

Benchmarking Self-Supervised Speech Models on Multilingual Nigerian Speech

Omotayo Omoyemi
University of Derby, UK
etayo25@gmail.com

Ifeoluwa Oladeni
National Open University of Nigeria
oladeniifeoluwa123@gmail.com

1 Abstract

2 Self-supervised speech models such as Whisper and wav2vec
3 2.0 have significantly advanced automatic speech recognition
4 (ASR) performance for high-resource languages. However,
5 their robustness and generalization to underrepresented Afri-
6 can languages remain insufficiently studied.

7 In this work, we present a systematic benchmark of modern
8 self-supervised ASR models on a multilingual Nigerian
9 speech corpus comprising English, Hausa, Igbo, and Yoruba.
10 Using the Nigerian Common Voice dataset (158 hours), we
11 evaluate zero-shot performance of pretrained models and
12 compare it with supervised adaptation using fine-tuning of
13 multilingual speech encoders. We report Word Error Rate
14 (WER) and Character Error Rate (CER) across languages and
15 analyze the effect of supervised adaptation and cross-lan-
16 guage transfer.

17 Our results show that zero-shot ASR performance is substan-
18 tially degraded for Nigerian languages compared to widely
19 represented benchmark languages. Supervised fine-tuning
20 consistently improves recognition accuracy, although the
21 magnitude of improvement varies across languages and de-
22 pends on the compatibility between the pretrained checkpoint
23 and the target language. In particular, adaptation from a
24 Hausa-pretrained XLS-R model yields strong gains for Hausa
25 but more limited improvements for Igbo, highlighting the im-
26 portance of language-specific training data.

27 These findings demonstrate that multilingual pretraining
28 alone is insufficient for reliable ASR in underrepresented Afri-
29 can languages and that supervised adaptation remains nec-
30 essary for robust deployment. The study provides reproduci-
31 ble benchmarks for multilingual ASR evaluation in African
32 contexts and offers practical guidance for adapting large-
33 scale speech models to underrepresented languages.

34

35 1 Introduction

36 Self-supervised and weakly supervised pretraining has be-
37 come a dominant paradigm in automatic speech recognition
38 (ASR), enabling strong performance with limited labeled
39 data and improving robustness across domains.

40 Representative examples include wav2vec 2.0, which learns
41 speech representations from raw audio and achieves compet-
42 itive word error rates after fine-tuning [Baevski et al., 2020],
43 and Whisper, which scales weakly supervised training to hun-
44 dreds of thousands of hours of web data and demonstrates
45 broad generalization across benchmarks [Radford et al.,
46 2022]. Cross-lingual pretraining further extends these gains
47 to multilingual settings, as shown by large-scale wav2vec-
48 style models trained on hundreds of thousands of hours span-
49 ning many languages [Babu et al., 2021; Conneau et al.,
50 2021].

51 Despite this progress, recent work shows that generalization
52 remains uneven for underrepresented languages, accents, and
53 demographic groups. Performance disparities in ASR are of-
54 ten linked to imbalanced data coverage and limited linguistic
55 diversity, with surveys highlighting persistent challenges for
56 underrepresented African languages [Nakatumba Nabende et
57 al., 2025]. In addition, accent and language labels in widely
58 used speech datasets are often self-reported and inconsis-
59 tently defined, which can complicate both training and evalu-
60 ation [Reid and Williams, 2023]. These issues motivate sys-
61 tematic benchmarking to determine where modern pretrained
62 ASR models succeed, where they fail, and how much adap-
63 tation is required for reliable deployment.

64 Nigeria provides a particularly relevant testbed for this anal-
65 ysis. It is linguistically diverse and includes several widely
66 spoken languages, such as English, Hausa, Igbo, and Yoruba,
67 that remain underrepresented in standard ASR benchmarks.
68 At the same time, recent work has begun to release African-
69 accented speech resources, including AfriSpeech-200 for
70 Pan-African accented English [Olatunji et al., 2023] and Af-
71 rispeech-Dialog for conversational African-accented English
72 [Sanni et al., 2025], highlighting the need for broader multi-
73 lingual evaluations that include multiple African languages
74 rather than English alone.

