# Building Text and Speech Benchmark Datasets and Models for Low-Resourced East African Languages: Experiences and Lessons

## 4 1 Introduction

- 5 Africa is home to over 2000 languages, yet most remain absent from the global NLP ecosystem. The
- 6 lack of text and speech datasets has hindered the development of machine translation, sentiment
- 7 analysis, and speech technologies for African contexts. In Uganda alone, more than 41 indigenous
- 8 languages are spoken, but only a few have any digital representation. Addressing this gap is critical
- 9 for inclusive technological development and language preservation. This work focuses on building
- open benchmark datasets and baseline models for five East African languages: Luganda, Runyankore-
- Rukiga, Lumasaba, Acholi, and Swahili through participatory, community-driven approaches.

# 2 Methodology

#### 13 2.1 Data collection

- 14 We adopted a participatory approach, mobilizing linguists, translators, and communities to contribute
- 15 text and speech data. Tools such as Mozilla Pontoon, Common Voice, and the custom Yogera
- application facilitated translation, crowdsourced voice contributions, and validation. Each translation
- 17 underwent linguistic validation by experts. Sentences followed guidelines ensuring readability,
- 18 semantic accuracy, and cultural neutrality. Speech data was reviewed for clarity, speaker diversity,
- 19 and dialectal representation.
- 20 We developed monolingual corpora consisting of approximately 400,000 sentences in Luganda,
- 21 200,000 in Swahili, and about 40,000 sentences each in Acholi, Runyankore-Rukiga, and Lumasaba.
- 22 In addition, a parallel corpus of 40,000 English sentences was translated into each of the five East
- 23 African languages, producing aligned bilingual datasets, alongside sentiment-tagged parallel corpora
- 24 in Luganda and Swahili. For speech data, the collection yielded 582 hours of Luganda audio (438
- 25 hours validated) and 1,100 hours of Swahili audio (393 hours validated), contributed by more than
- 26 1,700 speakers across a broad age range (18–80 years) and multiple dialects, thereby ensuring diversity
- 27 and representativeness.

### 28 3 Baseline Models

- 29 The baseline experiments underscore the utility of the developed datasets for a range of NLP tasks.
- 30 For machine translation, a Marian MT model achieved BLEU scores of 26.0 (English→Luganda)
- 32 Luganda sentences, while topic classification showed strong performance, with SVM (F1 score
- of 0.967) and BERT (F1 score of 0.983) as the best models. Sentiment classification on the
- 34 Luganda corpus also produced high performance with Stacking Classifier (F1 score of 0.90) and
- 35 MultinomialNB (F1 score of 0.88). In speech, Coqui STT trained on 300 hours of Luganda data
- achieved a WER of 23%, with further gains from Wav2vec2, underscoring the benefit of transfer
- 37 learning.

38

#### 4 Conclusion

- 39 This work shows that open, community-driven datasets can enable robust NLP baselines for African
- 40 languages despite challenges of scarcity and diversity. Community participation and transfer learning
- 41 were key to ensuring quality and improving performance. The released benchmark datasets and
- 42 baseline models for five East African languages provide a foundation for inclusive NLP research and
- 43 future cross-lingual applications.