
STATE OF NLP IN KENYA: A SURVEY

A PREPRINT

Cynthia Jayne Amol
Maseno University

Everlyn Asiko Chimoto
University of Cape Town
African Institute for Mathematical Sciences - AIMS

Rose Delilah Gesicho
United States International University-Africa
Zindi

Antony M. Gitau
African Institute for Mathematical Sciences - AIMS

Naome A. Etori
Jomo Kenyatta University of Agriculture and Technology - JKUAT
The University of Minnesota, Twin Cities

Caringtone Kinyanjui
University of Manchester
Songhai

Steven Ndung'u
University of Groningen
Stellenboch University

Lawrence Moruye
Jumia Group

Samson Otieno Ooko
Adventist University of Africa

Kavengi Kitonga
University of Nairobi

Brian Muhia
Independent Researcher

Catherine Gitau
African Institute for Mathematical Sciences - AIMS

Antony Ndolo
Deep Learning Indaba
Karadeniz Technical University

Lilian D. A. Wanzare
Maseno University

Albert Njoroge Kahira
Datawise Africa

Ronald Tombe
Kisii University
Future Africa - University of Pretoria

October 15, 2024

ABSTRACT

Kenya, known for its linguistic diversity, faces unique challenges and promising opportunities in advancing Natural Language Processing (NLP) technologies, particularly for its underrepresented indigenous languages. This survey provides a detailed assessment of the current state of NLP in Kenya, emphasizing ongoing efforts in dataset creation, machine translation, sentiment analysis, and speech recognition for local dialects such as Kiswahili, Dholuo, Kikuyu, and Luhya. Despite these advancements, the development of NLP in Kenya remains constrained by limited resources and tools, resulting in the underrepresentation of most indigenous languages in digital spaces. This paper uncovers significant gaps by critically evaluating the available datasets and existing NLP models, most notably the need for large-scale language models and the insufficient digital representation of Indigenous languages. We also analyze key NLP applications—machine translation, information retrieval, and sentiment analysis—examining how they are tailored to address local linguistic needs. Furthermore, the paper explores the governance, policies, and regulations shaping the future of AI and NLP in Kenya and proposes a strategic roadmap to guide future research and development efforts. Our goal is to provide a foundation for accelerating the growth of NLP technologies that meet Kenya's diverse linguistic demands.

Table 1: Main language groups in Kenya with Number of speakers

Bantu	Nilotic	Cushitic
Kikuyu 8.1 million		
Kamba 4.7 million	Dholuo 5.0 million	
Luhya 10 million	Kalenjin languages 4.6 million	Oromo (over 48 million incl. Ethiopia)
Gusii 2.7 million		Borana, 3.4 million speakers in 2010
Meru 2.0 million	Maasai 1.2 million	Orma, 659,000 speakers in 2015
Mijikenda/Giriama	Turkana 1.0 million	Somali 2.8 million (22 million incl. Ethiopia and Somalia)
ca. 1 million		

Keywords NLP, low-resource languages, Kenya, Kenyan languages, Kiswahili, datasets, machine translation, sentiment analysis, speech recognition, AI governance.

1 Introduction

The last few years have seen an explosion in the number of Natural Language Processing(NLP) applications and tools. Tools such as chatGPT are now ubiquitous in many aspects of our lives. NLP subfield of artificial intelligence (AI) that seeks to enable machines to understand, interpret, and generate human language. Although there have been significant advances in the field of NLP, most of the research today focuses on less than 1% of the world’s languages with a particular focus on ten or so languages, including English Joshi et al. (2020); Bender (2019). One of the main challenges of NLP, especially in low-resource settings, is the availability of requisite datasets to train models in low-resource languages. This is a big issue regarding the progress of NLP for low-resource languages Siminyu et al. (2021); Ruder and Korashy (2019). With the increasing number of research interests in NLP globally, multiple efforts have emerged in Africa to create both NLP tools and datasets that power NLP technology, reflecting the continent’s drive toward AI democratization while addressing concerns about digital colonization Etori et al. (2024) and the need for locally grounded solutions.

Africa is one of the most linguistically diverse regions in the world, with an estimated 3,000 languages spoken across the continent. However, efforts to forge cohesive nation-states from multiple ethnic groups have often resulted in the marginalization of several languages. Priority has typically been given to high-resource languages such as English and French, alongside regional or national languages like Kiswahili and Afrikaans Adebara (2024); Yan and Xu ([n. d.]). While these dominant languages are widely spoken, millions of Africans communicate in at least three languages, including one or two national and indigenous languages. Unfortunately, many of these indigenous languages remain underrepresented in the digital space despite their prevalence.

In Kenya, there has been limited progress in mapping the country’s languages and developing digital tools to support them, particularly for the indigenous languages spoken by millions. This highlights the need for a comprehensive assessment of the current state of NLP technologies in Kenya, considering its diverse geographical and linguistic contexts. Many of these languages remain underrepresented in digital spaces, underscoring the urgency for targeted efforts to bridge this gap and accurately reflect NLP advancements in the region.

Kenya is made up of multilingual communities. First, it is important to distinguish the difference between tribes and languages. A widely cited but controversial number is the 42 tribes. One would expect, therefore, that Kenya would have about 42 languages based on the tribes. However, the total number of tribes in Kenya is not well researched, neither are there clear distinctions on what is a tribe. According to ethnology Gordon Jr (2005), there are about 68 spoken languages in Kenya today. Most of the research attempts to study language groups instead of individual languages. Kenyan languages are generally fall into 3 groups; Bantu, Nilotic and Cushitic. Table 1 shows the most widely spoken languages in their language groups. With a population of 50 million and growing, doubled by a recent emergence of language and culture appreciation, the roadmap in this paper provides a path towards digital representation for this population.

Some efforts have started in Kenya to collect, clean and organise language datasets in several indigenous languages. These datasets have already been used to develop NLP tools in several of Kenyans languages¹. With this backdrop, this paper surveys on the state of NLP in Kenya, aiming to highlight ongoing work and bring out the challenges and

¹We highlight such efforts in subsequent sections of the paper

opportunities in the field. As these efforts continue, it is important to record and map them and create a road map for the future. This survey is an initial attempt at this.

To the best of our knowledge, this survey is the first attempt to consolidate all efforts by researchers, practitioners and linguists and provide a reference point for NLP research in Kenya. It is also the first major attempt to map Kenyan languages from an NLP perspective and provide an initial perspective of a prospective research road map for NLP research in Kenya. While the field of NLP is very broad and keeps growing, we aim to cover as many subtopics as possible.

The remainder of this paper is organised as follows; first, we look at similar surveys and works in Section 2. In Section 4, we survey all datasets on Kenyan languages. In section 5 we survey NLP applications developed in Kenya or using Kenyan languages. We touch on categories of applications such as machine translations, information retrieval, and text classification. In section 6 we look at governance, policies and regulations related to AI and NLP specifically. We highlight different efforts by governments and communities to draft AI policies and regulations around the technology. Finally in section 8 we discuss a possible roadmap for NLP research in Kenya.

2 Related Work

The systematic review by Chesire and Kipkebut (2024) highlights an expanding corpus of research in dataset creation, machine translation, and the development of multilingual pretrained models tailored for African languages. Prominent initiatives include the development of datasets such as MASAKHANEWS, which facilitates news classification in 16 African languages Adelani et al. (2023b), and AfriSenti, designed for sentiment analysis across 14 African languages Muhammad et al. (2023a). Advancements in machine translation are evident, with models like AfriByT5 Adelani et al. (2022) and MMTAfrica Emezue and Dossou (2022) catering to the linguistic diversity of the continent. Nonetheless, the authors point out that the African NLP landscape still confronts considerable hurdles, chiefly the dearth of extensive datasets and a limited number of NLP researchers within the region.

Adebara and Abdul-Mageed (2022) assess the Afrocentric NLP for the African languages. In their paper, they explore the linguistic diversity of African languages and their distinct challenges for NLP, including tonal systems, vowel harmony, and serial verb constructions. It emphasizes that these linguistic traits are underrepresented in most mainstream NLP languages, underscoring the necessity for specialized methodologies tailored to African languages. Furthermore, the paper sheds light on various sociopolitical factors, such as national language policies that prioritize foreign languages, low literacy rates in indigenous languages with individuals often only literate in foreign tongues, and the absence of standardized orthographies or the presence of inconsistent spelling norms, all of which add complexity to text processing and analysis. These issues profoundly affect the advancement of NLP for African languages.

There have been attempts to survey the state of NLP in Africa from regional and country-specific perspectives. From a regional perspective, the study by Mussandi and Wichert (2024) discusses some of the opportunities and challenges of building technologies for languages of African origin in their survey on NLP tools for African languages. It features corpora and task-specific language models developed by previous studies for these languages, exposing the scarcity of these NLP tools that contribute to African languages' 'technological delay'. AI4D - Artificial Intelligence for Development initiative Siminyu et al. (2021) noted the technological gap created by lack of datasets in African languages. This initiative, aimed at creating datasets for African languages through crowd-sourcing data, reported the setbacks of creating data resources for these languages. Hedderich et al. (2020) review NLP approaches for low-resource languages touching on issues that affect African languages, where most are low-resource.

From country specific perspectives, Azunre et al. (2021) in their study on the state of NLP in Ghana highlight some of the efforts of NLP Ghana in developing data sources and language tools for Ghanaian languages that are considered low-resource. Some of their contributions such as building annotated datasets, embedding models and translators and for most-widely spoken languages in Ghana are aimed at increasing availability of language resources for these languages. Marivate (2020) highlights the South African NLP landscape, zooming into the existing disparity when it comes to availability of language content in the various South African languages on Wikipedia. The study also looks at past experiences of language processing for South African languages and paints a way forward for African NLP through community building.

3 Methodology

This research follows a multi-step approach encompassing data collection, analysis of existing models, and evaluation of available resources for Kenyan languages to assess the current state of NLP technologies in Kenya. Our methods

involved reviewing locally and internationally published research and datasets, as long as they incorporated Kenyan languages, ensuring a comprehensive analysis.

We gathered relevant data by extensively reviewing publicly available datasets, research papers from Google Scholar, ACL, etc., and NLP tools. We focused on Kenyan languages such as Swahili, Dholuo, and Luhya. Our survey included datasets for text classification, machine translation, sentiment analysis, and question answering.

Lastly, we compared the available datasets and models with Kenya’s languages and identified significant gaps in data availability and technological adaptation. Our comparison followed a structured approach focused on five key areas: **language coverage**, **NLP tasks**, **NLP tools**, **NLP application tasks**, and **resource availability**. We assessed the representation of Kenyan languages in available datasets. The datasets were analyzed using the application of NLP tasks. We examined the practical application of these tools in functions such as translation and speech recognition. Lastly, we compared the availability of computational resources, including annotated datasets and pre-trained models, across Kenyan languages, identifying critical gaps.

4 Datasets

This section provides an overview of publicly available datasets in the specified languages, focusing on their multilingualism characteristics, encompassing speech and text modalities. The datasets are further analyzed based on their applicability to various downstream tasks, highlighting their relevance for computational linguistics and machine learning applications.

4.1 Speech Datasets

(a) Kenyan Languages Corpus for Machine Learning and Natural Language Processing (Kencorpus):

KenCorpus² Wanjawa et al. (2023b) is a comprehensive text and speech corpus for three Kenyan languages: **Swahili**, **Dholuo** and **Luhya (covering the Lumarachi, Logooli, and Lubukusu dialects)**. The dataset was collected and curated using contributions from native speakers across diverse sources, including language communities, schools, media outlets, and publishers. KenCorpus provides a rich resource for machine learning (ML) and NLP applications; specifically, it comprises 2,585 texts (approximately 1.8 million words) and 19 hours of speech for Swahili, 546 texts (about 1.3 million words) and 99 hours of speech for Dholuo, and 987 texts (about 2.2 million words) and 58 hours of speech for the three Luhya dialects. Additionally, the dataset includes KenSpeech, a transcribed corpus of roughly 27 hours of Swahili speech Awino et al. (2022). KenCorpus also features Part-of-Speech (POS) tagged sentences, covering approximately 50,000 words for Dholuo and 90,000 words for Luhya, further enhancing its utility for linguistic and computational research.