75 In this work, we benchmark modern pretrained ASR models
76 on a multilingual Nigerian speech corpus with official
77 train/dev/test splits. We evaluate (i) zero-shot recognition us-
78 ing a large pretrained model, and (ii) supervised adaptation
79 using fine-tuning of an XLS-R checkpoint on individual lan-
80 guages. We report Word Error Rate (WER) and Character Er-
81 ror Rate (CER) across languages and compare performance

82 under zero-shot and fine-tuned settings. Our results show that
 83 zero-shot multi-lingual ASR performs poorly for Nigerian
 84 languages, while supervised adaptation substantially reduces
 85 error rates, although the magnitude of improvement varies
 86 across languages. This study provides empirical evidence on
 87 the limits of multilingual pretraining and the importance of
 88 language-specific adaptation for underrepresented African
 89 ASR, using reproducible evaluation protocols based on
 90 widely used dataset tooling [Lhoest et al., 2021] and multi-
 91 lingual speech corpus design practices inspired by Common
 92 Voice [Ardila et al., 2020].

93 2 Related Work

94 2.1 Multilingual and Large-Scale ASR Bench- 95 marks

96 The evaluation of multilingual ASR systems has expanded
 97 beyond high-resource settings toward more inclusive and
 98 globally representative benchmarks. The ML-SUPERB 2.0
 99 Challenge introduces a large-scale multilingual evaluation
 100 suite covering over 200 language varieties and demonstrates
 101 that even state-of-the-art pretrained models exhibit signifi-
 102 cant performance variability across languages and dialects
 103 (Chen et al., 2025). These findings reinforce that zero-shot
 104 robustness cannot be assumed for underrepresented speech
 105 communities.

106 Complementing this direction, recent industrial research has
 107 released African-focused ASR benchmarks such as
 108 PazaBench, which provides standardized evaluation re-
 109 sources for 39 African languages and explicitly targets de-
 110 ployment in underrepresented settings (Muchai et al., 2026).
 111 These efforts collectively highlight the need for systematic,
 112 reproducible benchmarking of pretrained speech models in
 113 linguistically diverse African contexts.

114 2.2 Efficient Adaptation for Low-Resource Lan- 115 guages

116 Beyond benchmarking, methodological work has examined
 117 how to adapt large speech models to underrepresented lan-
 118 guages efficiently. Parameter-efficient fine-tuning ap-
 119 proaches, including adapter-based methods, have demon-
 120 strated competitive performance in extremely underrepre-
 121 sented scenarios without full model retraining (Mainzinger,
 122 2024). Similarly, efficient ASR training frameworks tailored
 123 for underrepresented languages have been proposed to reduce
 124 data and computational requirements while preserving accu-
 125 racy (Bandarupalli, 2025).

126 These studies suggest that adaptation strategy plays a critical
 127 role in determining whether large pretrained models can
 128 meaningfully generalize to underrepresented African lan-
 129 guages.

131 2.3 Accent Robustness and Fairness in ASR

132 A growing body of work has examined disparities in ASR
 133 performance across accents and demographic groups. Ac-
 134 cent-invariant modeling approaches, such as saliency-driven
 135 spectrogram masking, have been proposed to enhance robust-
 136 ness to dialectal variation and reduce word error rate gaps be-
 137 tween accent groups (Sameti et al., 2025). In parallel, fair-
 138 ness-aware fine-tuning techniques have shown that multilin-
 139 gual ASR models can be optimized to mitigate performance
 140 disparities across demographic attributes without substantial
 141 degradation in overall accuracy (Swain et al., 2024).

142 Systematic reviews focused specifically on African un-
 143 derrepresented languages further document persistent gaps in
 144 dataset coverage, evaluation consistency, and bias mitigation
 145 strategies, emphasizing the importance of transparent bench-
 146 marking practices (Imam et al., 2025).

