

SPEECH RECOGNITION DATASETS FOR LOW-RESOURCE CONGOLESE LANGUAGES

Ussen Kimanuka¹, Ciira wa Maina^{2,3}, Osman Büyük⁴,

¹ Pan African University Institute for Basic Sciences, Technology and Innovation, Kenya,

² Dedan Kimathi University of Technology, Kenya,

³ Centre for Data Science and Artificial Intelligence (DSAIL), Kenya,

⁴ Izmir Demokrasi University, Turkey,

Correspondence: abre.ussen@students.jkuat.ac.ke

ABSTRACT

Large pre-trained Automatic Speech Recognition (ASR) models have begun to perform better in low-resource languages due to the availability of data and transfer learning. However, a few languages have sufficient resources to benefit from transfer learning. This paper contributes to expanding speech recognition resources for under-represented languages. We release two new datasets to the research community: Lingala Read Speech Corpus consisting of 4 hours of labelled audio clips and Congolese Speech Radio Corpus containing 741 hours of unlabeled audio in 4 major spoken languages in the Democratic Republic of the Congo. Additionally, we obtain benchmark results for Congolese wav2vec2. We observe an average decrease of 2% in WER when a Congolese multilingual pre-trained model is used for finetuning on Lingala. Importantly, our study is the first attempt towards benchmarking speech recognition systems for Lingala and the first-ever multilingual model for 4 Congolese languages spoken by 65 million people. Our data and models will be publicly available, and we hope they help advance research in ASR for low-resource languages.

1 INTRODUCTION

Low-resource languages are an important and under-represented area in conversational AI. These languages suffer from a serious lack of training data for speech recognition systems, mostly developed for high-resource languages like English.

Automatic speech recognition (ASR), considered one of the primary building blocks of Conversational AI, uses training corpora created by researchers to construct acoustic models capable of speech-to-text conversion. However, the development of speech recognition technology needs vast datasets with annotations. Unfortunately, the languages spoken in Sub-Saharan Africa tend to fall into the “low-resource” category, which, in contrast to “high-resource” languages, has fewer datasets accessible. Transfer learning, or the transfer of representations gained on unrelated high-resource languages, has not been studied for most low-resource languages (Barnard et al., 2010; Besacier et al., 2014). Even with the assistance of transfer learning, labelled data are still required to construct viable models (Kunze et al., 2017).

It is challenging to create high-quality ASR models for such a broad and varied array of languages. For the development of high-quality ASR systems, many contemporary ASR models depend on vast quantities of labelled data for each language. Such techniques are costly and not scalable, restricting the use of ASR technology to a subset of languages and populations. In addition to these issues of the availability of labelled data, Congolese languages¹ confront several other issues that must be addressed. Most of these languages have less digital content on the web, which pose obstacles to adopting language models. Likewise, there are prospects from a united standpoint. As numerous sources may be shared across different languages, the collection of unlabelled data for pre-training may be accomplished collaboratively. A vastly overlapping phoneme inventory may be used to construct multilingual models in which transfer learning is successful (Javed et al., 2021).

¹https://en.wikipedia.org/wiki/languages_of_Africa

In this contribution, we report initial development work on a broadband automatic speech recognition (ASR) system for the Lingala language and build a multilingual acoustic model for Congolese languages. As in (Doumbouya et al., 2021), we employ unsupervised speech representation learning to solve the data scarcity issue and confirm representations acquired from radio archives, which are rich in many areas of the globe, which may be used for speech recognition in low-resource conditions. In this regard, it is essential to notice recent research demonstrating the advantages of unsupervised pretraining and multilingual fine-tuning for dramatically improving ASR quality for low-resource languages (Baevski et al., 2020; Conneau et al., 2020). In particular, the Wav2Vec2.0 model has produced two significant outcomes for English ASR. First, SOTA outcomes are established through an end-to-end DNN design that borrows the popular Transformer architecture from NLP. In addition, pretraining on a large dataset decreases the labelled fine-tuning time from hundreds of hours to a few hours or even tens of minutes. It is important to investigate if these results from English ASR may be applied to Lingala ASR, considering the variety and aforementioned difficulties associated with Niger-Congo Bantu languages. Our objective is to answer the following research questions:

- Can a model similar to Wav2Vec2.0 provide SOTA performance on existing benchmarks for Lingala?
- Could extensive pretraining eliminate the requirement to gather vast quantities of labelled data?
- Could pretraining a multilingual model spanning a few Congolese languages benefit these related languages?