(b) Building African Voices:

The Building African Voices³ Ogayo et al. (2022) is a curated text and speech dataset with 16 African languages featuring 4 Kenyan languages: **Dholuo**, **Suba** and **Kenyan English**. The dataset contains 13,897 utterances in Luo, 2,078 in Suba, and 1,150 in Kenyan English. Data in this corpus was scraped from books, websites, and social media posts.

(c) CMU Wilderness Multilingual Speech Dataset:

CMU Wilderness⁴ Black (2019) is a multilingual speech dataset containing several audio and textual data collected from the Bible. The dataset comprises audio in 700 different languages, with three (3) Kenyan languages represented, **Oromo**, **Somali** and **Sabaot**. On average, each dataset contains about 20 hours of aligned sentence-level text and word pronunciations.

²<https://kencorpus.maseno.ac.ke/corpus-datasets/>

³<https://www.africanvoices.tech/>

⁴http://festvox.org/cmu_wilderness/

Table 2: Published Kenyan Datasets for Text and Speech

#	Dataset Name	Available	Venue	Language
1	Kencorpus: Kenyan Languages	ML/NLP Dataset Wanjawa et al. (2023b)	Maseno	Swahili Luol Luhya
2	Building African Voices	Speech synthesis Dataset Ogayo et al. (2022)	Interspeech	Subal Luol Swahili
3	CMU Wilderness	Multilingual Speech Black (2019)	IEEE	Oromol Somali Sabaot
4	Bible TTS	Speech Dataset Meyer et al. (2022)	Interspeech	Kikuyul Luo
5	Building Low Resource Datasets	Text and Speech Dataset Babirye et al. (2022)	AfricaNLP	Kiswahili
6	XTREME-S	Cross-lingual Dataset Conneau et al. (2022)	Interspeech	Luol Somali Swahili etc
7	The African Story Book(ASb)	multilingual African Stories Stranger-Johannessen and Norton (2017)	SA	Kisii Kikuyul Swahili etc
8	Common Voice	Multilingual Dataset Ardila et al. (2019)	Mozilla	Luol Swahili Somali etc
9	IARPA Babel	Telephone speech Bills et al. (2022)	Appen	Dholuo
10	Kiswahili TTS dataset	TTS Dataset Rono (2021)	Mendeley	Kiswahili
11	Swahili audio mini-kit	Audio Dataset	TWB	Kiswahili
12	Bloom-lm	Audio Dataset Leong et al. (2022)	SIL-AI	Kiswahili Kambal Luo
13	AfriSpeech-200	ASR Dataset Olatunji et al. (2023)	TACL	Kenyan-English
14	AfroDigits	Speech Dataset Chinenye Emezue et al. (2023)	AfricaNLP	Oromo
15	1000 African Voices	Speech synthesis Dataset Ogun et al. (2024)	Interspeech	Kenyan-eng
16	Helsinki Corpus	Text Dataset Hurskainen (2004)	Kielipankki	Kiswahili
17	MasakhaNEWS	Text Classification Dataset Adelani et al. (2023b)	IJCNLP	Kiswahili
18	AFRIHG	Text Summarization Dataset Ogunremi et al. ([n. d.]	AfricaNLP	kiswahili Somali Oromo
19	Global Voices	Text Summarization Dataset Nguyen and Daumé III (2019)	EMNLP	kiswahili
20	Glott500	Scaling Multilingual Dataset ImaniGooghari et al. (2023)	ACL	kiswahili Somali Luo
21	LORELEI	disaster incidents Dataset Tracey et al. (2019)	DARPA	kiswahili Somali Oromo
22	SIB-200	Topic classification Dataset Adelani et al. (2023a)	ACL	kiswahili Kikuyul Dholuo
23	Kamba POS Tagger Memory Based	POS Dataset Kituku et al. (2015)	IJNLC	kamba
24	OSCAR	Common Crawl Dataset Abadji et al. (2022)	ACL	Kiswahili Somali
25	Language modeling	Language modeling Dataset Shikali and Refuoe (2019)	OSCAR	Kiswahili
26	Spell-checker for Gikuyu	Spell-checker Dataset Chege et al. (2010)	UON	Kikuyu
27	CommonCrawl	cross-lingual Dataset Conneau et al. (2019)	ACL	Kiswahili Somali Oromo
28	AI4D -African Language Program	cross-lingual Dataset Siminyu et al. (2021)	IDRC	Kiswahili
29	AfriMTE and AfriCOMET	human evaluation Dataset Wang et al. (2023)	NAACL	Kiswahili

(d) Bible TTS:

The Bible TTS⁵ is a speech dataset for 10 languages spoken in Sub-Saharan Africa, featuring **Kikuyu** and **Dholuo** languages Meyer et al. (2022). The dataset includes Bible recordings released by the Open.Bible project⁶. The corpus contains up to 86 hours of high quality, aligned single speaker recordings per language.

⁵<https://masakhane-io.github.io/bibleTTS/>

⁶<https://bibleproject.com/>

(e) AfriSpeech-200

AfriSpeech Olatunji et al. (2023), a Pan-African accented English speech dataset for clinical and general domain ASR, crowdsourced from 2,463 African speakers, 200.70 hrs with an average audio duration of 10.7 seconds. Kenya contributes 8,304 clips from 137 speakers, totaling 20.89 hours of audio. This data helps improve ASR performance for Kenyan-accented English, ensuring more accurate and inclusive speech recognition tools that cater to local healthcare systems and other domains where speech recognition can enhance productivity and accessibility.

(f) AfroDigits:

This is a Community-Driven Spoken Digit Dataset⁷ Chinenye Emezue et al. (2023) for African. It is an openly available dataset of spoken digits in 38 African languages, including **Oromo** and **Borana**, spoken widely in the northern region of Kenya.

(g) 1000 African Voices:

The Afro-TTS dataset Ogun et al. (2024) is a collection of English speech recordings featuring diverse African accents, curated through crowdsourcing. It includes 136 hours of audio from 747 contributors across 9 countries, with 86 different accents. The Kenyan subset comprises 5,307 samples from 58 speakers, highlighting the unique linguistic characteristics of Kenyan English. This open-source dataset aims to support research on African-accented English.

(h) Building Text and Speech Datasets for Low Resourced Languages: A Case of Languages in East Africa:

This corpus Babirye et al. (2022) contains text and speech data in 5 languages including **Kiswahili**. Data was collected from contributions by the community and pre-existing text data sources: news websites, published reports, storybooks and open-source Kiswahili datasets. The Kiswahili dataset has over 200K sentences and 100 hours of voice data crowdsourced using Mozilla Commonvoice⁸.

(i) XTREME-S: Evaluating Cross-lingual Speech Representations:

The XTREME-S project⁹ Conneau et al. (2022) created a speech dataset covering 102 languages, including **Somali**, **Swahili**, **Dholuo** and **Kamba** from Kenya. **HOW was is created, QUANTITY**

(j) The African Storybook Project (ASb) :

ASb¹⁰ Stranger-Johannessen and Norton (2017) is a children’s literacy project by the South Africa Institute of Distance Education¹¹. The project is a collection of 4,317 storybooks, poems, songs, rhymes, and picture books in 242 African languages including **Ekegusii**, **Gikuyu**, **Dholuo**, **Turkana**, **Somali**, **Oromo**, **Maasai**, **Samburu** and **Kipsigis** languages spoken in Kenya. The textual data is organized according to the size of words and paragraphs.

(k) Mozilla Common Voice :

Common voice¹² Ardila et al. (2019) is the largest open-source, multi-language speech dataset. The text and speech datasets are collected and validated through crowdsourcing. It contains datasets in several languages, including 973 hrs of speech and 101,669 sentences in **Swahili**. This constantly scaling dataset also has 266 sentences in **Somali**.

(k) IARPA Babel Dholuo Language Pack:

This is a speech dataset¹³ Bills et al. (2022) containing approximately 204 hours of audio data in South Nyanza and Trans-Yala **Dholuo** dialects. The data is sourced from telephone speech.

⁷<https://huggingface.co/datasets/chrisjay/crowd-speech-africa>

⁸<https://commonvoice.mozilla.org/en>

⁹https://hf.co/datasets/google/xtreme_s

¹⁰<https://www.africanstorybook.org/>

¹¹<https://www.devex.com/organizations/south-african-institute-for-distance-education-saide-22060>

¹²<https://commonvoice.mozilla.org/en>

¹³<https://abacus.library.ubc.ca/dataset.xhtml?persistentId=hdl:11272.1/AB2/HSAU9N>

(m) A Kiswahili TTS dataset:

This dataset¹⁴ Rono (2021), available on Mendeley Data, contains 1,570 text files (23,487 words) and 1,570 audio files in **Kiswahili**. The textual data was sourced from newspaper articles, stories, and novels.

(n) Swahili audio mini-kit:

This dataset¹⁵ contains 4,700 samples from **Swahili** mini-kit recorded by a Kenyan male speaker and their transcriptions.

(o) Bloom-lm:

The Bloom library¹⁶ Leong et al. (2022) has stories in 363 languages including the Kenyan languages **Taveta, Samburu, Swahili, Rendile, Ekegusii, Nyole, Turkana, Suba, Kikuyu, Oromo, Okiek, Kitharaka, Marakwet, Bukusu, Kamba, Pokomo** and **Meru**. The library, which has a mean of 32 stories and a median of 2 stories per language, was created by SIL International¹⁷ to empower communities with low-resource languages to create literature for children.

4.2 Text Datasets**(a) Helsinki Corpus of Swahili (HCS):**

HCS Hurskainen (2004) is the most refined **Swahili** corpus containing 25 million words collected from news sources, stories, and legislative assemblies. The corpus is available in the Language Bank of Finland (Kielipankki) in two versions: annotated and not annotated. The not-annotated version is available openly via <https://korp.csc.fi/download/HCS/na-v2/>. In contrast, the annotated version, which contains words annotated with individual **lemma, Part of Speech (PoS), morphological and syntactic tags** Steimel et al. (2023) is access-restricted via <http://urn.fi/urn:nbn:fi:lb-201608301..>

(b) MasakhaNEWS:

MasakhaNEWS Adelani et al. (2023b) is the largest news classification dataset covering 16 African languages, including three languages in Kenya. The dataset has 7782, 6431, and 2915 news articles in **Oromo, Swahili, and Somali**, respectively. The news articles were analyzed in six categories: business, entertainment, health, politics, sports, and technology.

(c) AFRIHG:

AFRIHG Ogunremi et al. ([n. d.]) is the largest abstractive summarization dataset for headline generation covering 16 African languages, including three languages in Kenya. The dataset has 11,998, 17,931, and 20,276 new articles and their respective headlines in **Somali, Oromo, and Swahili**, respectively. The dataset was scrapped from BBC News and incorporates news articles from Adelani et al. (2023b).

(d) Global Voices

Nguyen and Daumé III (2019) is a multilingual collection for evaluating cross-lingual summarization methods. It includes articles and their summaries in 15 languages. The dataset contains news articles and summaries in **Swahili**. The summaries are gathered automatically from social network descriptions and manually through crowdsourcing, ensuring quality through human ratings.

(e) Glot500: Scaling Multilingual Corpora and Language Models to 500 Language:

The Glot500¹⁸ project ImaniGooghari et al. (2023) is an effort to scale NLP to support as many of the world’s languages and cultures as possible. It created an LLM trained on 700GB of text in 2266 low-resource ‘tail’ languages, including **Swahili, Dholuo, Somali** and **Oromo** from Kenya. The text data was collected by crawling websites and compiling data from around 150 different datasets.

¹⁴<https://data.mendeley.com/datasets/rbn6nmygcn/1>

¹⁵<https://gamayun.translatorswb.org/download/swahili-audio-mini-kit/>

¹⁶<https://huggingface.co/datasets/sil-ai/bloom-lm>

¹⁷<https://www.sil.org/>

¹⁸<https://github.com/cisnlp/Glot500>.

(f) Corpus Building for Low Resource Languages in the DARPA LORELEI Program:

The LORELEI (Low Resource Languages for Emergent Incidents)¹⁹ program Tracey et al. (2019) developed large volumes of both monolingual and parallel language packs to improve technologies capable of providing awareness in disasters or emergent incidents. The two language packs, representative and incident, focus on low-resource languages and cover 23 languages, including **Oromo, Somali** and **Swahili**. The project’s target was over 1 million words for the parallel datasets and over 2 million for the monolingual datasets. Data was sourced from formal news, social media, blogs, discussion forums, and reference materials such as Wikipedia. The textual data contains semantic, morphosyntactic, entity, and parallel annotations.

