
French	<b>66.56</b>	64.40	63.72	58.61	57.19	54.12	64.94*
Fulani	<b>27.56</b>	23.82	16.25	–	–	–	–
Ga	<b>20.08</b>	19.09	–	–	–	–	–
Hausa	<b>48.48</b>	29.81	13.47	13.29	7.78	–	–
Igbo	<b>32.10</b>	25.40	18.52	–	–	–	–
Kinyarwanda	<b>37.69</b>	23.62	–	–	–	–	–
Luganda	<b>44.23</b>	35.56	44.21	–	–	–	–
Pedi	<b>34.63</b>	27.51	–	–	–	–	–
Sesotho	<b>38.00</b>	26.71	–	–	–	–	–
Shona	<b>42.07</b>	33.56	21.65	15.59	12.76	–	–
Swahili	<b>61.74</b>	55.90	53.39	30.00	22.13	–	–
Tswana	<b>35.52</b>	25.11	–	–	–	–	–
Twi	<b>24.22</b>	23.15	–	–	–	–	–
Xhosa	<b>48.82</b>	28.54	–	–	–	–	–
Yoruba	38.45	28.37	<b>40.53</b>	14.29	10.45	–	–
Zulu	<b>52.76</b>	31.54	32.79	–	–	–	–

Table 19: **CHrF scores for each model–language pair on the Multilingual African Speech translation dataset;** the highest (best) CHrF score per language is shown in bold. "-" indicates the language is not supported by the model.

Language	Gemini-2.0 flash		GPT-4o audio-preview		SeamlessM4T-v2 Large		Whisper Large		Whisper Medium	
	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF
Amharic	<b>29.44</b>	<b><u>62.09</u></b>	5.60	33.25	21.24	50.16	1.20	19.06	1.08	16.30
Arabic	<b>33.25</b>	<b><u>66.44</u></b>	30.66	63.85	33.86	62.88	18.83	50.45	18.07	48.54
Fulani	<b>2.41</b>	<b><u>27.56</u></b>	1.05	23.82	0.58	16.25	–	–	–	–
Hausa	<b>17.68</b>	<b><u>50.09</u></b>	6.07	34.25	0.48	16.79	0.16	15.18	0.22	10.13
Igbo	<b>5.54</b>	<b><u>34.91</u></b>	2.48	27.37	1.17	17.99	–	–	–	–
Luganda	<b>13.79</b>	<b><u>44.23</u></b>	7.77	35.56	15.97	44.21	–	–	–	–
Pedi	<b>6.30</b>	<b><u>36.41</u></b>	2.95	28.84	–	–	–	–	–	–
Shona	<b>12.20</b>	<b><u>43.54</u></b>	6.15	34.43	2.67	25.44	0.79	17.46	0.55	14.62
Swahili	<b>30.70</b>	<b><u>62.10</u></b>	23.89	55.24	28.41	57.03	4.48	29.04	2.54	20.40
Xhosa	<b>20.09</b>	<b><u>51.51</u></b>	4.19	29.77	–	–	–	–	–	–
Yoruba	10.21	40.15	4.23	30.70	<b>13.25</b>	<b><u>41.04</u></b>	0.62	16.73	0.41	12.20
Zulu	<b>21.54</b>	<b><u>53.45</u></b>	5.86	33.00	7.67	34.19	–	–	–	–

Table 20: **BLEU & CHrF scores for each model–language pair on the FLEURS subset of the Multilingual African Speech translation dataset**; the highest (best) BLEU & CHrF score per language is shown in bold with the CHrF score further underlined. "-" indicates the language is not supported by the model.