In our journey to answer these questions, we offer the following contributions:

- **We assemble two novel datasets:** The first is the Congolese Speech Radio Corpus, while the second is the Lingala Read Speech Corpus. In the first dataset, we gather 741 hours of raw audio data for pretraining across four languages from the Niger-Congo Bantu language family, making it one of the biggest and most varied collections of Niger-congo B language data. The second dataset is a preliminary attempt to produce linguistic resources for speech recognition systems in Lingala, 4.35 hours of labelled data including a text corpus, a speech corpus, a lexical dictionary, and a language model.
- **We provide a baseline for the collected datasets:** Starting with the Supervised Learning (SL) method, we perform contrastive studies against Self-Supervised Learning (SSL) paradigms, particularly the Wav2Vec 2.0 model, and then we conduct extensive experiments to arrive at a training and decoding rule that is effective for the Lingala language.
- To the best of our knowledge, this is the first work carried out towards creating baselines for Lingala and the only created multilingual acoustic model combining four Congolese languages (i.e. Kikongo, Tshiluba, Lingala and Swahili).
- **To encourage further research** in the field of underrepresented languages such as Niger-Congo Bantu ASR, we will make all our artefacts available. This comprises (a) sources of pretraining data as well as scripts for their collection and preprocessing, (b) pretraining, finetuning, and decoding scripts, (c) language models, and (d) our top ASR models.

2 RELATED WORK

Previous work in Natural Language Processing has made it possible to gather data, construct monolingual and multilingual ASR models, and push ASR research further for Niger-Congo Bantu languages.

The current state of speech datasets for African languages. Efforts have been made to compile speech datasets in low-resource African languages. The Niger-Congo Bantu language family have 8 documented dataset collection: Basaa, Myene, and Embosi (Adda et al., 2016), Godard et al. (2017) collected speech corpus for a realistic language documentation process for Mboshi, Öktem et al. (2020) collected Congolese Swahili for developing humanitarian response against Covid-19, Doumbouya et al. (2021) collected the West African Virtual assistant speech recognition corpus for Maninka, Susu, and Pular. Amharic, Swahili, and Wolof are only a few African

Table 1: Numbers of speakers, phonological features and orthographic conventions

Languages	ISO	Population	Tones	Diacritical	Example
Lingala	ln	40M	high,low	Yes	lɛlo 'today'
Kongo/Kikongo	kon	6.9M	high,low	Yes	mbó 'to hit'
Congolese Swahili	swc	11M	high,low	No	mtoto 'child'
Luba-Kasai/Tshiluba	lua	7M	high,mid,low	Yes	kutúla 'finished'

languages for which open datasets are available (Abate et al., 2005; Gelas et al., 2012; Gauthier et al., 2016a). Additionally, East African languages (Babirye et al., 2022), as well as South African languages (Badenhorst et al., 2011; Dr et al., 2013; de Wet et al., 2017), have been the focus of several efforts. The open sourcing of some of these African datasets has been a driving force in ASR research, as it encourages experimentation, training, and the development of better models. To this end, we could identify a dataset (Conneau et al., 2022) that included only one of the four Niger-Congo B languages we focused on. However, the dataset was noisy and lacked curation, which prevented further testing.

Utilizing “found” data, such as radio archives (Doubouya et al., 2021) investigated the use of found data, including ASR data and radio broadcasting, for three Niger-Congo Bantu languages. However, the data mostly included West African languages and did not include the central African Niger-Congo Bantu languages. (Cooper, 2019) studied radio news and audiobooks for low-resource language text-to-speech synthesis. However, the radio broadcast is good quality and in English, instead of being noisy and having limited resources. Some research has focused on speech processing using “found” data in Indian languages (Mendels et al., 2015; Baljekar et al., 2018; Gupta et al., 2021; Javed et al., 2021). None of these found data projects includes noisy radio data for low-resource Niger-Congo languages in the Central African Regions.

Self-supervised speech representation learning Self-supervised speech representation learning approaches such as Wav2vec 2.0 aim to learn speech representation from unlabeled data to improve accuracy on downstream tasks such as phoneme classification, speaker recognition, sentiment analysis, speech recognition, and phoneme recognition while using fewer training data points (Baevski et al., 2020; Ritchie et al., 2022) This paper compares the multilingual Wav2vec 2.0 (XLS-R) large model pre-trained on 128 languages - a large (436 thousands of hours) corpus (Babu et al., 2021) - to its counterpart trained on a small (741.5 hours) dataset of noisy radio broadcasting archives in Central African languages for speech recognition on Niger-Congo B languages.

3 CONGOLESE LANGUAGES

There are more than 200 languages spoken in the Democratic Republic of the Congo (Eberhard et al., 2021). The top four languages in terms of global use are Lingala, Kongo or Kikongo, Luba-Kasai or Tshiluba and Congolese Swahili. These are Bantu languages, particularly a branch of the Niger-Congo B language family. They employ prefixes, suffixes, and infixes to show grammatical connections, have a sophisticated system of noun classes and verb conjugation, and also have a tonal system, agglutinative morphology, Subject Verb Object word order, complex syllable structure, and vowel harmony (Sene-Mongaba, 2015). All three languages in this study are written using the standard 26 letters of the Latin alphabet, except for Lingala, which possesses two additional special characters: the open vowels ϵ and \circ . Three languages may use diacritical marks to indicate certain linguistic features, such as tone, vowel length, or emphasis. A summary of the number of speakers, some phonological features and orthographic conventions used for the 4 languages are given in Table 1.