(g) SIB-200: A Simple, Inclusive, and Big Evaluation Dataset for Topic Classification in 200+ Languages and Dialects :

SIB-200²⁰ Adelani et al. (2023a) is a benchmark topic classification dataset based on the FLORES-200 corpus that covers over 200 languages and dialects. The dataset contains 1,004 sentences and is annotated at sentence level in seven categories: science/technology, travel, politics, sports, health, entertainment, and geography. It contains data from Kenya in **Swahili, Gikuyu** and **Dholuo** languages.

(h) Kamba Part- of-Speech Tagger Using Memory Based Approach:

This corpus contains approximately 30K words manually annotated with PoS tags. The words in the dataset were collected from online sources and documents written in **Kikamba** language Kituku et al. (2015).

(i) OSCAR: Open Super-large Crawled Aggregated corpus:

The OSCAR project²¹ Abadji et al. (2022) provides open-source resources and large, unannotated datasets. The current version contains web data in 166 languages, including 1,670 documents and 164,510 words in **Swahili** and **Somali**.

(j) Language modeling data for Swahili:

This dataset²² contains 28,000 unique words and a total of 9.81million **Swahili** words Shikali and Refuoe (2019). It was developed specifically for language modeling tasks.

(k) Developing an Open source Spell-checker for Gikuyu:

This corpus Chege et al. (2010) is a collection of 19,000 words derived from pre-existing datasets in **Gikuyu**. It contains text from religious material, poems, short stories, novels, and the Internet. The data was manually annotated into parts of speech (PoS).

(l) Unsupervised Cross-lingual Representation Learning at Scale:

The study Conneau et al. (2019) resulted in CommonCrawl²³ Corpus²⁴in 100 languages including **Oromo, Swahili** and **Somali**. This large corpus has over 2.5GB of data with 8 million tokens in Oromo, 275 million in Swahili, and 62 million in Somali.

(m) Building a database for Kiswahili language in Africa:

This document classification dataset contains 10K instances in **Swahili** and is a product of AI4D – African Language Program ²⁵Siminyu et al. (2021) which is working towards boosting the integration of African languages on digital platforms.

¹⁹<https://www.darpa.mil/program/low-resource-languages-for-emergent-incidents>

²⁰<https://github.com/dadelani/sib-200>

²¹<https://oscar-project.org/>

²²<https://zenodo.org/records/3553423>

²³<https://commoncrawl.org/>

²⁴[https://github.com/facebookresearch/\(fairseq-py,pytext,xlm\)](https://github.com/facebookresearch/(fairseq-py,pytext,xlm))

²⁵<https://www.k4all.org/project/language-dataset-fellowship/>

(n) AfriMTE and AfriCOMET: Enhancing COMET to Embrace Under-resourced African Languages:

Wang et al. (2023) enhances machine translation (MT) evaluation for under-resourced African languages. It addresses the limitations of n-gram metrics like BLEU and introduces AfriCOMET, an evaluation metric with a higher correlation to human judgements. The dataset includes high-quality human evaluation data using simplified guidelines for error detection and direct assessment (DA) scoring across 13 diverse African languages, including Kenya; the language used is **Swahili**.

4.3 Parallel Corpora

Table 3: Published Kenyan Parallel, Question and Answer, Sentiment Datasets, Hate Speech and NER Datasets

#	Dataset Name	Available	Venue	Language
1	ParaCrawl: Web-Scale Acquisition	Parallel Dataset Bañón et al. (2020)	ACL	Swahili Somalia
2	MAFAND-MT	MT Parallel Dataset Adelani et al. (2022)	ACL	Swahili Luo
3	WikiMatrix	Parallel Dataset Schwenk et al. (2019)	ACL	Swahili
4	KenTrans	Parallel Dataset Wanjawa et al. (2023b)	Maseno	Swahili Luo Luhyia
5	PanLex	lexemes pairwise Dataset Kamholz et al. (2014)	LREC	Kiswahili Luhyia etc
6	FLORES-200	MT Parallel Dataset Conneau et al. (2022)	ACL	Swahili Kikuyuletc
7	NLLB	multilingual African Stories Costa-jussà et al. (2022)	Meta	Kikuyul Swahililetc
8	CCAligned	Parallel Dataset El-Kishky et al. (2020)	ACL	Swahili Oromo
9	GoURMET	Parallel speech Espla-Gomis et al. (2019)	ACL	Swahili
10	The SAWA	Parallel Dataset De Pauw et al. (2009)	EACL	Swahili
11	A Knowledge-Light Approach	Trilingual Dataset De Pauw et al. (2010)	LREC	Luo
12	Tatoeba	Parallel Dataset Raine (2018)	Helsinki	Swahili Somali Rendile
13	Very LR Sentence Alignment	Parallel Dataset Chimoto and Bassett (2022)	ACL	Swahili Luhyia
14	English–Bukusu Automatic MT	Parallel Dataset Ngoni (2022)	UON	Bukusu Kenyan-eng
15	TICO-19	Pairwise Dataset Anastasopoulos et al. (2020)	ACL	Swahili Somali
16	Tanzil	Parallel Dataset ²⁶	Quran	Swahili Somali
17	AfriCLIRMatrix	I and R Dataset Ogundepo et al. (2022)	EMNLP	Swahili Somali
18	CIRAL	I and R Dataset Adeyemi et al. (2024)	ACM	Swahili
19	KenSwQuAD	Q and A Dataset Wanjawa et al. (2023b)	ACM	kiswahili
20	AfriQA	cross-lingual Q and A Dataset Ogundepo et al. (2023)	EMNLP	Swahili
21	TYDI	Q and A Dataset Clark et al. (2020)	TACL	Swahili
22	AfriSenti	Sentiment Dataset Muhammad et al. (2023a)	EMNLP	Swahili Oromo
22	RideKE	Sentiment Dataset Shikali and Refuoe (2019)	ACL	Swahili Sheng
23	Hate Speech	Hate speech Dataset Ombui et al. (2019)	IEEE	Swahili Kenyan-eng
24	XTREMESPEECH	Hate speech Dataset Maronikolakis et al. (2022)	ACL	Swahili Kenyan-eng
25	MasakhaNER	NER Dataset Adelani et al. (2021)	TACL	Swahili Luo
26	Mining Wikidata	NER Dataset Sälevä and Lignos (2021)	AfricaNLP	Swahili Oromo
25	Cross-lingual Tagging and Linking	NER and entity linking Dataset Pan et al. (2017)	ACL	Swahili Kikuyul Somali

(a) ParaCrawl: Web-Scale Acquisition of Parallel Corpora:

ParaCrawl²⁷ Bañón et al. (2020) is the largest parallel corpora curated by crawling the websites. It comprises 223 million unique sentence pairs in 23 languages, including 132,517 and 14,879 sentences in **Swahili** and **Somali**, respectively.

(b) Masakhane Anglo & Franco Africa News Dataset for Machine Translation:

MAFAND-MT Adelani et al. (2022) is a translated news corpus of 16 African languages, including **Swahili** and **Luo** sourced from news websites from local newspapers. The dataset contains 31K and 872K sentences in Swahili and Luo, respectively.

(c) WikiMatrix: Mining 135M Parallel Sentences in 1620 Language Pairs from Wikipedia:

Wikimatrix²⁸ Schwenk et al. (2019) is a parallel corpus mined from Wikipedia containing 135M parallel sentences for 1620 different language pairs. The textual data is in 85 languages, including **Swahili**.

(d) KenTrans:

On top of the individual monolingual datasets, the Kencorpus Wanjawa et al. (2023b) project has parallel corpora between **Swahili** and **Dholuo** (1,500 sentences) and between Swahili and Luhya (11,900 sentences).

(e) PanLex: Building a Resource for Panlingual Lexical Translation:

PanLex²⁹ Kamholz et al. (2014) is a documentation of 20million lexemes in 9,000 language varieties . It covers several Kenyan Languages, including **Swahili**, **Somali**, **Samburu**, **Kipsigis**, **Makonde**, **Turkana**, **Oromo**, **Kuria**, **Suba** and the **Luhya** dialects Bukusu and Lulogooli.

(f) FLORES-200 Evaluation Benchmark for Low-Resource and Multilingual Machine Translation:

The FLORES 200³⁰ Goyal et al. (2022) is a high-quality dataset consisting of 3,001 sentences mined from English Wikipedia and professionally translated in 204 languages. Some Kenyan languages covered by the dataset are **Swahili**, **Gikuyu** and **Dholuo**.

(g) No language left behind: Scaling human-centered machine translation:

NLLB³¹ Costa-jussà et al. (2022) The dataset has texts from **Somali**, **Oromo**, **Kamba**, **Gikuyu**, **Dholuo** and **Swahili** Kenyan languages.

(h) CCAligned: A Massive Collection of Cross-lingual Web-Document Pairs:

CCAligned corpus El-Kishky et al. (2020) comprises parallel web document and sentence pairs in 137 languages, including the Kenyan languages **Swahili** and **Oromo** aligned with English. The data contains over 100 million aligned documents created by performing language identification on raw web documents and can be used for machine translation tasks.

(i) Global Under - Resourced MEDIA Translation (GoURMET):

The GoURMET³² project curated a parallel dataset Espla-Gomis et al. (2019) in 16 languages including 3.5 million sentences in **Swahili**. The data was sourced from web crawls. Language pairing is done to and from English, depending on the task.

²⁷<https://paracrawl.eu/>

²⁸<https://github.com/facebookresearch/LASER/tree/main/tasks/WikiMatrix>

²⁹<http://panlex.org>

³⁰<https://github.com/facebookresearch/flores>

³¹<https://github.com/facebookresearch/fairseq/tree/nllb>

³²<https://gourmet-project.eu/>

(j) The SAWA Corpus: a Parallel Corpus English - Swahili

SAWA corpus De Pauw et al. (2009) is a parallel corpus of English - Swahili curated to bootstrap a data-driven machine translation system for **English - Swahili**. The corpus comprises 542.1K and 442.9K words in English and Swahili, respectively. The textual data was collected from various sources, including the New Testament section of the Bible, Quran, kamusi (dictionary), reports, translators and movie subtitles.

(k) A Knowledge-Light Approach to Luo Machine Translation and Part-of-Speech Tagging:

This study De Pauw et al. (2010) resulted in a trilingual corpus **English - Swahili - Luo (Dholuo)** curated by utilizing the New Testament data of the SAWA corpus De Pauw et al. (2009) to construct a trilingual parallel corpus (English - Luo - Swahili). The dataset contains 192K 156K and 170K token counts in English, Swahili and Luo respectively. Data is annotated with Part of Speech (PoS) tags.

(l) Tatoeba:

Tatoeba³³ is a collection of over 11million sentences and translations in 446 languages including **Swahili, Somali** and **Rendile** Raine (2018).

(m) Very Low Resource Sentence Alignment: Luhya and Swahili:

This is the first digital parallel corpus in Luhya-English for the **Marama** dialect. The corpus comprises 7,952 parallel sentences in Marama generated by aligning the Bible’s New Testament in Luhya and English Chimoto and Bassett (2022).

(n) English–Bukusu Automatic Machine Translation for Digital Services Inclusion in E-governance:

This study resulted in a parallel **English-Bukusu** dataset containing 5,146 sentences. The data was collected from three sources: the New Testament section of the Bukusu version of the Bible, electronic texts in Bukusu and the English-Bukusu dictionary Ngoni (2022).

(o) Translation Initiative for COVID-19:

TICO-19 Translation Benchmark³⁴ is a translation of COVID-19-related terms from English to various languages. This benchmark aims to include 30 documents (3071 sentences, 69.7k words) translated from English into 36 languages, including **Somali** and **Swahili**. Anastasopoulos et al. (2020)

(p) Tanzil Dataset :

Tanzil³⁵ is a collection of Quran translations in 42 languages including **Swahili** and **Somali**.

4.4 Question Answering Datasets**(a) AfriCLIRMatrix:**

AfriCLIRMatrix³⁶ Ogundepo et al. (2022) is an information retrieval (Question Answering) test collection comprising queries and documents in 15 African languages including **Swahili** and **Somali**. This dataset, sourced from Wikipedia, contains 9860 documents from Somalia and 70808 from Swahili.