148 2.4 Implications for African Multilingual ASR

149 Collectively, recent benchmarking initiatives, adaptation
 150 techniques, and fairness-aware modeling strategies under-
 151 score a central insight: while large pretrained speech models
 152 have dramatically advanced ASR, their performance in mul-
 153 tilingual African contexts remains uneven and insufficiently
 154 characterized. The literature increasingly calls for reproduc-
 155 ible, language-specific evaluation frameworks grounded in re-
 156 alistic deployment settings (Chen et al., 2025; Muchai et al.,
 157 2026).

158 In this work, we respond to this gap by conducting a system-
 159 atic evaluation of modern pretrained ASR models on multi-
 160 lingual Nigerian speech, analyzing zero-shot generalization,
 161 supervised adaptation, and cross-language transfer across
 162 four languages.

163 3 Dataset

Language	#Ut- ter- ances	Hours	Mean Dur (s)	Me- dian Dur (s)	P95 Dur (s)	Mean Chars
Eng- lish	3,402	11.08	11.72	6.7	39.0	127.3
Hausa	9,008	10.76	4.30	4.1	6.7	68.4
Igbo	5,714	8.60	5.42	5.2	9.1	45.8
Yo- ruba	4,171	6.82	5.88	5.8	8.3	57.6

164 **Table 1:** Summary Statistics of the Nigerian Multilingual Speech
 165 Corpus (All Splits Combined)

166 We conduct our experiments using a multilingual Nigerian speech
 167 corpus comprising English, Hausa, Igbo, and Yoruba. The dataset
 168 provides official train, validation, and test splits for each language.
 169 Statistics reported in Table 1 are computed across all splits com-
 170 bined.

171 The corpus contains a total of approximately 37.25 hours of tran-
 172 scribed speech distributed unevenly across languages. Hausa con-
 173

174 tains the largest number of utterances (9,008), while Yoruba con- 229
175 tains the smallest total duration (6.82 hours). English exhibits nota- 230
176 bly longer utterances on average (mean duration 11.72 seconds) 231
177 compared to the other languages (4–6 seconds). The 95th percentile 232
178 duration for English (39.0 seconds) is substantially higher than for 233
179 Hausa, Igbo, and Yoruba (all below 10 seconds), indicating struc- 234
180 tural differences in utterance length distributions across languages. 235
181 All audio recordings were resampled to 16 kHz to ensure compati- 236
182 bility with pretrained ASR models and to maintain consistent acous- 237
183 tic representation across languages. Following the dataset documen- 238
184 tation, transcripts were normalized by removing surrounding quota- 239
185 tion marks and ensuring sentence-final punctuation. No additional 240
186 text cleaning or filtering was performed. 241
187 Table 1 also reports transcript length statistics. English transcripts 242
188 are considerably longer on average (127.3 characters) than the other 243
189 languages, with Igbo exhibiting the shortest mean transcript length 244
190 (45.8 characters). These differences in duration and transcript length 245
191 are expected to influence recognition difficulty and are considered 246
192 in the interpretation of experimental results. 247
193 The dataset version used in this study was frozen after preprocessing 248
194 to ensure reproducibility. 249
195

196 4 Method

197 4.1 Model Selection

198 We evaluate automatic speech recognition performance using 253
199 Whisper-small as a zero-shot baseline. Whisper is a large- 254
200 scale encoder–decoder model trained on multilingual 255
201 speech–text data using weak supervision. We select the 256
202 “small” variant to balance computational feasibility on CPU 257
203 hardware and model capacity. For subsequent supervised ex- 258
204 periments, we fine-tune multilingual speech encoders (e.g., 259
205 XLS-R) using a CTC objective. 260