4 DATA COLLECTION

In this research, we present two datasets: the Lingala Read Speech Corpus, which can be used to create the large vocabulary supervised speech recognition module for different applications, and the Congolese Radio Corpus, which is designed for unsupervised speech representation learning for downstream tasks involving 4 Congolese languages.

Table 2: General characteristics of the Lingala Read Speech corpus

Subset	Duration (hours)	Utterances
Train	4	2465
Dev	0.2	204
Test	0.1	180
Total	4.3	2849

Table 3: The token counts for the two sets of text sources used to create the language models task. LM_{small} relate to texts originating from the transcript of the corpus, and LM_{large} are texts from additional online resources.

Language model	Sentences	No. of unique tokens
LM_{small}	2557	3980
LM_{large}	22 000	22261

4.1 LINGALA READ SPEECH CORPUS

Description The corpus contains a total of more than 4.3 hours of voice data. In order to allow future experimentation, we offer consistent training, development, and test divides. There is no speaker overlap between the subgroups. All recordings utilise Waveform Audio File Format (WAVE), and each file is encoded with a sample rate of 16 kHz and a single track (mono).

Collection of the dataset To create the Lingala Read Speech corpus, we utilised the Lig-Aikuma app (Gauthier et al., 2016b) for speech recording. The software’s elicitation mode was used to capture their readings of text scripts that had been tokenized at the sentence level. Numerous scholars have utilized the Lig-Aikuma in their own similar works (Sikasote & Anastasopoulos, 2021; Blachon et al., 2016; Gauthier et al., 2016c).

Speakers Lingala Read Speech recording speakers were from Kinshasa. In all, there are 2849 recordings recorded by 32 distinct speakers (13 male and 19 female). The ability to speak and read Lingala was a requirement for all speakers. Noting that the recordings in this corpus were not created under perfectly controlled conditions. The speakers’ voices were recorded in a variety of conditions and with a varied range of accents. Consequently, it is anticipated that some utterances will include some ambient noise. This enables us to train and evaluate ASR systems in environments that more closely mimic the real world than a studio setting.

Preprocessing All utterances were validated against the transcripts to ensure the accuracy of the data. According to the utterances, the text was modified such that all numbers, dates, and times were converted into their textual equivalents. We also used speaker ID to sort the recorded audio files into groups based on the speaker. We also renamed all the audio files by appending the speaker ID to the utterance IDs. Table 2 shows the distribution of the dataset.

Text sources The recorded phrases and sentences are drawn from a wide variety of Lingala language sources, including Lingala literature. Table 3 summarises the textual contributions to the Lingala Read Speech corpus and the distribution of text used for language modelling tasks. Each phrase ranges in length from one to twenty words.

Availability The corpus will be accessed at our GitHub project repository and will be shared with the research community under a creative commons license.

4.2 CONGOLESE SPEECH RADIO CORPUS

Description There are 741 hours of audio clips in the Congolese Radio Corpus, with the longest being 25 seconds long. These recordings were taken from the archives of four different broadcast stations in the Democratic Republic of the Congo. News and Radio programmes are aired in many Niger-Congo B languages, including Lingala, Tshiluba, Kikongo, and Congolese Swahili. There

Table 4: The number of Congolese Speech Radio Corpus hours per language post preprocessing.

Language	Radio archives(hours)	Youtube (hours)	Total(hours)
Lingala	187.6	86.4	274.1
Kikongo	174.4	-	174.4
Tshiluba	112.8	-	112.8
Congolese Swahili	191.1	-	191.1

may be a variety of sounds and music playing in the background and front of certain radio content. Crawling the websites of the Radio archives (R.A) to get the URLs of audio recordings in 4 Congolese languages simplified the gathering procedure. We not only curated material from the radio’s official website but also from YouTube clips that the stations had released underneath a Creative Commons License. Although this drastically reduced the quantity of data we could gather, particularly in the low-resource Niger-Congo language family, it was necessary to guarantee that we may freely share the URLs with the scientific world, assuring the repeatability of our study.

Preprocessing The data was curated in the same way as in [8] after the URLs were bought to light. To begin, we grabbed the audio clips from the websites of the various radio broadcasts; for those files hosted on YouTube, we utilised the Youtube-dl library². Secondly, the information was not always mono channel, and the sampling frequency ranged from 8 kilohertz to 44 kilohertz since the data was picked from various sources. We used the FFmpeg library³ to upsample/downsample the data that was captured at a frequency less than/ higher than 16 kHz and then decrease the number of audio channels to 1. Thirdly, we further improved the data by excluding extended periods of quietness from the audio files using the Py-webrtcvad library⁴, a python interface to the widely used WebRTC VAD(Voice Activity Detection) module built by Google. The VAD algorithm screens out noise and lets us choose a harshness parameter (a number between 0 and 3), which determines how strictly to apply the filter on speech (0 is the least aggressive about filtering out non-speech, and 3 is the most aggressive). Similar to the findings in (Javed et al., 2021), we found that changing this value to 2 provided the best results for our data. Also, we utilised Waveform Amplitude Distribution Analysis (WADA-SNR) (Kim & Stern, 2008) to exclude audio samples with a signal-to-noise ratio (SNR) below 15 dB from our dataset. Using a sample of the audio files, a threshold was optimal. Finally, we followed industry standards by chunking audio recordings to a maximum of 25 seconds in length. Table 4 summarises the data acquired in this way.