(b) Cross-lingual information retrieval:

CLIR³⁷ Adeyemi et al. (2024) is an information retrieval (Question Answering) test collection comprising queries and documents in 15 African languages including **Swahili** and **Somali**. This dataset, sourced from Wikipedia, contains 9860 documents from Somalia and 70808 from Swahili.

³³<https://opus.nlpl.eu/Tatoeba.php>

³⁴<https://tico-19.github.io/>

³⁵<https://tanzil.net/trans/>

³⁶<https://github.com/castorini/africlirmatrix>

³⁷<https://github.com/ciralproject/ciral>

(c) KenSwQuAD—A Question Answering Dataset for Swahili Low-resource Language:

KenSwQuAD Wanjawa et al. (2023a) is a Question Answering dataset containing about 1,440 **Swahili** texts labeled with at least 5 questions each, totaling about 7,526 QA pairs. The dataset is based on the KenCorpus Wanjawa et al. (2023b) project.

(d) AfriQA: Cross-lingual Open-Retrieval Question Answering for African Languages:

AfriQA³⁸ Ogundepo et al. (2023) is the first cross-lingual QA dataset containing 10 African languages, including 1,134 questions in **Swahili**. Native speakers of each language were tasked with data collection and annotation. The process involved question elicitation, translation into a pivot language (English or French), answer labelling and answer translation back to the source language.

(e) TYDI QA: A Benchmark for Information-Seeking Question Answering in Typologically Diverse Languages:

TYDI QA³⁹ Clark et al. (2020) is a QA dataset covering 11 languages including **Kiswahili**. The dataset contains 204K QA pairs collected by generating questions based on short prompts from Wikipedia articles and then pairing each question with a Wikipedia article.

4.5 Sentiment Analysis**(a) AfriSenti:**

Afrisenti Muhammad et al. (2023a) is a Twitter sentiment analysis dataset consisting of annotated textual data in 14 African languages, including **Swahili** and **Oromo** spoken in Kenya. The dataset contains 3,014 and 2,491 texts labelled as negative or positive in Swahili and Oromo, respectively.

(b) RideKE:

RideKE Etori and Gini (2024) is a Twitter-based corpus containing over 29,000 code-switched entries in Kenyan-accented English, Swahili, and Sheng. It is designed explicitly for sentiment and emotion analysis within the ride-hailing service domain. The dataset includes 553 labelled entries for supervised training, 2,000 human-annotated entries for testing, and over 27,000 unlabeled entries used in a semi-supervised learning loop.

4.6 Hate Speech**(a) Hate Speech Detection in Code-switched Text Messages:**

The study by Ombui et al. (2019) curated a dataset with 260k tweets in **Swahili**, **English** and other Native African languages. The data was sourced from Twitter and is annotated in three categories: offensive, hate, and neither.

(b) XTREMESPEECH: Listening to affected communities to define extreme speech:

XTREMESPEECH⁴⁰ Maronikolakis et al. (2022) is a hate speech dataset in 6 languages including **Swahili**. It contains 20,297 passages collected from social media and annotated in three categories: derogatory, exclusionary, and dangerous.

4.7 Named Entity Recognition**(a) MasakhaNER:**

MasakhaNER⁴¹ Adelani et al. (2021) is the first, publicly available Named Entity Recognition dataset in 10 African Languages sourced from local news. The dataset contains 921 and 3,006 sentences in **Luo** and **Swahili** respectively. The sentences are annotated in four categories: personal name, location, organization and date & time.

³⁸<https://github.com/masakhane-io/afriqa>

³⁹<https://github.com/google-research-datasets/tydiqa>

⁴⁰<https://github.com/antmarakis/xtremespeech>

⁴¹<https://git.io/masakhane-ner>

(b) Mining Wikidata for Name Resources for African Languages:

This project⁴² Sälevä and Lignos (2021) contains a list of approximately 1.9million names in 28 African languages mined from Wikipedia. It contains 23,0614 **Swahili** names, 26,215 **Oromo** names and 28,215 **Somali** names.

(c) Cross-lingual Name Tagging and Linking for 282 Languages:

The WikiAnn project⁴³ developed cross-lingual name tagging and linking for 282 languages in Wikipedia including the Kenyan languages **Swahili**, **Kikuyu**, **Somali** and **Oromo**. The data is annotated for three entity types: PER, ORG and GPE/LOC. It contains 9.3K, 6.5K, 1.0K and 709 names Swahili, Somali, Kikuyu and Oromo respectively Pan et al. (2017).

4.8 Dictionaries**(a) Glosbe Dictionary:**

Glosbe⁴⁴ is the most extensive online community-built dictionary that provides free dictionaries with in-context translations. The dictionary supports 6,000 languages, including the Kenyan languages **Swahili**, **Kikuyu**, **Luo**, **Gusii**, **Kalenjin**, **Meru** among others. Aside from sentence translation, the dictionary contains phrase illustrations, audio recordings and pronunciations, translated sentences, and automatic translators for long sentences. It has over 2 billion translations, 400K audio recordings, and 1 billion sentence examples.

(b) Mandla African language dictionary:

Mandla⁴⁵ is a free, multilingual dictionary that translates in over 100 African languages. The crowd-sourced dictionary, which has both text and audio options, gives both the native and Latin script definitions of words and has over 75K users. Some Kenyan languages Mandla supports are **Swahili**, **Kipsigis**, **Oromo**, **Makonde**, **Turkana**, **Kamba**, **Somali**, **Oromo** and **Dholuo**.

(c) A Lexicon of Key Words in Kiswahili:

This is a dictionary⁴⁶ with 52 **Swahili** words and phrases relating to technology translated to English Nyabola (2022). The dictionary was created to spark conversation on digital rights in languages other than English and is labelled according to Part of Speech.

(d) IPA-dict: Monolingual wordlists with pronunciation information in IPA:

IPA-dict⁴⁷ is the first standardized series of dictionaries of wordlists with accompanying phonemic pronunciation information in IPA - International Phonetic Alphabet transcription Doherty (2019). The data is in 23 languages, including **Swahili**.

(e) Common Swahili Slangs:

This dataset⁴⁸ contains 188 Swahili slangs and their respective proper words Masasi (2020).

5 Applications

This section provides an overview of various NLP tasks and the publicly available models for each. We also examine the current state of Kenyan languages about these tasks and techniques, highlighting the challenges and progress in developing NLP resources for these languages.

⁴²<https://github.com/bltllab/africanlp2021-wikidata-names>

⁴³<https://elisa-ie.github.io/wikiann/>

⁴⁴<https://glosbe.com/>

⁴⁵<https://dictionary.sebmita.com/>

⁴⁶<https://pacscenter.stanford.edu/publication/a-lexicon-of-key-words-in-kiswahili/>

⁴⁷<https://github.com/open-dict-data/ipa-dict>

⁴⁸<https://data.mendeley.com/datasets/b8tc96xf3h/>

5.1 Machine Translation

Rule-Based approach: This is one of the earliest methods in the field based on a deep understanding of the linguistic properties of both source and target languages. It combines expert-crafted grammar rules and dictionaries, focusing on specific linguistic aspects such as morphology, syntax, and lexical semantics. All Kenyan languages are still considered low-resource; thus, they do not have curated grammar rules for rule-based machine translation. On the other hand, all languages listed in Table 1 have available dictionaries online⁴⁹. Collectively, even with these resources, are few to no publicly available rule based MT models for Kenyan languages.

Statistical Machine Translation: This approach employs statistical techniques such as probability models to facilitate translation between source and target languages. This approach is based on the analysis of large corpora of bilingual text data, where the system learns how words, phrases and sentences in one language correspond to those in another language. It assigns probability scores to words or phrases in each target sentence, with those scoring highest considered the best translations. De Pauw et al.⁵⁰, developed a Swahili-English Statistical Machine translation system

Neural Machine Translation: This approach is currently regarded as the state of art in the field as it has shown improved performance compared to the other techniques. It employs deep learning techniques to infer high-level semantics from language translations. A prominent method in this approach is the transformer-based model with encoder-decoder architecture introduced by Vaswani et al.(2017) Many African languages including Kenyan languages have not benefited from these developments as parallel data is scarce for these languages. However, there have been several efforts to build NMT models or include them in multilingual MT models. We will break this down according to the various languages:

1. Swahili: This stands as the most resourced language in Kenya, prominently featured in both national and international contexts. Several models cover Swahili, namely: bilingual model, multilingual model, fine-tuned models.

Bilingual Models :

- Rogendo’s English-Swahili Model: Available at Hugging Face, this model represents a key resource for English-Swahili translation⁵¹.
- HPLT’s Swahili-English Model: Trained using data from OPUS and HPLT, this model facilitates Swahili to English translation⁵².
- Helsinki’s English-Swahili and Finnish-Swahili Models: These models are part of Helsinki-NLP’s offerings, demonstrating a broader linguistic reach^{53,54}.
- Masakhane English-Swahili models: These models are found on GitHub and form part of the machine translation for Africa project.

Multilingual Models :

- **NLLB, MMTAfrica, m2m-100:** These models incorporate Swahili as part of their multilingual capabilities, highlighting the language’s significant presence in global language models.
2. Kikuyu: Integrated into the NLLB model, and subsequent fine-tuning efforts. Additionally, a Swahili-Kikuyu translation model is hosted on GitHub⁵⁵, along with a Masakhane’s English-Kikuyu model.
 3. Kamba: Similarly to Kikuyu, included in both NLLB and Masakhane’s models available on GitHub.
 4. Luo: Featured in multilingual models like NLLB and specific bilingual models by Helsinki-NLP⁵⁶.
 5. Luhya: Despite available parallel datasets, there are no publicly available machine translation models for Luhya.
 6. Somali, Oromo, Maasai, Nandi (Kalenjin): Included in the MADLAD-400 model, with Oromo supported by Helsinki’s bilingual models and Somali featured in NLLB.

⁴⁹ can be found on this site: <https://glosbe.com/en>

⁵⁰<https://aclanthology.org/www.mt-archive.info/MTMRL-2011-DePauw.pdf>

⁵¹<https://huggingface.co/Rogendo/en-sw>

⁵²https://huggingface.co/HPLT/translate-sw-en-v1.0-hplt_opus

⁵³<https://huggingface.co/Helsinki-NLP/opus-mt-en-sw>

⁵⁴<https://huggingface.co/Helsinki-NLP/opus-mt-fi-sw>

⁵⁵https://github.com/starnleymbote/Kikuyu_Kiswahili-translation

⁵⁶<https://huggingface.co/Helsinki-NLP/opus-mt-en-luo>

7. Gusii, Meru, Giriama/Mijikenda, Turkana, Borana, Orma: No models are publicly available for these languages.

Closed MT Systems There are several closed machine translation systems that feature Kenya languages. Namely: Google Translate features Swahili and Oromo, machinetranslate.org features Swahili, Kikuyu, Kamba, Somali and Oromo, rephrasely features Kikuyu.

5.2 Information Retrieval

Information retrieval is still an underexplored task in the Kenyan context. We see the growing research initiatives such as Masakhane⁵⁷, and the advancement of deep learning methodologies such as transfer learning help in building new African datasets by adapting existing multilingual pretrained models in high-resource languages Ogueji et al. (2021); Ogundepo et al. (2023) that promote improving natural language processing and information retrieval for African languages Wanjawa et al. (2023a); Abedissa et al. (2023). However, these previously mentioned works only cover Swahili, leaving opportunity for exploration for the rest of the Kenya languages.

5.3 Text Classification

5.3.1 Sentiment analysis

In this subsection, we review of the state of sentiment analysis in Kenya.

Sentiment analysis for Kenyan languages has been explored limitedly in multilingual settings. For instance, Muhammad et al. (2023a) collected Afrisenti geographically distributed over the African continent and included Oromo and Kiswahili, spoken in Kenya. Other sentiment analysis applications in Africa documented include Sentimentr Cannon et al. (2022), VADERBotchway et al. (2020) and NVIVO II Ochieng (2017).

Although limited below are some sentiment analysis applications developed for the Kenyan context.