Language	Gemini-2.0 flash		GPT-4o audio-preview		SeamlessM4T-v2 Large		Whisper Large		Whisper Medium	
	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF
Afrikaans	<b>38.76</b>	<b><u>64.33</u></b>	31.59	56.39	27.62	56.13	23.20	50.33	19.39	45.58
Akan	<b>5.15</b>	<b><u>29.86</u></b>	2.44	25.01	–	–	–	–	–	–
Amharic	<b>16.45</b>	<b><u>45.29</u></b>	1.39	22.12	6.07	29.50	0.12	13.29	0.31	7.98
Arabic	<b>24.75</b>	<b><u>55.28</u></b>	21.98	52.07	15.99	44.95	13.55	41.54	10.78	36.94
French	<b>32.49</b>	<b><u>60.96</u></b>	28.99	57.45	20.07	50.06	23.95	53.37	21.31	51.01
Ga	<b>1.06</b>	<b><u>20.08</u></b>	0.49	19.09	–	–	–	–	–	–
Hausa	<b>23.18</b>	<b><u>48.70</u></b>	6.48	28.76	0.19	11.88	0.16	12.52	0.15	6.34
Igbo	<b>5.69</b>	<b><u>29.50</u></b>	2.99	23.62	2.05	17.18	–	–	–	–
Kinyarwanda	<b>10.91</b>	<b><u>37.69</u></b>	1.99	23.62	–	–	–	–	–	–
Pedi	<b>6.40</b>	<b><u>31.04</u></b>	3.61	24.81	–	–	–	–	–	–
Sesotho	<b>11.23</b>	<b><u>38.00</u></b>	4.11	26.71	–	–	–	–	–	–
Shona	<b>12.98</b>	<b><u>40.15</u></b>	7.55	32.42	1.15	16.26	0.23	13.34	0.25	10.40
Swahili	<b>30.45</b>	<b><u>58.71</u></b>	23.52	51.43	19.82	49.07	6.51	30.33	4.00	21.80
Tswana	<b>9.59</b>	<b><u>35.52</u></b>	3.72	25.11	–	–	–	–	–	–
Twi	<b>2.48</b>	<b><u>24.22</u></b>	2.83	23.15	–	–	–	–	–	–
Xhosa	<b>19.76</b>	<b><u>46.48</u></b>	5.11	27.47	–	–	–	–	–	–
Yoruba	<b>14.37</b>	39.68	5.61	27.77	14.01	<b><u>40.44</u></b>	0.11	12.72	0.08	8.35
Zulu	<b>24.01</b>	<b><u>52.14</u></b>	7.17	30.20	8.60	31.48	–	–	–	–

Table 21: **BLEU & CHrF scores for each model–language pair on the Intron-AfriVox subset of the Multilingual African Speech translation dataset**; the highest (best) BLEU & CHrF score per language is shown in bold with the CHrF score further underlined. "-" indicates the language is not supported by the model.



Language	Canary1b		Qwen2.5	
	BLEU	ChrF	BLEU	ChrF
French	13.78	44.46	<b>41.40</b>	<b><u>64.94</u></b>

Table 21: **BLEU & ChrF scores for each model–language pair on the Intron-AfriVox subset of the Multilingual African Speech translation dataset**; the highest (best) BLEU & ChrF score per language is shown in bold with the ChrF score further underlined. "-" indicates the language is not supported by the model.

Language	Gemini		GPT-4o-audio preview		SeamlessM4T v2 Large		Whisper Large		Whisper Medium		Qwen Omni	
	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF
Hausa	<b>19.15</b>	<b><u>44.84</u></b>	5.61	25.34	0.17	12.69	0.17	12.52	0.11	8.29	0.25	13.19
Igbo	<b>6.97</b>	<b><u>28.67</u></b>	4.35	22.91	4.22	22.80	–	–	–	–	0.26	12.59
Yoruba	9.92	32.57	4.88	24.32	<b>16.34</b>	<b><u>39.61</u></b>	0.11	11.52	0.11	10.33	0.24	13.12

Table 22: **BLEU and ChrF scores for each model–language pair on the NaijaVoices subset of the Multilingual African Speech Translation dataset**. The highest (best) BLEU and ChrF score per language is shown in bold, with the ChrF score further underlined. "–" indicates the language is not supported by the model.

Language	Gemini-2.0 flash		GPT-4o audio-preview		SeamlessM4T-v2 Large		Whisper Large		Whisper Medium	
	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF
Swahili	<b>37.22</b>	<b><u>65.60</u></b>	33.74	62.25	25.15	57.15	4.32	30.09	1.68	23.38

Table 23: **BLEU & ChrF scores for each model–language pair on the IWSLT\_LRST subset of the Multilingual African Speech translation dataset**; the highest (best) BLEU & ChrF score per language is shown in bold with the ChrF score further underlined. "-" indicates the language is not supported by the model.

Language	Gemini-2.0 flash		GPT-4o audio-preview		SeamlessM4T-v2 Large		Whisper Large		Whisper Medium	
	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF	BLEU	ChrF
Arabic	<b>51.72</b>	<b><u>70.78</u></b>	45.97	64.50	37.07	62.11	30.92	54.18	28.03	50.48
French	<b>44.40</b>	<b><u>66.91</u></b>	42.19	64.83	34.35	64.56	29.32	58.98	27.84	57.57

Table 24: **BLEU & ChrF scores for each model–language pair on the Covost subset of the Multilingual African Speech translation dataset**; the highest (best) BLEU & ChrF score per language is shown in bold with the ChrF score further underlined. "-" indicates the language is not supported by the model.

Language	Canary-1b		QWEN	
	BLEU	ChrF	BLEU	ChrF
French	25.03	54.72	<b>41.40</b>	<b><u>64.94</u></b>

Table 24: **BLEU & ChrF scores for each model–language pair on the Covost subset of the Multilingual African Speech translation dataset**; the highest (best) BLEU & ChrF score per language is shown in bold with the ChrF score further underlined. "-" indicates the language is not supported by the model.