Availability The audio file URLs, as well as the scripts used to collect the data and clean it up, can be accessed on the following link⁵

5 EXPERIMENTS

In this section, we describe the experiments to ascertain the usefulness of the two novel speech corpus for ASR applications in supervised, weakly-supervised and self-supervised learning configurations.

5.1 SUPERVISED MODELS

The following experiment includes two supervised baseline models; These are traditional GMM-HMM based followed by a Time-Delay neural network (T-DNN) model (Peddinti et al., 2015) and the end-to-end architectural framework (Chang et al., 2020).

The T-DNN architecture consists of 4 hidden layers with layerwise temporal contexts of [-2,2],-1,2,-3,3 and -7,2, utilising a P-norm input dimension of 3000 and a group size of 10. The architecture follows the implementation of Kaldi recipe (Povey et al., 2011; Kimanuka & Büyük, 2018). The Deep Speech2-like architecture is an end-to-end sequence-to-sequence model with two primary neural networks modules: N-blocks of Residual Convolutional Neural Networks (rCNNs) and combined

²<https://github.com/tpikonon/youtube-dl>

³<https://ffmpeg.org/>

⁴<https://github.com/wiseman/py-webrtcvad>

⁵https://github.com/ussenuk/Congolese_Wav2Vec_datasets_experiment

M-blocks of Bidirectional Long Short Term Memory (BiLSTMs) and Bidirectional Gated Recurrent Units (Bi-GRUs). The rCNNs function as feature extractors by translating speech into spectrograms, while the combined BiLSTMs and BiGRUs serve as encoder and decoder. We adhere to the same design as (Meyer, 2019). The final hidden state of the last Bi-GRU block is utilised to calculate the attention weights and the context vector, concatenated with the BI-GRU output to serve as the final outputs. These two supervised experiments implement the 4.3 hours of supervised Lingala Read Speech corpus.

5.2 WEAKLY-SUPERVISED MODEL

The weakly-supervised Whisper model (Radford et al., 2022) is an encoder-decoder Transformer with an end-to-end auto-regressive architecture. This model has been pre-trained using a total of 680 000 hours of supervised data from various languages and tasks. The encoder uses a pair of convolution functions of GELU, the first of which has a filter width of 3 and an activation function of GELU, and the second has a stride of 2. The encoder and decoder are architecturally identical, with the same width and number of transformer blocks. The Whisper model was pre-trained in a weakly supervised manner using varying sizes. We used the medium-sized model for this experiment and performed supervised finetuning using the 4.3 hours of supervised Lingala Read Speech corpus. Since Whisper models are pre-trained for multitasking, it is important to specify the task before finetuning the operation. Hence we selected the "Transcribe task" before finetuning.

5.3 SELF-SUPERVISED MODEL

Using the unlabelled audio data given in the Congolese Speech Radio corpus, we implemented the pretraining of a self-supervised ASR model for Congolese languages. We strictly adhere to the process outlined in (Baevski et al., 2020) for developing a self-supervised ASR system.

For this self-supervised experiment, we selected Wav2Vec2.0 architecture which is an end-to-end model consisting primarily of a feature encoder for converting raw audio into a sequence of T latent representations, a Transformer for learning contextualised representations for each of the T units and a quantizer for discretizing the representations learned by the feature encoder. Training this model involves masking some T input representations before passing them to the context network. The quantizer's inputs, however, are not masked since they constitute the target of the Self-learning procedure. The objective is achieved by minimizing a contrastive loss between the context network's output and the quantized representation.

Pretraining a Wav2Vec2 model from scratch (Congolese Wav2vec / CdWav2Vec) We pre-train only on the BASE model similar to that of [9], with 7 convolutional layers, 512 channels per model, strides of (5,2,2,2,2,2), and kernel widths of (10,3,3,3,3,2,2). This BASE model comprises 12 transformer blocks, 768 model dimensions, 3072 FFN dimensions, and 8 attention heads. In the quantization module of this design, we use $G=2$ (codebooks) with $V=320$ elements per codebook. We have used the pre-trained checkpoint of the equivalent (BASE) English Wav2Vec 2.0 model to kick off our pretraining process. We next use our carefully chosen data on the Congolese Speech Radio archive dataset to further pre-train the model. Shorter audio clips of 15.6 seconds in length (256k samples) are used in the BASE model. No changes are set to any of the hyperparameters from the standard configuration of the Wav2Vec 2.0 source code.⁶