- The mapping of sentiment of the Government of Kenya's state service, *Huduma Kenya* Ng'ang'ira (2018). Ng'angira develops a prototype system to obtain, preprocess and make inferences about public sentiment on government services. For a use case, Ng'angira tests the system on 8001 tweets; 3307, 2658 and 2036 positive and neutral respectively.
- Sauer et al have analysed the Kenyan public's reaction to China's Belt and Road Initiative (BRI) Morrison et al. (2022). This is a project by China for the building of large scale infrastructure in 150 countries in Africa, Middle and Far East and Europe. In Kenya, the flagship project of the BRI is the Standard Gauge Railway (SGR). The Authors use VADER to analyse a set of multilingual tweets in English and Kiswahili constituting a commentary corpus on the SGR. For Kiswahili, the authors first translate the tweets into English then carry the sentiment analysis on the English translation. For error analysis, 200 tweets were randomly sampled for each language and sentiment class. While there's a worry that there could be an error amplification from translation, the author's demonstrate that the Swahili and English error rates do not demonstrate significant divergences.
- Moge Noor analyses sentiment reaction to 1122 tweets on the Demonetisation of legal tender using multinomial naive Bayes classification Noor (2020).
- Evanega et al on the other hand compare the attitudes in Kenya and other African countries (Uganda, Nigeria) towards the introduction of Genetically Modified Organisms (GMO) in the country Evanega et al. (2022). From a baseline of around 67% positive or neutral views towards GMO, there has been a slight increase towards positive attitudes towards GMO. They show that this level is relatively stable from January 2018 to December 2020.
- Noor (tion) developed a model to analyse sentiment of demonetization in Kenya. The data was retrieved from twitter using a web scraper with the help of advanced search. The data was collected between June and October, 1128 tweets collected. Some techniques for data processing were used to reduce noise and dimensionality in the data. The SA was performed using the Multinomial Naive Bayes Algorithm and the development of lexicon analysis which contains a word list that is negative and positive. For model evaluation, the author used a confusion matrix to summarize the performance of the prediction and its results. To validate the results, precision, recall and accuracy was used. The overall score maintained was an accuracy of 0.704.

⁵⁷<https://www.masakhane.io/>

- Jytte (2015) used sentiment analysis to help with identification, detection and tracking of terrorists activities in a more timely manner. Data was collected from twitter where the analysis was done. To accomplish the analysis ,analyst must have real time access to incident reports and make information scanning a daily process.
- Ngoge and Orero (2015) carried out mapping of sentiments, which is the process of linking classified tweets and geocoordinates.

5.4 Automatic Speech Recognition(ASR)

There is a significant opportunity to create and expand datasets and applications in local languages and dialects to improve the accuracy of ASR. Increased collaboration between local universities, research institutions, and international organizations can accelerate advancements in ASR technology. ASR technology can be applied in public service sectors such as education, healthcare, and government to improve accessibility and efficiency. Encouraging local startups to innovate in the ASR space can lead to the development of tailored solutions that address specific regional needs Reitmaier et al. (2022). However, as it stand, there are limited ASR applications for Kenyan languages.

5.5 Named Entity Recognition

We detail NER contributions in Kenya in this subsection:

The Swahili NER dataset is a Named Entity Recognition (NER) dataset generated from <https://huggingface.co/datasets/swahili> using back-translation techniques. This data has been cleaned using a couple of techniques and is ready for training a Spacy NER model without any modifications, with this data one is able to train a swahili-spacy-ner contributing to the NER space in Kenya.

Another application of NER is the Analysis of a machine learning-based algorithm used in name entity recognition. The work analysed various machine learning algorithms and implemented KNN which has been widely used in machine learning and remains one of the most popular methods to classify data. It was established by the researchers that no published study has presented Named entity recognition for the Kikuyu language using a machine learning algorithm. This research will fill this gap by recognising entities in the Kikuyu language. An evaluation was done by testing precision, recall, and F-measure. The experiment results demonstrate that using K-NN is effective in classification performance compared to other models such as Naive Bayes,Maximum Entropy Model,Support Vector Machine (SVM),Decision Tree Learning and Hybrid.With enough training data, researchers could perform an experiment and check the learning curve with accuracy that compares to state of art NER.Kamau et al. (2023)

5.6 Question Answering

The application of question-answering methodologies to indigenous Kenyan languages significantly lags behind rich-resource languages such as English. Recently, advancement have relied on large-scale language modeling, for which one example is available in Kiswahili, UlizaLlama, by JacarandaHealth, which uses a version of the ALPACA dataset translated into Kiswahili. version of the ALPACA dataset. Another initiative aimed at advancing question-answering for African languages is AfriQA

The need for Question Answering datasets in low resource languages also motivated the development of KenSwQuAD.The dataset was annotated from raw story texts of Swahili low resource language, which is a predominantly spoken in Eastern African and in other parts of the world. Question Answering (QA) datasets are important for machine comprehension of natural language for tasks such as internet search and dialogue systems. The research involved annotators in creating QA pairs from Swahili texts gathered by the Kencorpus project, which is a corpus of Kenyan languages. Out of a total of 2,585 texts, 1,445 were annotated, each with at least five QA pairs, resulting in a final dataset of 7,526 QA pairs. A quality assurance check on 12.5 percent of the annotated texts verified the correctness of all QA pairs. A proof of concept demonstrated that this dataset is suitable for QA tasks. Additionally, KenSwQuAD has enhanced the resources available for the Swahili language. The dataset used for generating KenSwQuAD was the Swahili portion of the data collected by the Kencorpus projectWanjawa et al. (2022). Kencorpus project collected primary data, both text and voice, in three low-resource languages of Swahili, Dholuo and Luhya.

5.7 Natural Language Inference (NLI) or Textual Entailment:

Natural Language Inference (NLI), or textual entailment, is a fundamental NLP task determining whether a given hypothesis logically follows a provided premise. This task typically classifies the relationship between the premise and hypothesis into one of three categories: entailment, contradiction, or neutral. While NLI models are well-developed

for high-resource languages like English, their application to low-resource languages, such as those spoken in Kenya, remains underexplored. Most NLI models that support Kenyan languages, particularly Swahili, are multilingual. These models often rely on translations of English NLI datasets into Swahili. One example is the Multilingual NLI dataset⁵⁸, which supports Swahili and has been used in models such as `mDeBERTa-v3-base-xnli-multilingual-nli-2mil7`. This model provides coverage for Swahili, leveraging cross-lingual capabilities to address the lack of Swahili-specific NLI datasets. Additionally, the `afriXNLI` dataset includes a test set for various African languages, though Swahili is the only Kenyan language represented in this dataset. Other Kenyan languages, such as Kikuyu, Luo, and Luhya, currently lack dedicated NLI datasets or models, which presents a significant gap in the research landscape.

5.8 Pretrained Language Models (PTMs)

In the AI revolution, language modeling is fundamental to NLP. At its core, language modeling involves prediction of the next word in a sequence given a probability distribution. These models are used in various applications from autocomplete to more complex tasks like classification and machine translation. The growth of language models has been catalysed by the availability of large datasets of textual data and development of deep learning architectures such as transformers.

Despite these advancements in language modelling, African languages have lagged behind due to various factors. Namely: data scarcity, technological infrastructure and commercial interest. Below we discuss the coverage of Kenyan languages in pretrain language models and how they are being utilised.

- BERT-models, GPT-models, BLOOM-models, Claude, LLaMA-models.
BERT -models which was developed by Google, has been widely used in language modelling . It currently supports 100 languages however only one Kenyan language, Swahili, is supported. Research (2024) The GPT models by OpenAI doesn't support Kenya languages yet however it understands around 60 indigenous languages from Kenya with limited proficiency. BLOOM-models developed by The BigScience project supports 46 natural languages and of those Swahili is the only Kenyan language supported with a percentage contribution of 0.02. BigScience (2023) Claude developed by Anthropic is known to include support for widely spoken African languages such as Swahili, which is a major language in Kenya however there is no mention of other Kenyan languages. LLaMA-models , support several Kenyan languages, focusing primarily on Swahili. Swahili is one of the most widely spoken languages in Kenya. LLaMA models aim to extend support to other underutilized Kenyan dialects, although the specifics about these dialects are not detailed.
- Africa AI o-specific LMs (AfriBERTa, Afro-XLMR, AfroLM, SERENGETI, mBERT)
AfriBERTa Ogueji et al. (2021) model was trained on 11 languages which includes support for Swahili and Somali, two major languages spoken in Kenya.
AfroXLMR Alabi et al. (2022) model was trained on 17 languages which includes representation of Swahili and Somali which is spoken in Kenya.
AfroLM Dossou et al. (2022) model is trained on 23 African languages including the Kenyan languages Luo and Swahili.
The SERENGETI Adebara and Abdul-Mageed (2022); Adebara et al. (2022) is a multilingual model covering 517 African languages including Luo, Bukusu, Embu, Chidigo, Ekegusii, Oromo and Borana-Arsi-Gurji, Gikuyu, Kalenjin, Kikamba, Kuria, Maasai, Giriyama, Pokomo (Kipfokomo), Rendile, Samburu, Somali, Swahili, Teso, Turkana and Tharaka. Datasets in the model are classified as Nilo-Saharan, Niger-Congo or Afro-Asiatic
- LMs finetuned on Kenyan specific languages (UlizaLLaMA, SwahBERT). These models have been finetuned with Swahili data.

Application Areas of PTMs

UlizaLLaMA, developed by Jacaranda Health⁵⁹, is finetuned to converse fluently in Swahili. This model helps to improve healthcare.

The AfroXLMR model has been used by Muhammad et al. (2023b) in sentiment analysis for 14 African languages including Swahili and Oromo. The model is able to perform both monolingual and multilingual classification of text in three classes: positive, negative or neutral.

The study by Alabi et al. (2022) used AfriBERTa in multilingual adaptive finetuning (MAFT) for 17 African languages including Oromo, Swahili and Somali. In the Kenyan context, this model can be used in cross-lingual transfer learning.

⁵⁸<https://huggingface.co/datasets/MoritzLaurer/multilingual-NLI-26lang-2mil7>

⁵⁹<https://jacarandahealth.org/jacaranda-launches-first-in-kind-swahili-large-language-model/>

mBERT which is a BERT model has multilingual capabilities to handle tasks such as Name Entity Recognition for African languages as seen in Adelani et al. (2021). The model covers 10 African languages including Swahili and Luo.

6 Governance, Policies, and Regulations

Governance, policy, ethics, and responsible NLP are interconnected concepts that relate to the oversight, regulation, and moral considerations of applying Natural Language Processing technologies. Siau and Wang (2018) studied the general concepts of AI governance, policies, and regulations. In their work, they argue that AI and other related technologies are expeditiously advancing. Hence, there is a need to discuss governance, policy, and regulatory issues in the field. The paper calls on relevant researchers, policymakers and government officials to take these matters into consideration.

6.1 Assessing AI governance

Progress on AI governance, NLP included, has been slow not only in Kenya but also in the wider continent. The ALT Advisory 2022 report Advisory (2022) on AI governance, provides a useful set of indicators upon which to assess continental and country-level progress on AI governance: existence of AI legislation, existence of Data protection integrating rights on automated decision making, existence of a national AI strategy, existence of a draft policy on AI and presence of a task force / commission to oversee AI adoption and AI integration in the national development plan.

6.2 State of AI governance in Kenya

6.2.1 Legislation and strategy

Kenya, like many countries in Africa, neither has a strategy nor regulation governing the AI sector Hankins et al. (2023). That said, the past year saw significant legislative activity. On September 8, 2023 The Ministry of Information, Communications, and the Digital Economy unveiled the Working Group on Policy and Legislative Reforms for the Information, Communications and the Digital Economy Sector. The mandate of the sector working group includes reviewing current policy, legal and institutional structures, identifying emerging technologies requiring oversight and making legislative and policy recommendations. As AI is encapsulated under emerging technologies, recommendations from this task force are of keen interest to stakeholders. Regarding AI regulation, the Robotics Society Of Kenya (RSK) introduced a draft bill seeking to regulate the Robotics and AI sector. The draft titled the "Kenya Robotics and Artificial Intelligence Society Bill, 2023" seeks to provide for the establishment of the Kenya Robotics and Artificial Intelligence Society whose envisioned mandate includes regulation, licensing, policy guidance, awareness creation and capacity building of Kenya (2023). Overall, the draft bill, currently in the petition stage, has attracted vehement opposition from practitioners on the basis of being punitive, bureaucratic and ill-timed. While it remains to be seen whether the draft bill will proceed beyond the petition, its mark on the governance discourse, albeit controversial, is cemented.