206 4.2 Zero-Shot Inference Protocol

207 Zero-shot evaluation is performed using the faster-whisper 262
208 implementation with int8 CPU inference. Audio inputs are 263
209 resampled to 16 kHz as described in Section 3. Decoding is 264
210 performed with voice activity detection enabled. 265
211 Where supported by the tokenizer, forced language decoding 266
212 is applied (English: “en”, Hausa: “ha”, Yoruba: “yo”). For 267
213 Igbo, language auto-detection is used because the faster- 268
214 whisper tokenizer does not include an Igbo language code. 269
215 Transcripts are normalized using the same preprocessing de- 270
216 scribed in Section 3 prior to computing evaluation metrics. 271

217 4.3 Fine-Tuning Strategy

218 For supervised adaptation, we fine-tune a multilingual speech 272
219 encoder using the Connectionist Temporal Classification 273
220 (CTC) objective. Experiments are conducted using the pre- 274
221 trained XLS-R checkpoint Mofe/xls-r-hausa-40, which was 275
222 originally trained for Hausa speech recognition. This check- 276
223 point is used as the starting point for all adaptation experi- 277
224 ments to study both language-specific fine-tuning and cross- 278
225 language transfer.
226 Fine-tuning is performed separately for each language using
227 the official training split, with model selection based on vali-
228 dation performance. In particular, we fine-tune the Hausa

checkpoint on the Hausa subset and additionally adapt the
same checkpoint to Igbo to evaluate cross-language adapta-
tion behavior.

Training is performed using AdamW with a linear learning-
rate schedule. Because experiments are conducted on CPU
hardware, we use small batch sizes together with gradient ac-
cumulation to ensure stable optimization. The feature en-
coder of the pretrained model is frozen during fine-tuning to
reduce memory usage and improve convergence stability.

To avoid invalid CTC alignments, examples where the label
sequence is longer than half the number of input frames are
filtered prior to training. This filtering rule ensures that the
CTC loss remains well-defined and prevents training insta-
bilities.

All experiments use fixed train/validation/test splits to ensure
comparability across languages and training regimes.

245 4.4 Evaluation Metrics

246 Performance is measured using Word Error Rate (WER) and
247 Character Error Rate (CER), computed on the official test
248 split. WER captures token-level transcription accuracy, while
249 CER provides robustness for morphologically rich or ortho-
250 graphically variable languages. Metrics are computed after
251 transcript normalization to ensure consistency across lan-
252 guages.

253 4.5 Implementation Details

254 All experiments are implemented in Python using the Hug-
255 ging Face Transformers and Datasets libraries, together with
256 faster-whisper for zero-shot inference. Zero-shot evaluation
257 is conducted on a CPU-only system (Intel i7-10510U, 16 GB
258 RAM) using int8 quantization. Fine-tuning experiments are
259 also performed on CPU hardware using gradient accumula-
260 tion to simulate larger batch sizes.

Dataset preprocessing, resampling to 16 kHz, manifest gen-
eration, and filtering of invalid CTC alignments are fixed
prior to experimentation to ensure reproducibility. All runs
use the same official train/validation/test splits provided with
the corpus.

266 5 Experiments

267 5.1 Zero-Shot Evaluation

268 We first evaluate Whisper-small in a zero-shot setting on the
269 official test split for each language. No supervised adaptation
270 is performed. Audio inputs are resampled to 16 kHz and de-
271 coded using faster-whisper with int8 CPU inference, as de-
272 scribed in Section 4.

Table 2 reports the average Word Error Rate (WER) and
Character Error Rate (CER) for English, Hausa, Igbo, and
Yoruba. Forced language decoding is applied for English,
Hausa, and Yoruba. For Igbo, language auto-detection is used
because the faster-whisper tokenizer does not include an Igbo
language code.

Language	Avg WER	Avg CER	Test Utterances	Decoding Language Setting
English	0.2739	0.1500	341	Forced (en)
Hausa	0.9546	0.3604	901	Forced (ha)
Igbo	1.4532	0.9965	572	Auto-detect (no 'ig' code)
Yoruba	1.5025	1.1076	418	Forced (yo)

279 **Table 2:** Zero-shot ASR performance of Whisper-small on the Ni-
280 gerian multilingual test split.