Finetuning a Wav2Vec2 model We update all of the network parameters besides the convolutional features encoder parameters during fine-tuning. In our finetuning experiment, we considered two scenarios: (1) Finetuning the Congolese Wav2vec (CdWav2Vec) model pre-trained on the 4 Congolese languages (Transfer learning); (2) Finetuning multilingual XLSR-53 (Conneau et al., 2020) and XLS-R (Babu et al., 2021) model, which is also a Wav2Vec 2.0 based model trained on freely accessible speech audio data in 53 and 128 languages respectively. We explored experiments on the XLS-R-0.3B model with 600 million parameters. Both the fine-tuning experiments used the supervised Lingala Read Speech corpus implemented using either the HuggingFace transformer library (Wolf et al., 2020) or the Fairseq toolkit (Ott et al., 2019). The hyperparameters used for pretraining and finetuning are similar to the BASE model in (Baevski et al., 2020).

⁶<https://github.com/pytorch/fairseq/tree/master/examples/wav2vec#pretrained-models>

Transfer learning From Congolese languages to Lingala From the Congolese Wav2Vec pre-trained model, we tried implementing cross-lingual transfer learning in our quest to find a model with improved performance. To do this, we finetuned the multilingual Congolese model pre-trained on the 4 Congolese languages on our supervised Lingala Read Speech corpus using Adam optimizer, with a learning rate of $1e-4$ and a three-stage learning rate schedule for 80 000 steps.

6 RESULTS

The results of our experimental evaluation are summarised in Table 5. The model with the best performance was the CdWav2Vec multilingual model obtained by pretraining on the Congolese speech radio corpus and then fine-tuning on the Lingala Read Speech corpus. We were able to get a WER of 21.4% using this model. CdWav2Vec model consistently outperforms the data greedy end-to-end architecture and the multilingual self-supervised XLS-R model with an exception on weak supervision⁷, which tends to provide competitive results. Transferring learning from the multilingual pre-trained CdWav2Vec model to the Lingala Read speech corpus means we obtained 2% relative WER improvement over the multilingual XLS-R models. The improvement in the performance of the CdWav2Vec model compared to the supervised models demonstrates the superior cross-lingual transfer ability of multilingually pre-trained models over supervised (trained from scratch) models in low resource scenarios cases from three Congolese languages to Lingala. These results indicate that multilingual pre-trained models perform better in cross-lingual transfers than supervised monolingual models.

6.1 IMPACT ON FINETUNING

Using checkpoints from monolingual and multilingual pretraining, we fine-tune 4.3 hours of Lingala data. Even with a small amount of supervised data, Lingala benefits from multilingual pretraining. Table 6 demonstrates further that fine-tuning monolingual and multilingual pre-trained models resulted in a 1% WER difference when decoding with or without a language model.

6.2 TEST SETS RESULTS

Next, we focused on the CdWav2Vec model and tested it on an out-of-domain TICO-19 test set, which mostly consists of read speech Congolese Swahili recorded from the Devset of the machine translation benchmark for the COVID-19 domain (Anastasopoulos et al., 2020). Table 7 in the Appendix section shows that our CdWav2vec multilingual model outperforms an end-to-end model pre-trained in English from (Öktem, 2021)

6.3 IMPACT OF CROSS-LINGUAL REPRESENTATIONS

It has been explained earlier that the first stage of the CdWav2vec model entails computing representations of audio frames from a learned codebook. All four languages are represented in the model using the same codebook vectors. In the same attempt as (Gupta et al., 2021), we try to answer whether these codebooks are independent or shared across languages and the degree to which the related languages share codebook entries increases with linguistic distance. To answer the question, we created a plot which displays the quantized speech representations for all the languages in our study. To do so, 200 samples are chosen randomly from each language for a grand total of 5 hours of data. For the purpose of creating codebook vectors, we first run the audio through the feature encoder and then through the quantizer. For each sequence, the quantizer module’s output vectors. For each language, vectors are normalised to form $V \times G$ vectors of size, followed by K-means clustering and principal component analysis to reduce the dimensions of these vectors (PCA). As seen in Figure 1, groupings of languages most phonetically similar tend to cluster together. This proves that multilingual pretraining may promote representation learning that is transferable across linguistically related languages.

⁷It worth noting that The weakly supervised Whisper model contained the non-curated Lingala subset of Fleurs dataset not considered in this study

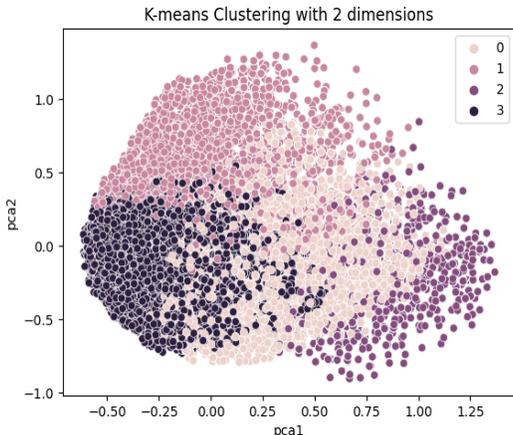


Figure 1: Quantized speech representation where 0:Kikongo, 1: Lingala, 2: Congolese Swahili and 3: Tshiluba)

Table 5: Results summary reporting the best performance of each model.