6.2.2 Data protection integrating rights on automated decision making

In 2019, the Data Protection Act No. 24 of 2019 came into effect. Although not explicit on AI governance, it is of great relevance to stakeholders within the AI ecosystem, given that data is the cornerstone of AI technology development. Key tenets of the act related to data protection include:

- incorporating operational and technical systems that include data protection principles, enforceability mechanisms, risk management, cyber-security measures, access security, physical security, and de-identification measures in all AI software and applications like chatbots,
- implementing security safeguards to ensure that personal data is accessed only by authorised persons. This includes technical safeguards for encryption, personnel vetting, and procedural safeguards like restricted access control and continuous database backup,
- limiting the amount of personal data collected for a given purpose. Thus, avoiding the processing of personal data altogether when possible and ensuring the data collected is minimised,
- enforcing transparency and lawfulness. That is, use clear and plain language to communicate with data and related subjects,

6.2.3 Presence of task force, commission to oversee AI adoption

Increasing interest in the opportunities and challenges brought about by AI technologies has elicited significant interest from stakeholders within the AI ecosystem, leading to the formation of task forces with sector-specific mandates. To illustrate, the Kenya Blockchain and AI Task Force, formed in 2018, was mandated to study emerging technologies, their use cases for development goals, and opportunities for legislation Ministry of Information and Technology (2023). The task force outlined important use cases for AI technologies within Kenya’s development framework, and these suggestions are officially integrated into Kenya’s economic development blueprints through the following:

- **Data Science Africa (DSA)** - This is a non-profit organisation founded in 2019 by Kenyan researchers, practitioners, and entrepreneurs passionate about using AI to solve real-world problems. Its main mission is to work with institutions of higher learning and other government authorities to promote AI and data science for academia and R&D, and support the adoption of data science and AI in businesses and government. The researchers are also advocating for ethical use of AI. DSA works with the government to influence AI policies. DSA ([n. d.]).
- **Information & Communications Technology Authority (ICTA)** - It was established in 2013. It is advocating for e-Government services, development and use of ICT, efficient and effective provision of ICT services and applications, rights of ICT users, and cyber-security. The authority is working to streamline the management of AI and ICT functions within the national and county governments ICTA ([n. d.]).
- **Office of the Data Protection Commissioner (ODPC)** - This is government body that is entrusted safeguard personal data. It was established under the Data Protection Act of parliament in 2019. Its mandates are available in ODPC ([n. d.]) and on this basis advocates for safe and ethical use of private data in AI research and innovation.
- **The Communications Authority of Kenya (CA)** serves as the regulatory body that oversees the country’s communications sector. It was founded in 1999 under the Kenya Information and Communications Act of 1998. Its mandate covers the advancement of various segments within the fields of information and communication, such as broadcasting, cyber security, multimedia, telecommunications, e-commerce and post and courier services C.A. ([n. d.]). In mid-2024, that authority operationalized the regulatory sandbox to test, monitor, and govern emerging technologies, innovations, and applications for AI-driven services, including chatboxes and BERTs Communication Authority of Kenya ([n. d.]).

Both ICTA and CA are domiciled in Ministry of Information Communication and Technology.

7 Indigenous Knowledge Systems and NLP

Indigenous or local knowledge systems (IKS) refer to a body of knowledge, skills, innovations, value and belief systems passed down through generations in a given cultural locality and acquired through the accumulation of experiences, relationships with the surrounding environment, and traditional community rituals, practices and institutions Marshall (2019). The indigenous knowledge systems are not mere collections of facts but are dynamic, living libraries of understanding, embodying the wisdom of centuries. IKS encompass a spectrum of skills, innovations, and deeply held beliefs rooted in the rhythm of local ecosystems, nurtured by communal rituals and enduring practices. Natural Language Processing (NLP) stands as a bridge between this ancient wisdom and the digital age. It is a tool that has the potential to decode and digitise these verbal and non-verbal languages, transforming them into accessible data for the global stage.

8 Discussion

Natural Language Processing (NLP) has garnered growing attention over the past decade, highlighting the undeniable need for continued research to enhance these techniques. This section provides insight and discusses the progress made in NLP in Kenya so far.

Section 4 presents that the NLP ecosystem in Kenya is rapidly expanding. These can be understood in the context of the African continent’s wider digitisation and technologisation. Progress in North, South, West, and East Africa has provided a context and complement to NLP efforts in Kenya. Internally, the development of datasets such as Kencorpus Wanjawa et al. (2023b), KenSwQuADWanjawa et al. (2023a) and Mozilla Common Voice⁶⁰ have helped create the necessary input infrastructure for applications such as the case in machine translation, information retrieval,

⁶⁰<https://commonvoice.mozilla.org/en>

sentiment analysis and in general language modelling tasks. As a consequence, various application models such as Serengeti Adebara et al. (2022), Afribert Ogueji (2022), and AfroXLMR Azime et al. (2023) have been developed for the Kenya NLP application contexts.

This has been accentuated significantly by international work seeking to broaden the scope of open-source or commercial products to the African continent. International efforts in the creation of datasets in SwahiliChimoto and Bassett (2022), OromoPan et al. (2017), SomaliOgunremi et al. ([n.d.]) and other languages have triggered the digital enrichment of various Kenyan languages. Several languages in Kenya, including most recently Dholuo, have been incorporated into Google Translate Goyal et al. (2022).

However, this relatively rapid expansion of the NLP ecosystem in Kenya is not without qualification. From an inspection of the covered datasets, it can be seen that the language with the most resources is Swahili, followed by Oromo, Somali, and Luo. Most of these languages are transboundary, hence their attractiveness for digitisation. This is also reproduced in the wider continent, where Bambara (spoken in Mali, Ivory Coast and other West African Countries) Diallo et al. (2021), Arabic in the North Fourati et al. (2020), and Nguni languages (IsiZulu) in the South are widely captured in NLP datasetsMarivate (2020). There is a danger of "tail of tail" differentiation where the least resourced African languages in need of the most urgent digital preservation remain unpreserved, leading to a potential loss of culture in the continent.

Concomitant to this, there is a possibility that these under-resourced languages will be considered as an international "terra nullius" - open for the preservation (and possible exploitation) by any global institutions or companies with the resources to carry out the digitisation and profit from it financially and socially.

One should look to the entrenchment of the gains made in the last decade, taken together with correcting the gaps identified, whether infrastructural, environmental or incidental. Recognising that the current wave of AI activity has been formed on the basis of an NLP attractor, it is essential to find ways of institutionalising the NLP ecosystem in Kenya to ensure the independence and robustness of Kenya's technological development.

This calls for carefully considering governance and regulation of the country's AI and data ecosystem. In this developmental phase of technology, there is a need to develop flexible regulations whose aim is not to bureaucratise, stifle or gatekeep the industry but to guide it into the possibility of full flourishing. This has been discussed in the previous section, where the data protection act and the proposed AI and Robotics Bill have been seen as tentative steps taken by the Government of Kenya in this direction. There is an opportunity here for the Government to involve stakeholders more organically in ensuring robust legal infrastructure for the growth and takeoff of NLP and AI in the country.

8.1 Research Opportunities: A road map to the Future

8.1.1 IKS and NLP in Climate Action

The effectiveness of climate action depends on the extent to which available knowledge systems have been considered in the climate change discourse and their inclusivity in addressing the negative impacts of climate change. A fertile ground for revolutionising climate action strategies lies at the crossroads of Indigenous Knowledge Systems (IKS) and Natural Language Processing (NLP). In Kenya, Indigenous communities, stewards of their lands, offer a wealth of ecological wisdom shaped by centuries of sustainable living and intimate interactions with nature. Their practices, from agroforestry to water conservation, are a testament to a profound understanding of biodiversity and ecosystem management—key to combating climate change. However, a significant hurdle remains the need for datasets capturing these rich IKS worldviews, which hinders the potential of NLP applications. This data gap limits NLP technologies' ability to fully comprehend and utilise Indigenous languages and knowledge, excluding these vital perspectives from global climate action dialogues. In regions where indigenous knowledge has been integrated into technology, for instance, in the Amazon Stejskal (2022), NLP has been used to document and translate local knowledge on medicinal plants, directly influencing sustainable practices and conservation efforts. Similarly, in Australia, Aboriginal fire management techniques, once overlooked, are now being studied and modelled through NLP Varma et al. (2024), providing insights into natural disaster management that are both ancient and innovative. By bridging these gaps, particularly in Kenya, through dedicated research into creating comprehensive datasets of IKS, NLP can unlock a treasure trove of environmental strategies. Fusing traditional wisdom and modern technology can enrich our approaches to climate challenges and ensure a more inclusive and effective global response.

9 Conclusions

This paper underscores the importance of a concerted, long-term strategy for NLP in Kenya. This strategy includes expanding data resources, developing context-aware models, and engaging with Kenya's many languages' cultural

and linguistic dimensions. Our findings reveal significant gaps in data resources, language models, and the adaptation of existing NLP technologies to the unique linguistic contexts of Kenyan languages. Finally, we offer targeted recommendations for addressing these gaps, aiming to foster the development of more inclusive, accessible, and contextually relevant NLP solutions in Kenya. To address these issues, we propose the following recommendations: a) a concerted effort should be made to expand the digital footprint of Kenyan languages through collaborative data collection and the creation of language-specific datasets. b) investments in local research and infrastructure are essential to empower Kenyan NLP Researchers and ensure that innovations are locally driven and sustainable. c) Integrating Indigenous knowledge systems into NLP workflows can provide critical linguistic and cultural context, enhancing model performance and ensuring the development of more accurate, context-aware, and culturally sensitive language technologies.

10 Limitation

Our work is derived from a literature survey on NLP research and publications about Kenyan languages. While we have covered many key publications, this review is not exhaustive. Due to the scope of our review, other relevant works and advancements in the field may not have been included. Further studies should aim to capture a more comprehensive view of the evolving landscape of NLP in Kenya.

References

- Julien Abadji, Pedro Ortiz Suarez, Laurent Romary, and Benoît Sagot. 2022. Towards a cleaner document-oriented multilingual crawled corpus. *arXiv preprint arXiv:2201.06642* (2022).
- Tilahun Abedissa, Ricardo Usbeck, and Yaregal Assabie. 2023. Amqa: Amharic question answering dataset. *arXiv preprint arXiv:2303.03290* (2023).
- Ifeoluwanimi Adebara. 2024. *Towards Afrocentric natural language processing*. Ph.D. Dissertation. University of British Columbia.
- Ife Adebara and Muhammad Abdul-Mageed. 2022. Towards afrocentric NLP for African languages: Where we are and where we can go. *arXiv preprint arXiv:2203.08351* (2022).
- Ife Adebara, AbdelRahim Elmadany, Muhammad Abdul-Mageed, and Alcides Alcoba Inciarte. 2022. Serengeti: Massively multilingual language models for africa. *arXiv preprint arXiv:2212.10785* (2022).
- David Ifeoluwa Adelani, Jade Abbott, Graham Neubig, Daniel D’souza, Julia Kreutzer, Constantine Lignos, Chester Palen-Michel, Happy Buzaaba, Shruti Rijhwani, Sebastian Ruder, et al. 2021. MasakhaNER: Named entity recognition for African languages. *Transactions of the Association for Computational Linguistics* 9 (2021), 1116–1131.
- David Ifeoluwa Adelani, Jesujoba Oluwadara Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruitter, Dietrich Klakow, Peter Nabende, Ernie Chang, et al. 2022. A few thousand translations go a long way! leveraging pre-trained models for african news translation. *arXiv preprint arXiv:2205.02022* (2022).
- David Ifeoluwa Adelani, Hannah Liu, Xiaoyu Shen, Nikita Vassilyev, Jesujoba O Alabi, Yanke Mao, Haonan Gao, and Annie En-Shiun Lee. 2023a. SIB-200: A Simple, Inclusive, and Big Evaluation Dataset for Topic Classification in 200+ Languages and Dialects. *arXiv preprint arXiv:2309.07445* (2023).
- David Ifeoluwa Adelani, Marek Masiak, Israel Abebe Azime, Jesujoba Alabi, Atnafu Lambebo Tonja, Christine Mwase, Odunayo Ogundepo, Bonaventure FP Dossou, Akintunde Oladipo, Doreen Nixdorf, et al. 2023b. Masakhanews: News topic classification for african languages. *arXiv preprint arXiv:2304.09972* (2023).
- Mofetoluwa Adeyemi, Akintunde Oladipo, Xinyu Zhang, David Alfonso-Hermelo, Mehdi Rezagholizadeh, Boxing Chen, Abdul-Hakeem Omotayo, Idris Abdulmumin, Naome A Etori, Toyib Babatunde Musa, et al. 2024. CIRAL: A Test Collection for CLIR Evaluations in African Languages. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 293–302.
- ALT Advisory. 2022. *AI Governance in Africa*. Retrieved Nov 11, 2023 from <https://ai.altadvisory.africa/wp-content/uploads/AI-Governance-in-Africa-2022.pdf>
- Jesujoba O Alabi, David Ifeoluwa Adelani, Marius Mosbach, and Dietrich Klakow. 2022. Adapting pre-trained language models to African languages via multilingual adaptive fine-tuning. *arXiv preprint arXiv:2204.06487* (2022).
- Antonios Anastasopoulos, Alessandro Cattelan, Zi-Yi Dou, Marcello Federico, Christian Federman, Dmitriy Genzel, Francisco Guzmán, Junjie Hu, Macduff Hughes, Philipp Koehn, et al. 2020. TICO-19: the translation initiative for COvid-19. *arXiv preprint arXiv:2007.01788* (2020).

- Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber. 2019. Common voice: A massively-multilingual speech corpus. *arXiv preprint arXiv:1912.06670* (2019).
- Ebbie Awino, Lilian Wanzare, Lawrence Muchemi, Barack Wanjawa, Edward Ombui, Florence Indede, Owen McOnyango, and Benard Okal. 2022. Phonemic Representation and Transcription for Speech to Text Applications for Under-resourced Indigenous African Languages: The Case of Kiswahili. *arXiv preprint arXiv:2210.16537* (2022).
- Israel Abebe Azime, Sana Sabah Al-Azzawi, Atnafu Lambebo Tonja, Iyanuoluwa Shode, Jesujoba Alabi, Ayodele Awokoya, Mardiyyah Oduwole, Tosin Adewumi, Samuel Fanijo, Oyinkansola Awosan, et al. 2023. Masakhane-Afrisenti at SemEval-2023 Task 12: Sentiment Analysis using Afro-centric Language Models and Adapters for Low-resource African Languages. *arXiv preprint arXiv:2304.06459* (2023).
- Paul Azunre, Salomey Osei, Salomey Addo, Lawrence Asamoah Adu-Gyamfi, Stephen Moore, Bernard Adabankah, Bernard Opoku, Clara Asare-Nyarko, Samuel Nyarko, Cynthia Amoaba, et al. 2021. Nlp for ghanaian languages. *arXiv preprint arXiv:2103.15475* (2021).
- Claire Babirye, Joyce Nakatumba-Nabende, Andrew Katumba, Ronald Ogwang, Jeremy Tusubira Francis, Jonathan Mukiibi, Medadi Ssentanda, Lilian D Wanzare, and Davis David. 2022. Building text and speech datasets for low resourced languages: A case of languages in east africa. (2022).
- Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, et al. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. Association for Computational Linguistics (ACL).
- Emily Bender. 2019. The# benderrule: On naming the languages we study and why it matters. *The Gradient* 14 (2019), 34.
- BigScience. 2023. BLOOM: Training Data. <https://huggingface.co/bigscience/bloom#training-data>. Accessed: 2024-07-26.
- Aric Bills, Thomas Conners, Miriam Corris, Anne David, Eyal Dubinski, Jonathan G. Fiscus, Ketty Gann, Mary Harper, Michael Kazi, Nicolas Malyska, Jennifer Melot, Jessica Ray, Anton Rytting, and Bushra Zawaydeh. 2022. IARPA Babel Dholuo Language Pack IARPA-babel403b-v1.0b. <https://doi.org/11272.1/AB2/HSAU9N>
- Alan W Black. 2019. CMU Wilderness Multilingual Speech Dataset. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 5971–5975. <https://doi.org/10.1109/ICASSP.2019.8683536>
- Raphael Kwaku Botchway, Abdul Bashiru Jibril, Zuzana Komínková Oplatková, and Miloslava Chovancová. 2020. Deductions from a Sub-Saharan African Bank’s Tweets: A sentiment analysis approach. *Cogent Economics & Finance* 8, 1 (2020), 1776006.
- C.A. [n. d.]. Communication Authority of Kenya. <https://www.ca.go.ke/who-we-are>
- Brendon J Cannon, Mikiyasu Nakayama, and Dominic R Pkalya. 2022. Understanding African views of China: analyses of student attitudes and elite media reportage in Kenya. *Journal of Eastern African Studies* 16, 1 (2022), 92–114.
- Kamau Chege, Peter Wagacha, Guy De Pauw, Lawrence Muchemi, and W Ng’ang’a. 2010. Developing an Open source Spell-checker for Gikuyu. *AfLaT 2010* (2010), 31.
- Emmanuel Kigen Chesire and Andrew Kipkebut. 2024. Current State, Challenges and Opportunities for Natural Language Processing Research and Development in Africa: A Systematic Review. In *5th Workshop on African Natural Language Processing*. <https://openreview.net/forum?id=9CsL0PvDDV>
- Everlyn Asiko Chimoto and Bruce A Bassett. 2022. Very low resource sentence alignment: Luhya and Swahili. *arXiv preprint arXiv:2211.00046* (2022).
- Chris Chinenye Emezue, Sanchit Gandhi, Lewis Tunstall, Abubakar Abid, Joshua Meyer, Quentin Lhoest, Pete Allen, Patrick Von Platen, Douwe Kiela, Yacine Jernite, et al. 2023. AfroDigits: A Community-Driven Spoken Digit Dataset for African Languages. *arXiv e-prints* (2023), arXiv–2303.
- Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. Tydi qa: A benchmark for information-seeking question answering in ty pologically di verse languages. *Transactions of the Association for Computational Linguistics* 8 (2020), 454–470.
- title =Regulatory Sandbox Framework Communication Authority of Kenya. [n. d.]. <https://www.ca.go.ke/regulatory-sandbox>

- Alexis Conneau, Ankur Bapna, Yu Zhang, Min Ma, Patrick von Platen, Anton Lozhkov, Colin Cherry, Ye Jia, Clara Rivera, Mihir Kale, et al. 2022. Xtreme-s: Evaluating cross-lingual speech representations. *arXiv preprint arXiv:2203.10752* (2022).
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116* (2019).
- Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672* (2022).
- Guy De Pauw, Naomi Maajabu, and Peter Waiganjo Wagacha. 2010. A knowledge-light approach to Luo machine translation and part-of-speech tagging. In *Proceedings of the Second Workshop on African Language Technology (AfLaT 2010)*. Valletta, Malta: European Language Resources Association (ELRA). 15–20.
- Guy De Pauw, Peter Waiganjo Wagacha, and Gilles-Maurice de Schryver. 2009. The SAWA corpus: a parallel corpus English-Swahili. In *Proceedings of the First Workshop on Language Technologies for African Languages*. 9–16.
- Mountaga Diallo, Chayma Fourati, and Hatem Haddad. 2021. Bambara language dataset for sentiment analysis. *arXiv preprint arXiv:2108.02524* (2021).
- L Doherty. 2019. ipa-dict-Monolingual wordlists with pronunciation information in IPA. (2019).
- Bonaventure FP Dossou, Atnafu Lambebo Tonja, Oreen Yousuf, Salomey Osei, Abigail Oppong, Iyanuoluwa Shode, Oluwabusayo Olufunke Awoyomi, and Chris Chinenye Emezue. 2022. AfroLM: A self-active learning-based multilingual pretrained language model for 23 African languages. *arXiv preprint arXiv:2211.03263* (2022).
- DSA. [n. d.]. Data Science Africa. <https://www.datascienceafrica.org/>
- Ahmed El-Kishky, Vishrav Chaudhary, Francisco Guzmán, and Philipp Koehn. 2020. CCAligned: A Massive Collection of Cross-lingual Web-Document Pairs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP 2020)*. Association for Computational Linguistics, Online, 5960–5969. <https://doi.org/10.18653/v1/2020.emnlp-main.480>
- Chris C Emezue and Bonaventure FP Dossou. 2022. MMTAfrica: Multilingual machine translation for African languages. *arXiv preprint arXiv:2204.04306* (2022).
- Co-authors Miquel Espla-Gomis, Juan Antonio Pérez-Ortiz, Victor M Sánchez-Cartagena, Felipe Sánchez-Martinez, and Reviewers Alexandra Birch. 2019. Global Under-Resourced Media Translation (GoURMET). *H2020 Research and Innovation ActionNumber: 825299-D1. 1–Survey of relevant low-resource languages* (2019).
- Naome Etori, Maurice Dawson, and Maria Gini. 2024. Double-edged sword: Navigating AI Opportunities and the Risk of Digital Colonization in Africa. (2024).
- Naome Etori and Maria Gini. 2024. RideKE: Leveraging Low-resource Twitter User-generated Content for Sentiment and Emotion Detection on Code-switched RHS Dataset.. In *Proceedings of the 14th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis*. 234–249.
- Sarah Evanega, Joan Conrow, Jordan Adams, and Mark Lynas. 2022. The state of the ‘GMO’ debate-toward an increasingly favorable and less polarized media conversation on ag-biotech? *GM Crops & Food* 13, 1 (2022), 38–49.
- Chayma Fourati, Abir Messaoudi, and Hatem Haddad. 2020. Tunizi: a tunisian arabizi sentiment analysis dataset. *arXiv preprint arXiv:2004.14303* (2020).
- Raymond G Gordon Jr. 2005. Ethnologue, languages of the world. <http://www.ethnologue.com/> (2005).
- Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics* 10 (2022), 522–538.
- Emma Hankins, Pablo Fuentes Nettel, Livia Martinescu, Grau Gonzalo, and Rahim Sulmaan. 2023. *The Government AI Readiness Index 2023*. Retrieved Jul 12, 2023 from <https://oxfordinsights.com/wp-content/uploads/2023/12/2023-Government-AI-Readiness-Index-1.pdf>
- Michael A Hedderich, Lukas Lange, Heike Adel, Jannik Strötgen, and Dietrich Klakow. 2020. A survey on recent approaches for natural language processing in low-resource scenarios. *arXiv preprint arXiv:2010.12309* (2020).
- Arvi Hurskainen. 2004. Helsinki corpus of Swahili. *Compilers: Institute for Asian and African Studies (University of Helsinki) and CSC* (2004).
- ICTA. [n. d.]. Informatio & Communicaiton Technology Authority. <https://icta.go.ke/>