281 **Note:** WER/CER are computed on the official test split after
282 transcript normalization and 16 kHz resampling. For Igbo,
284 faster-whisper does not provide an "ig" language code, so de-
285 coding used language auto-detection.

286 Zero-shot performance varies substantially across languages.
287 English achieves comparatively low error rates, while Hausa
288 exhibits significantly higher WER despite forced language
289 decoding. Performance on Igbo and Yoruba is extremely
290 poor, with WER exceeding 1.45 and CER approaching or ex-
291 ceeding 1.0. These results indicate that Whisper-small gener-
292 alizes unevenly to Nigerian languages, particularly for lan-
293 guages that are likely underrepresented in the model’s pre-
294 training data.

295 The sharp degradation for Igbo and Yoruba suggests that
296 multilingual pretraining alone is insufficient to ensure robust
297 recognition for underrepresented African languages. These
298 findings motivate the supervised fine-tuning experiments de-
299 scribed in the next subsection.

300 5.2 Supervised Fine-Tuning (Hausa)

301 To examine whether supervised adaptation can mitigate the
302 performance gap observed in the zero-shot setting, we fine-
303 tune a pretrained Hausa ASR checkpoint based on XLS-R us-
304 ing the Hausa portion of the Nigerian multilingual speech
305 corpus. Training is performed on the official training split,
306 model selection is based on validation performance, and final
307 evaluation is conducted on the held-out test split.

308 Table 3 summarizes the results. After one epoch of fine-tun-
309 ing, the validation word error rate (WER) reaches 0.3563,
310 with a character error rate (CER) of 0.0912. After the second
311 epoch, the validation WER improves slightly to 0.3485 with
312 a CER of 0.0888, indicating stable convergence during train-
313 ing.

314 On the Hausa test split, the fine-tuned model achieves a WER
315 of 0.5306 and a CER of 0.5268. Compared with the zero-shot
316 Whisper-small baseline (WER 0.9546), supervised fine-tun-
317 ing reduces the WER by 44.4% relative. This substantial im-
318 provement demonstrates that multilingual pretrained speech

19 models benefit strongly from language-specific adaptation
20 when applied to underrepresented African languages.

21 Although the fine-tuned model does not reach the perfor-
22 mance typically observed for high-resource languages, the re-
23 sults confirm that zero-shot transfer alone is insufficient for
24 reliable recognition in Hausa, while supervised fine-tuning
25 significantly improves recognition accuracy. This finding
26 supports the hypothesis that underrepresented languages re-
27 quire targeted adaptation even when using large multilingual
28 pretrained models.

Model / Setting	Split	WER	CER	Notes
Whisper-small (zero-shot)	Test	0.9546	0.3604	No supervised adaptation
XLS-R Hausa (fine-tuned)	Validation	0.3563	0.0912	Epoch 1
XLS-R Hausa (fine-tuned)	Validation	0.3485	0.0888	Epoch 2
XLS-R Hausa (fine-tuned)	Test	0.5306	0.5268	Verified evaluation

329 **Table 3:** Supervised fine-tuning results for Hausa ASR. The
330 adapted XLS-R model significantly outperforms the zero-shot
331 Whisper-small baseline on the Hausa test split

332 5.3 Cross Language Adaptation (Igbo)

333 To examine whether the gains observed for Hausa generalize
334 to another Nigerian language, we further adapt the pretrained
335 XLS-R Hausa checkpoint on the Igbo portion of the corpus.
336 Table 4 reports the validation and test results. After one epoch
337 of fine-tuning, the validation WER reaches 0.9597, improv-
338 ing slightly to 0.9463 after the second epoch. On the Igbo test
339 split, the adapted model achieves a WER of 0.9464 and a
340 CER of 0.4341.