Model	WER(%)	CER(%)
<i>Traditional Supervised ASR</i>		
GMM-HMM	39.9	15.9
T-DNN	33.4	13.4
<i>End-to-End Supervised ASR</i>		
rCNN + BiLSTMs + BiGRUs	100	95
<i>Weakly-supervised model</i>		
Whisper fine-tuned(*)	20.9	5.8
<i>multilingual self-supervised models</i>		
XLSR-53	23.3	6.8
XLS-R-0.3B	25.8	7.0
CdWav2Vec	21.4	6.8
CdWav2Vec+ LM _{small}	19.1	-
CdWav2Vec+ LM _{large}	18.4	6.3

Table 6: Effect of multilingual and monolingual pretraining on WER

Pretraining	Finetuning	Decoding	WER	CER
monolingual	Lingala	Viterbi	22.1	6.8
multilingual	Lingala	Viterbi	21.4	6.8

7 CONCLUSION

Building speech recognition systems require very large datasets to train the models. We created two novel datasets in this paper: a 4.3-hour labelled Lingala Read Speech Corpus for Lingala and a 741-hour unlabeled Congolese Speech Radio Corpus for four major Congolese languages. We provide an extensive evaluation using both supervised and self-supervised learning settings. Our experimental results show the potential of multilingual pretraining approaches like wav2vec 2.0 for low-resource Congolese languages. We created the first benchmark for speech recognition systems for the Lingala language and the first-ever multilingual model for 4 major Congolese languages (CdWav2Vec). We have answered research questions and provided a starting point for further research, and we will publicly release the models and datasets to promote further speech technology research.

ACKNOWLEDGMENTS

We thank Radio Okapi for donating radio archives and the 32 speakers who contributed their voices and time to the speech recognition corpus. We thank Yannis Assael for providing critical guidance that made this research possible. We also thank the Africa PhD Fellowship Program for the resources provided during the research.

REFERENCES

- Solomon Teferra Abate, Wolfgang Menzel, and Bairu Tafila. An Amharic speech corpus for large vocabulary continuous speech recognition. In *Proc. Interspeech 2005*, pp. 1601–1604, 2005. doi: 10.21437/Interspeech.2005-467.
- Gilles Adda, Sebastian Stüker, Martine Adda-Decker, Odette Ambouroué, Laurent Besacier, David Blachon, Hélène Bonneau-Maynard, Pierre Godard, Fatima Hamlaoui, Dmitry Idiatov, Guy-Noël Kouarata, Lori Lamel, Emmanuel-Moselly Makasso, Annie Rialland, Mark Van de Velde, François Yvon, and Sabine Zerbian. Breaking the unwritten language barrier: The bulb project. *Procedia Computer Science*, 81:8–14, 2016. ISSN 1877-0509. doi: <https://doi.org/10.1016/j.procs.2016.04.023>. URL <https://www.sciencedirect.com/science/article/pii/S1877050916300370>. SLTU-2016 5th Workshop on Spoken Language Technologies for Under-resourced languages 09-12 May 2016 Yogyakarta, Indonesia.
- Antonios Anastasopoulos, Alessandro Cattelan, Zi-Yi Dou, Marcello Federico, Christian Federmann, Dmitriy Genzel, Francisco Guzmán, Junjie Hu, Macduff Hughes, Philipp Koehn, Rosie Lazar, Will Lewis, Graham Neubig, Mengmeng Niu, Alp Öktem, Eric Paquin, Grace Tang, and Sylwia Tur. TICO-19: the Translation initiative for COvid-19, 2020. arXiv:2007.01788.
- Claire Babirye, Joyce Nakatumba-Nabende, Jeremy Franics, Jonathan Mukiibi, Andrew Katumba, Ronnie Ogwang, Medadi Sentanda, Lilian Wanzare, and Davis David. Building text and speech datasets for low resourced languages: A case of languages in east africa. 03 2022.
- Arun Babu, Changhan Wang, Andros Tjandra, Kushal Lakhota, Qiantong Xu, Naman Goyal, Kritika Singh, Patrick von Platen, Yatharth Saraf, Juan Pino, Alexei Baevski, Alexis Conneau, and Michael Auli. Xls-r: Self-supervised cross-lingual speech representation learning at scale, 2021. URL <https://arxiv.org/abs/2111.09296>.
- Jaco Badenhorst, Charl Johannes van Heerden, Marelie Hattigh Davel, and Etienne Barnard. Collecting and evaluating speech recognition corpora for 11 south african languages. *Language Resources and Evaluation*, 45:289–309, 2011.
- Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *CoRR*, abs/2006.11477, 2020. URL <https://arxiv.org/abs/2006.11477>.
- Pallavi Baljekar, SaiKrishna Rallabandi, and Alan W Black. An Investigation of Convolution Attention Based Models for Multilingual Speech Synthesis of Indian Languages. In *Proc. Interspeech 2018*, pp. 2474–2478, 2018. doi: 10.21437/Interspeech.2018-1869.
- Etienne Barnard, Johan Schalkwyk, Charl van Heerden, and Pedro J. Moreno. Voice search for development. In *Proc. Interspeech 2010*, pp. 282–285, 2010. doi: 10.21437/Interspeech.2010-111.
- Laurent Besacier, Etienne Barnard, Alexey Karpov, and Tanja Schultz. Automatic speech recognition for under-resourced languages: A survey. *Speech Communication*, 56:85–100, 2014. ISSN 0167-6393. doi: <https://doi.org/10.1016/j.specom.2013.07.008>. URL <https://www.sciencedirect.com/science/article/pii/S0167639313000988>.
- David Blachon, Elodie Gauthier, Laurent Besacier, Guy-Noël Kouarata, Martine Adda-Decker, and Annie Rialland. Parallel speech collection for under-resourced language studies using the ligaikuma mobile device app. *Procedia Computer Science*, 81:61–66, 2016. ISSN 1877-0509. doi: <https://doi.org/10.1016/j.procs.2016.04.030>. URL <https://www.sciencedirect.com/science/article/pii/S1877050916300448>. SLTU-2016 5th Workshop on Spoken Language Technologies for Under-resourced languages 09-12 May 2016 Yogyakarta, Indonesia.
- Xuankai Chang, Wangyou Zhang, Yanmin Qian, Jonathan Le Roux, and Shinji Watanabe. End-to-end multi-speaker speech recognition with transformer, 2020. URL <https://arxiv.org/abs/2002.03921>.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Un-supervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual*