- Ayyoob ImaniGooghari, Peiqin Lin, Amir Hossein Kargaran, Silvia Severini, Masoud Jalili Sabet, Nora Kassner, Chunlan Ma, Helmut Schmid, André FT Martins, François Yvon, et al. 2023. Glot500: Scaling Multilingual Corpora and Language Models to 500 Languages. *arXiv preprint arXiv:2305.12182* (2023).
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the NLP world. *arXiv preprint arXiv:2004.09095* (2020).
- [Author’s First Name] Jytte. 2015. Tweeting the Jihad: Social Media Networks of Western Foreign Fighters in Syria and Iraq. [*Journal Name*] [Volume Number], [Issue Number] (2015), [Page Range].
- Francis M. Kamau, Kennedy O. Ogada, and Cheruiyot W. Kipruto. 2023. Analysis of Machine-Based Learning Algorithm Used in Named Entity Recognition. *Informing Science: The International Journal of an Emerging Transdiscipline* 26 (2023), 69–84. <https://doi.org/10.28945/5073>
- David Kamholz, Jonathan Pool, and Susan M Colowick. 2014. PanLex: Building a Resource for Panlingual Lexical Translation.. In *LREC*. 3145–3150.
- B Kituku, G Musumba, and P Wagacha. 2015. Kamba Part of Speech Tagger Using Memory-Based Approach. *International Journal on Natural Language Computing* 4, 2 (2015), 43–53.
- Colin Leong, Joshua Nemecek, Jacob Mansdorfer, Anna Filighera, Abraham Owodunni, and Daniel Whitenack. 2022. Bloom library: Multimodal datasets in 300+ languages for a variety of downstream tasks. *arXiv preprint arXiv:2210.14712* (2022).
- Vukosi Marivate. 2020. Why African natural language processing now? A view from South Africa# AfricaNLP. *Leap 4.0: African Perspectives on the Fourth Industrial Revolution* (2020), 126.
- Antonis Maronikolakis, Axel Wisioerek, Leah Nann, Haris Jabbar, Sahana Udupa, and Hinrich Schütze. 2022. Listening to affected communities to define extreme speech: Dataset and experiments. *arXiv preprint arXiv:2203.11764* (2022).
- Larry Marshall. 2019. AI+ Indigenous knowledge a powerful tool posing critical questions.
- Noel Masasi. 2020. Common swahili slangs. <https://data.mendeley.com/datasets/b8tc96xf3h/1>
- Josh Meyer, David Ifeoluwa Adelani, Edresson Casanova, Alp Öktem, Daniel Whitenack Julian Weber, Salomon Kabongo, Elizabeth Salesky, Iroro Orife, Colin Leong, Perez Ogayo, et al. 2022. BibleTTS: a large, high-fidelity, multilingual, and uniquely African speech corpus. *arXiv preprint arXiv:2207.03546* (2022).
- Communications Ministry of Information and Technology. 2023. *Emerging Digital Technologies for Kenya*. Retrieved Jul 12, 2023 from <http://repository.ca.go.ke/handle/123456789/1044>
- Michelle Morrison, Jeffery Sauer, Henry Overos, Juan Newlands, Kathleen Stewart, Tess Wood, and ARLIS. 2022. Analyzing multilingual discussions of the Standard Gauge Railway project in Kenya using natural language processing and social media analysis. (2022).
- Shamsuddeen Hassan Muhammad, Idris Abdulmumin, Abinew Ali Ayele, Nedjma Ousidhoum, David Ifeoluwa Adelani, Seid Muhie Yimam, Ibrahim Sa’id Ahmad, Meriem Beloucif, Saif M Mohammad, Sebastian Ruder, et al. 2023a. Afrisenti: A twitter sentiment analysis benchmark for african languages. *arXiv preprint arXiv:2302.08956* (2023).
- Shamsuddeen Hassan Muhammad, Idris Abdulmumin, Seid Muhie Yimam, David Ifeoluwa Adelani, Ibrahim Sa’id Ahmad, Nedjma Ousidhoum, Abinew Ayele, Saif M Mohammad, Meriem Beloucif, and Sebastian Ruder. 2023b. SemEval-2023 task 12: sentiment analysis for african languages (AfriSenti-SemEval). *arXiv preprint arXiv:2304.06845* (2023).
- Joaquim Mussandi and Andreas Wichert. 2024. NLP Tools for African Languages: Overview. In *Proceedings of the 16th International Conference on Computational Processing of Portuguese-Vol. 2*. 73–82.
- Lucas A. Ngoge and Joseph Onderi Orero. 2015. Mapping of Terrorist Activities in Kenya using Sentiment Analysis. Faculty of Information Technology, Strathmore University. <mailto:lucas.ngoge@strathmore.edu>
- Velma N Ngoni. 2022. *English–Bukusu Automatic Machine Translation for Digital Services Inclusion in E-governance*. Ph.D. Dissertation. university of nairobi.
- Khanh Nguyen and Hal Daumé III. 2019. Global Voices: Crossing borders in automatic news summarization. *arXiv preprint arXiv:1910.00421* (2019).
- Judy Nyakairu Ng’ang’ira. 2018. *A Prototype for mapping of tweets on state services for decision support: a case of Huduma Kenya*. Ph.D. Dissertation. Strathmore University.
- Ibrahim Moge Noor. 2020. Sentiment analysis on new currency in kenya using Twitter dataset. In *Proceeding International Conference on Science and Engineering*, Vol. 3. 237–240.

- Ibrahim Moge Noor. Year of Publication. Sentiment Analysis on New Currency in Kenya Using Twitter Dataset. *Where it was published, if applicable* Volume number, if applicable, Issue number, if applicable (Year of Publication), Page range, if applicable. Department of engineer, Faculty of Computer engineer.
- Nanjala Nyabola. 2022. A Lexicon of Key Words in Kiswahili. <https://pacscenter.stanford.edu/publication/a-lexicon-o>
- Robi Koki Ochieng. 2017. An exploration of gender based violence in online print media stories on prominent women journalists in Kenya. *Extracted from 'Media coverage on Online Violence against Women Journalists in Kenya'*AMWIK (2017).
- ODPC. [n. d.]. Office of the Data Protection Commissioner. <https://www.odpc.go.ke/functions-of-the-office/>
- Robotics Society of Kenya. 2023. *The Kenya Robotics and Artificial Intelligence Society Bill, 2023*. Retrieved Jul 12, 2024 from https://www.dataguidance.com/sites/default/files/the_kenya_robotics_and_artificial_intelligen
- Perez Ogayo, Graham Neubig, and Alan W Black. 2022. Building African Voices. *arXiv preprint arXiv:2207.00688* (2022).
- Kelechi Ogueji. 2022. *AfriBERTa: Towards viable multilingual language models for low-resource languages*. Master's thesis. University of Waterloo.
- Kelechi Ogueji, Yuxin Zhu, and Jimmy Lin. 2021. Small data? no problem! exploring the viability of pretrained multilingual language models for low-resourced languages. In *Proceedings of the 1st Workshop on Multilingual Representation Learning*. 116–126.
- Sewade Ogun, Abraham T Owodunni, Tobi Olatunji, Eniola Alese, Babatunde Oladimeji, Tejumade Afonja, Kayode Olaleye, Naome A Etori, and Tosin Adewumi. 2024. 1000 African Voices: Advancing inclusive multi-speaker multi-accent speech synthesis. *arXiv preprint arXiv:2406.11727* (2024).
- Ogunayo Ogundepo, Tajuddeen R Gwadabe, Clara E Rivera, Jonathan H Clark, Sebastian Ruder, David Ifeoluwa Adelan, Bonaventure FP Dossou, Abdou Aziz DIOP, Clayton Sikasote, Gilles Hacheme, et al. 2023. AfriQA: Cross-lingual Open-Retrieval Question Answering for African Languages. *arXiv preprint arXiv:2305.06897* (2023).
- Ogunayo Ogundepo, Xinyu Zhang, Shuo Sun, Kevin Duh, and Jimmy Lin. 2022. AfriCLIRMatrix: Enabling cross-lingual information retrieval for african languages. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 8721–8728.
- Toyib Ogunremi, Serah sessi Akojenu, Anthony Soronnadi, Olubayo Adekanmbi, and David Ifeoluwa Adelan. [n. d.]. AfriHG: News Headline Generation for African Languages. In *5th Workshop on African Natural Language Processing*.
- Tobi Olatunji, Tejumade Afonja, Aditya Yadavalli, Chris Chinenye Emezue, Sahib Singh, Bonaventure FP Dossou, Joanne Osuchukwu, Salomey Osei, Atnafu Lambebo Tonja, Naome Etori, et al. 2023. Afrispeech-200: Pan-african accented speech dataset for clinical and general domain asr. *Transactions of the Association for Computational Linguistics* 11 (2023), 1669–1685.
- Edward Ombui, Lawrence Muchemi, and Peter Wagacha. 2019. Hate speech detection in code-switched text messages. In *2019 3rd International Symposium on Multidisciplinary Studies and Innovative Technologies (ISMSIT)*. IEEE, 1–6.
- Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. Cross-lingual name tagging and linking for 282 languages. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 1946–1958.
- Paul Raine. 2018. Building sentences with web 2.0 and the tatoeba database. *Accents Asia* 10, 2 (2018), 2–7.
- Thomas Reitmaier, Electra Wallington, Dani Kalarikalayil Raju, Ondrej Klejch, Jennifer Pearson, Matt Jones, Peter Bell, and Simon Robinson. 2022. Opportunities and Challenges of Automatic Speech Recognition Systems for Low-Resource Language Speakers. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (, New Orleans, LA, USA,) (*CHI '22*). Association for Computing Machinery, New York, NY, USA, Article 299, 17 pages. <https://doi.org/10.1145/3491102.3517639>
- Google Research. 2024. BERT Multilingual. <https://github.com/google-research/bert/blob/master/multilingual.md> Accessed: 2024-07-21.
- Kiptoo Rono. 2021. A Kiswahili dataset. <https://data.mendeley.com/datasets/rbn6nmygcn/1>
- Sebastian Ruder and H Korashy. 2019. The 4 biggest open problems in NLP. *Ain Shams Eng. J* (2019).
- Jonne Sälevä and Constantine Lignos. 2021. Mining Wikidata for Name Resources for African Languages. *arXiv preprint arXiv:2104.00558* (2021).

- Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2019. Wikimatrix: Mining 135m parallel sentences in 1620 language pairs from wikipedia. *arXiv preprint arXiv:1907.05791* (2019).
- Shivachi Casper Shikali and Mokhosi Refuoe. 2019. Language modeling data for swahili. (2019).
- Keng Siau and Weiyu Wang. 2018. Artificial Intelligence: A Study on Governance, Policies, and Regulations.
- Kathleen Siminyu, Godson Kalipe, Davor Orlic, Jade Abbott, Vukosi Marivate, Sackey Freshia, Prateek Sibal, Bhanu Neupane, David I Adelani, Amelia Taylor, et al. 2021. AI4D–African Language Program. *arXiv preprint arXiv:2104.02516* (2021).
- Kenneth Steimel, Sandra Kübler, and Daniel Dakota. 2023. Towards a Swahili Universal Dependency Treebank: Leveraging the Annotations of the Helsinki Corpus of Swahili. In *Proceedings of the Fourth workshop on Resources for African Indigenous Languages (RAIL 2023)*. 86–96.
- Clara Stejskal. 2022. Sources of innovation in the Brazilian Amazon rainforest. (2022).
- Espen Stranger-Johannessen and Bonny Norton. 2017. The African storybook and language teacher identity in digital times. *The Modern Language Journal* 101, S1 (2017), 45–60.
- Jennifer Tracey, Stephanie Strassel, Ann Bies, Zhiyi Song, Michael Arrigo, Kira Griffitt, Dana Delgado, Dave Graff, Seth Kulick, Justin Mott, et al. 2019. Corpus building for low resource languages in the DARPA LORELEI program. In *Proceedings of the 2nd Workshop on Technologies for MT of Low Resource Languages*. 48–55.
- Smita Varma, Soumendu Shekar Roy, and Praveen Kumar Rai. 2024. Machine Learning for Forest Fire Risk and Resilience. In *Geospatial Technology to Support Communities and Policy: Pathways to Resiliency*. Springer, 171–184.
- Jiayi Wang, David Ifeoluwa Adelani, Sweta Agrawal, Marek Masiak, Ricardo Rei, Eleftheria Briakou, Marine Carpuat, Xuanli He, Sofia Bourhim, Andiswa Bukula, et al. 2023. AfriMTE and AfriCOMET: Enhancing COMET to Embrace Under-resourced African Languages. (2023).
- B. Wanjawa, L. Wanzare, F. Indede, O. McOnyango, L. Muchemi, and E. Ombui. 2022. Kencorpus - Kenyan languages corpus. <https://kencorpus.co.ke/>. Accessed May 05, 2022.
- Barack Wanjawa, Lilian Wanzare, Florence Indede, Owen McOnyango, Edward Ombui, and Lawrence Muchemi. 2023b. Kencorpus: A kenyan language corpus of swahili, dholuo and luhya for natural language processing tasks. *Journal for Language Technology and Computational Linguistics* 36 (2023), 1–27.
- Barack W Wanjawa, Lilian DA Wanzare, Florence Indede, Owen McOnyango, Lawrence Muchemi, and Edward Ombui. 2023a. KenSwQuAD—A Question Answering Dataset for Swahili Low-resource Language. *ACM Transactions on Asian and Low-Resource Language Information Processing* 22, 4 (2023), 1–20.
- Nan Yan and Cheng Xu. [n. d.]. Decolonizing African NLP: A Survey on Power Dynamics and Data Colonialism in Tech Development. In *5th Workshop on African Natural Language Processing*.