Model / Setting	Split	WER	CER	Notes
Whisper-small (zero-shot)	Test	1.4532	0.9965	No supervised adaptation
XLS-R Hausa checkpoint adapted to Igbo	Validation	0.9597	0.4584	Epoch 1
XLS-R Hausa checkpoint adapted to Igbo	Validation	0.9463	0.4385	Epoch 2
XLS-R Hausa checkpoint adapted to Igbo	Test	0.9464	0.4341	Final test evaluation

342 **Table 4:** Supervised fine-tuning results for Igbo ASR. Cross-lan-
343 guage adaptation from a Hausa XLS-R checkpoint improves per-
344 formance over the zero-shot baseline on the Igbo test split.
345

346
347 Compared with the zero-shot Whisper-small baseline (WER
348 1.4532), supervised adaptation reduces the Igbo test WER by
349 34.9% relative. Although the gains are smaller than those ob-
350 served for Hausa, the results still indicate that language-spe-
351 cific or cross-language supervised adaptation improves
352 recognition performance substantially over zero-shot transfer
353 alone. This suggests that pretrained multilingual speech mod-
354 els can benefit from further adaptation even when the starting
355 checkpoint was optimized for a different African language.

357 6. Discussion

358 The experimental results highlight several important obser-
359 vations about multilingual speech recognition for underrepre-
360 sented Nigerian languages. First, the zero-shot evaluation
361 demonstrates that performance varies substantially across
362 languages, even when using the same large pretrained model.
363 While Whisper-small performs reasonably well for English,
364 recognition accuracy drops sharply for Hausa, Igbo, and Yo-
365 ruba, with error rates exceeding those typically observed for
366 high-resource languages. This suggests that multilingual pre-
367 training alone is insufficient to ensure reliable performance
368 for underrepresented languages.

369 Second, supervised fine-tuning consistently improves recog-
370 nition accuracy across languages. For Hausa, adapting an
371 XLS-R Hausa checkpoint on the Nigerian Hausa corpus re-
372 duces the test WER from 0.9546 in the zero-shot setting to
373 0.5306 after fine-tuning, corresponding to a relative reduction
374 of 44.4%. A similar trend is observed for Igbo, where adapt-
375 ing the same checkpoint reduces the test WER from 1.4532
376 to 0.9464, a relative reduction of 34.9%. These results con-
377 firm that language-specific training data remains crucial even
378 when starting from large multilingual pretrained models.

379 However, the magnitude of improvement differs between
380 languages. The gains obtained for Hausa are larger than those
381 for Igbo, which may be explained by differences in training
382 data size, phonetic characteristics, or the mismatch between
383 the pretrained checkpoint and the target language. In this
384 study, the adaptation experiments use a checkpoint originally
385 trained for Hausa, which may provide a better initialization
386 for Hausa than for Igbo. This indicates that cross-language
387 transfer is possible but not uniform, and the effectiveness of
388 adaptation depends on both the similarity between languages
389 and the amount of available training data.

390 Overall, the results suggest that multilingual speech models
391 provide a strong starting point, but reliable recognition for
392 underrepresented African languages still requires targeted ad-
393 aptation. Zero-shot transfer alone is not sufficient, and super-
394 vised fine-tuning can substantially reduce error rates, alt-
395 hough performance remains below that typically achieved for
396 high-resource languages. These findings highlight the need
397 for more language-specific resources and better multilingual
398 training strategies to support speech recognition in un-
399 derrepresented languages.

400 This study has several limitations. First, experiments were
401 conducted on CPU hardware, which constrained the size of

402 models and the number of adaptations runs that could be per-
403 formed. Second, adaptation experiments were limited to a
404 single pretrained checkpoint, and results may differ when us-
405 ing other multilingual models. Third, the dataset size for
406 some languages is relatively small, which may affect the sta-
407 bility of training and evaluation. Future work should investi-
408 gate larger multilingual corpora, additional pretrained check-
409 points, and parameter-efficient adaptation methods to further
410 improve ASR performance for African languages.