- Meeting of the Association for Computational Linguistics*, pp. 8440–8451, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.747. URL <https://aclanthology.org/2020.acl-main.747>.
- Alexis Conneau, Ankur Bapna, Yu Zhang, Min Ma, Patrick von Platen, Anton Lozhkov, Colin Cherry, Ye Jia, Clara Rivera, Mihir Kale, Daan Van Esch, Vera Axelrod, Simran Khanuja, Jonathan H. Clark, Orhan Firat, Michael Auli, Sebastian Ruder, Jason Riesa, and Melvin Johnson. Xtreme-s: Evaluating cross-lingual speech representations, 2022. URL <https://arxiv.org/abs/2203.10752>.
- Erica Cooper. Text-to-speech synthesis using found data for low-resource languages. 2019.
- Febe de Wet, Neil Kleynhans, Dirk Van Compernelle, and Reza Sahraeian. Speech recognition for under-resourced languages: Data sharing in hidden markov model systems. *South African Journal of Science*, 113:1–9, 2017.
- Moussa Doumbouya, Lisa Einstein, and Chris Piech. Using radio archives for low-resource speech recognition: Towards an intelligent virtual assistant for illiterate users. *CoRR*, abs/2104.13083, 2021. URL <https://arxiv.org/abs/2104.13083>.
- Henselmans Dr, Thomas R. Niesler, and D. Vanleeuwen. Baseline speech recognition of south african languages using lwazi and ast. 2013.
- David M. Eberhard, Gary F. Simons, and Charles D. Fennig. *Ethnologue: Languages of the world*. twenty-third edition., 2021. URL <http://www.ethnologue.com>.
- Elodie Gauthier, Laurent Besacier, Sylvie Voisin, Michael Melese, and Uriel Pascal Elingui. Collecting resources in sub-Saharan African languages for automatic speech recognition: a case study of Wolof. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pp. 3863–3867, Portorož, Slovenia, May 2016a. European Language Resources Association (ELRA). URL <https://aclanthology.org/L16-1611>.
- Elodie Gauthier, David Blachon, Laurent Besacier, Guy-Noël Kouarata, Martine Adda-Decker, Annie Rialland, Gilles Adda, and Grégoire Bachman. Lig-aikuma: A mobile app to collect parallel speech for under-resourced language studies. In *Interspeech*, 2016b.
- Elodie Gauthier, David Blachon, Laurent Besacier, Guy-Noel Kouarata, Martine Adda-Decker, Annie Rialland, Gilles Adda, and Grégoire Bachman. LIG-AIKUMA: a Mobile App to Collect Parallel Speech for Under-Resourced Language Studies. In *Interspeech 2016 (short demo paper)*, San-Francisco, United States, September 2016c. URL <https://hal.science/hal-01350062>.
- Hadrien Gelas, Laurent Besacier, and François Pellegrino. Developments of Swahili resources for an automatic speech recognition system. In *Proc. 3rd Workshop on Spoken Language Technologies for Under-Resourced Languages (SLTU 2012)*, pp. 94–101, 2012.
- Pierre Godard, Gilles Adda, Martine Adda-Decker, Juan Benjumea, Laurent Besacier, Jamison Cooper-Leavitt, Guy-Noël Kouarata, Lori Lamel, Hélène Maynard, Markus Müller, Annie Rialland, Sebastian Stüker, François Yvon, and Marcelly Zanon Boito. A very low resource language speech corpus for computational language documentation experiments. *CoRR*, abs/1710.03501, 2017. URL <http://arxiv.org/abs/1710.03501>.
- Anirudh Gupta, Harveen Singh Chadha, Priyanshi Shah, Neeraj Chimmwal, Ankur Dhuriya, Rishabh Gaur, and Vivek Raghavan. CLSRIL-23: cross lingual speech representations for indic languages. *CoRR*, abs/2107.07402, 2021. URL <https://arxiv.org/abs/2107.07402>.
- Tahir Javed, Sumanth Doddapaneni, Abhigyan Raman, Kaushal Santosh Bhogale, Gowtham Ramesh, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh M. Khapra. Towards building ASR systems for the next billion users. *CoRR*, abs/2111.03945, 2021. URL <https://arxiv.org/abs/2111.03945>.
- Chanwoo Kim and Richard M. Stern. Robust signal-to-noise ratio estimation based on waveform amplitude distribution analysis. In *Interspeech*, 2008.