411 7 Conclusion

412 This paper investigated multilingual automatic speech recog-
413 nition for Nigerian languages using zero-shot inference and
414 supervised fine-tuning with pretrained speech models. Exper-
415 iments were conducted on a Nigerian multilingual speech
416 corpus containing English, Hausa, Igbo, and Yoruba. In the
417 zero-shot setting, Whisper-small achieved reasonable perfor-
418 mance for English but produced high error rates for the Nige-
419 rian languages, with test WER values exceeding 0.95 for
420 Hausa and above 1.4 for Igbo. These results indicate that mul-
421 tilingual pretraining alone does not guarantee reliable recog-
422 nition for underrepresented languages.

423 To address this limitation, we evaluated supervised adapta-
424 tion using an XLS-R checkpoint. Fine-tuning on the Hausa
425 subset reduced the test WER from 0.9546 in the zero-shot
426 setting to 0.5306, corresponding to a relative improvement of
427 44.4%. A similar experiment on Igbo reduced the test WER
428 from 1.4532 to 0.9464, yielding a relative improvement of
429 34.9%. These results show that supervised fine-tuning con-
430 sistentlly improves recognition accuracy, although the magni-
431 tude of improvement varies across languages.

432 The difference in improvement between Hausa and Igbo sug-
433 gests that cross-language transfer from a pretrained check-
434 point is not uniform. In this study, the starting checkpoint was
435 trained for Hausa, which likely provided a better initialization
436 for Hausa than for Igbo. This highlights the importance of
437 language-specific training data and indicates that multilin-
438 gual speech models still require targeted adaptation for un-
439 derrepresented languages.

440 Overall, the experiments demonstrate that zero-shot multilin-
441 gual ASR remains insufficient for reliable recognition of Ni-
442 gerian languages, but supervised fine-tuning can substantially
443 reduce error rates. Future work should explore larger multi-
444 lingual training sets, language-balanced pretraining, and ad-
445 aptation methods designed specifically for underrepresented
446 African languages.

449 References

451 [Ardila et al., 2020] Rosana Ardila, Megan Branson, Kelly
452 Davis, Michael Henretty, Michael Köhler, Josh Meyer,
453 Reuben Morais, Lindsay Saunders, Francis M. Tyers,
454 and Gregor Weber. Common Voice: A Massively-Mul-
455 tilingual Speech Corpus. In Proceedings of LREC,
456 2020.

458 [Baeovski et al., 2020] Alexei Baeovski, Henry Zhou, Ab- 516
459 delrahman Mohamed, and Michael Auli. wav2vec 2.0: 517
460 A Framework for Self-Supervised Learning of Speech 518
461 Representations. In *Advances in Neural Information* 519
462 *Processing Systems (NeurIPS)*, 2020. 520
463 521

464 [Babu et al., 2021] Arun Babu, Changhan Wang, Andros 522
465 Tjandra, Kushal Lakhota, Qiantong Xu, Naman Goyal, 523
466 Kritika Singh, Patrick von Platen, Yatharth Saraf, Juan 524
467 Pino, and Alexei Baeovski. XLS-R: Self-supervised 525
468 Cross-lingual Speech Representation Learning at Scale. 526
469 arXiv preprint, 2021. 527
470 528

471 [Conneau et al., 2021] Alexis Conneau, Alexei Baeovski, Ronan 529
472 Collobert, Abdelrahman Mohamed, and Michael Auli. 530
473 Unsupervised Cross-Lingual Representation Learning 531
474 for Speech Recognition. In *Proceedings of Interspeech*, 532
475 2021. 533
476 534

477 [Lhoest et al., 2021] Quentin Lhoest, Albert Villanova del 535
478 Moral, Yacine Jernite, Abhishek Thakur, Patrick von 536
479 Platen, Suraj Patil, Julien Chaumond, Mariama Drame, 537
480 Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, 538
481 Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, 539
482 Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas 540
483 Patry, Angelina McMillan-Major, Philipp Schmid, Syl- 541
484 vain Gugger, Clément Delangue, Théo Matussière, Ly- 542
485 sandre Debut, Stas Bekman, Pierrick Cistac, Thibault 543
486 Goehringer, Victor Mustar, François Lagunas, Alexan- 544
487 der M. Rush, and Thomas Wolf. Datasets: A Commu- 545
488 nity Library for Natural Language Processing. arXiv 546
489 preprint, 2021. 547
490 548