- Ussen Abre Kimanuka and Osman Büyük. Turkish speech recognition based on deep neural networks. *Süleyman Demirel Üniversitesi Fen Bilimleri Enstitüsü Dergisi*, 2018.
- Julius Kunze, Louis Kirsch, Ilija Kurenkov, Andreas Krug, Jens Johannsmeier, and Sebastian Stober. Transfer learning for speech recognition on a budget. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pp. 168–177, Vancouver, Canada, August 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-2620. URL <https://aclanthology.org/W17-2620>.
- Gideon Mendels, Erica Cooper, Víctor Soto, Julia Hirschberg, Mark John Francis Gales, Kate Knill, Anton Ragni, and Haipeng Wang. Improving speech recognition and keyword search for low resource languages using web data. In *Interspeech*, 2015.
- Josh Meyer. Multi-task and transfer learning in low-resource speech recognition. 2019.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pp. 48–53, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-4009. URL <https://aclanthology.org/N19-4009>.
- Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur. A time delay neural network architecture for efficient modeling of long temporal contexts. In *Proc. Interspeech 2015*, pp. 3214–3218, 2015. doi: 10.21437/Interspeech.2015-647.
- Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukáš Burget, Ondrej Glembek, Nagendra Kumar Goel, Mirko Hannemann, Petr Motlíček, Yanmin Qian, Petr Schwarz, Jan Silovský, Georg Stemmer, and Karel Veselý. The kaldi speech recognition toolkit. 2011.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision, 2022. URL <https://arxiv.org/abs/2212.04356>.
- Sandy Ritchie, You-Chi Cheng, Mingqing Chen, Rajiv Mathews, Daan van Esch, Bo Li, and Khe Chai Sim. Large vocabulary speech recognition for languages of africa: multilingual modeling and self-supervised learning, 2022. URL <https://arxiv.org/abs/2208.03067>.
- Bienvenu Sene-Mongaba. The making of lingala corpus: An under-resourced language and the internet. *Procedia - Social and Behavioral Sciences*, 198:442–450, 2015. ISSN 1877-0428. doi: <https://doi.org/10.1016/j.sbspro.2015.07.464>. URL <https://www.sciencedirect.com/science/article/pii/S1877042815044651>. Current Work in Corpus Linguistics: Working with Traditionally- conceived Corpora and Beyond. Selected Papers from the 7th International Conference on Corpus Linguistics (CILC2015).
- Claytone Sikasote and Antonios Anastasopoulos. Bembaspeech: A speech recognition corpus for the bamba language. *CoRR*, abs/2102.04889, 2021. URL <https://arxiv.org/abs/2102.04889>.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6. URL <https://aclanthology.org/2020.emnlp-demos.6>.
- Alp Öktem. Congolese swahili speech-to-text model. <https://gamayun.translatorswb.org/download/swc-stt-model/>, 2021.
- Alp Öktem, Muhannad Albayk Jaam, Eric DeLuca, and Grace Tang. Gamayun - language technology for humanitarian response. In *2020 IEEE Global Humanitarian Technology Conference (GHTC)*, pp. 1–4, 2020. doi: 10.1109/GHTC46280.2020.9342939.

A APPENDIX

Table 7: Additional experiments on Congolese Swahili test set

Pretraining	Architecture	pretraining	Train/Finetuning	WER
End-to-End (Öktem, 2021)	AED	English	SWC TICO-19	18.3
CdWav2vec	Wav2Vec 2.0	Congolese Radio Corpus	SWC TICO-19	13.7