491 [Nakatumba Nabende et al., 2025] Joyce Nakatumba Na- 549
492 bende, Sulaiman Kagumire, Caroline Kantono, and Pe- 550
493 ter Nabende. A Systematic Literature Review on Bias 551
494 Evaluation and Mitigation in Automatic Speech Recog- 552
495 nition Models for Low-Resource African Languages. 553
496 *ACM Computing Surveys*, 58(4), 2025. 554
497 555

498 [Olatunji et al., 2023] Tobi Olatunji, Tejumade Afonja, Ad- 556
499 itya Yadavalli, Chris Chinenye Emezue, Sahib Singh, 557
500 Bonaventure F. P. Dossou, Joanne Osuchukwu, Sa- 558
501 lomey Osei, Atnafu Lambebo Tonja, Naome Etori, and 559
502 Clinton Mbataku. AfriSpeech-200: Pan-African Ac- 560
503 cented Speech Dataset for Clinical and General Domain 561
504 ASR. *Transactions of the Association for Computa- 562
505 tional Linguistics*, 2023. 563
506 564

507 [Radford et al., 2022] Alec Radford, Jong Wook Kim, Tao 565
508 Xu, Greg Brockman, Christine McLeavey, and Ilya 566
509 Sutskever. Robust Speech Recognition via Large-Scale 567
510 Weak Supervision. arXiv preprint, 2022. 568
511 569

512 [Reid and Williams, 2023] Kathy Reid and Dominic Wil- 570
513 liams. Common Voice and Accent Choice: Data Con- 571
514 tributors Self-Describe Their Spoken Accents in Di- 572
515 verse Ways. In *Proceedings of ACM Conference*, 2023. 573

[Sanni et al., 2025] Mardhiyah Sanni, Tassallah Abdullahi, 574
Devendra D. Kayande, Emmanuel Ayodele, Naome A. 575
Etori, Michael S. Mollel, Moshood Yekini, Chibuzor 576
Okocha, Lukman E. Ismaila, Folafunmi Omofoye, Bo- 577
luwatife A. Adewale, and Tobi Olatunji. AfriSpeech-Di- 578
alog: A Benchmark Dataset for Spontaneous English 579
Conversations in Healthcare and Beyond. In *Proceed- 580
ings of NAACL*, 2025. 581

[Chen et al., 2025] William Chen, Chutong Meng, et al. The 582
ML-SUPERB 2.0 Challenge: Towards Inclusive ASR 583
Benchmarking for All Language Varieties. arXiv pre- 584
print, 2025. 585

[Muchai et al., 2026] Mercy Muchai, Kevin Chege, et al. 586
Paza: Introducing Automatic Speech Recognition 587
Benchmarks and Models for Low Resource Languages. 588
Microsoft Research, 2026. 589

[Bandarupalli, 2025] S. Bandarupalli. Efficient ASR for 590
Low-Resource Languages. arXiv preprint, 2025. 591

[Sameti et al., 2025] Mohammad Hossein Sameti, Sepehr 592
Harfi Moridani, Ali Zarean, and Hossein Sameti. Ac- 593
cent-Invariant Automatic Speech Recognition via Sali- 594
ency-Driven Spectrogram Masking. arXiv preprint, 595
2025. 596

[Swain et al., 2024] M. Swain, et al. On Mitigating Perfor- 597
mance Disparities in Multilingual ASR via Fairness- 598
Aware Fine-Tuning. In *Proceedings of EMNLP*, 2024. 599

[Imam et al., 2025] Sukairaj Hafiz Imam, et al. Automatic 600
Speech Recognition for African Low-Resource Lan- 601
guages: A Systematic Literature Review. *ACM Comput- 602
ing Surveys*, 2025. 603

[Mainzinger, 2024] J. Mainzinger. Fine-Tuning ASR Models 604
for Very Low-Resource Languages. In *Proceedings of 605
ACL Student Research Workshop*, 2024. 606