

# MasakhaNER: Named Entity Recognition for African Languages

David Ifeoluwa Adelani<sup>1\*</sup>, Jade Abbott<sup>2\*</sup>, Graham Neubig<sup>3</sup>, Daniel D’souza<sup>4\*</sup>,  
Julia Kreutzer<sup>5\*</sup>, Constantine Lignos<sup>6\*</sup>, Chester Palen-Michel<sup>6\*</sup>, Happy Buzaaba<sup>7\*</sup>  
Shrutri Rijhwani<sup>3</sup>, Sebastian Ruder<sup>8</sup>, Stephen Mayhew<sup>9</sup>, Israel Abebe Azime<sup>10\*</sup>  
Shamsuddeen H. Muhammad<sup>11,12\*</sup>, Chris Chinenye Emezue<sup>13\*</sup>, Joyce Nakatumba-Nabende<sup>14\*</sup>  
Perez Ogayo<sup>15\*</sup>, Anuluwapo Aremu<sup>16\*</sup>, Catherine Gitau\*, Derguene Mbaye\*  
Jesujoba Alabi<sup>17\*</sup>, Seid Muhie Yimam<sup>18</sup>, Tajuddeen Gwadabe<sup>19\*</sup>, Ignatius Ezeani<sup>20\*</sup>  
Rubungo Andre Niyongabo<sup>21\*</sup>, Jonathan Mukiibi<sup>14</sup>, Verrah Otiende<sup>22\*</sup>  
Iroro Orife<sup>23\*</sup>, Davis David\*, Samba Ngom\*, Tosin Adewumi<sup>24\*</sup>  
Paul Rayson<sup>20</sup>, Mofetoluwa Adeyemi\*, Gerald Muriuki<sup>14</sup>, Emmanuel Anebi\*  
Chiamaka Chukwuneke<sup>20</sup>, Nkiruka Odu<sup>25</sup>, Eric Peter Wairagala<sup>14</sup>  
Samuel Oyerinde\*, Clemencia Siro\*, Tobius Saul Bateesa<sup>14</sup>, Temilola Oloyede\*  
Yvonne Wambui\*, Victor Akinode\*, Deborah Nabagereka<sup>14</sup>, Maurice Katusiime<sup>14</sup>  
Ayodele Awokoya<sup>26\*</sup>, Mouhamadane MBOUP\*, Dibora Gebreyohannes\*, Henok Tilaye\*  
Kelechi Nwaike\*, Degaga Wolde\*, Abdoulaye Faye\*, Blessing Sibanda<sup>27\*</sup>  
Orevaoghene Ahia<sup>28\*</sup>, Bonaventure F. P. Dossou<sup>29\*</sup>, Kelechi Ogueji<sup>30\*</sup>  
Thierno Ibrahima DIOP\*, Abdoulaye Diallo\*, Adewale Akinfaderin\*  
Tendai Marengereke\*, and Salomey Osei<sup>10\*</sup>

- \* Masakhane, <sup>1</sup> Spoken Language Systems Group (LSV), Saarland University, Germany  
<sup>2</sup> Retro Rabbit, <sup>3</sup> Language Technologies Institute, Carnegie Mellon University  
<sup>4</sup> ProQuest, <sup>5</sup> Google Research, <sup>6</sup> Brandeis University, <sup>8</sup> DeepMind, <sup>9</sup> Duolingo  
<sup>7</sup> Graduate School of Systems and Information Engineering, University of Tsukuba, Japan.  
<sup>10</sup> African Institute for Mathematical Sciences (AIMS-AMMI), <sup>11</sup> University of Porto  
<sup>12</sup> Bayero University, Kano, <sup>13</sup> Technical University of Munich, Germany  
<sup>14</sup> Makerere University, Kampala, Uganda, <sup>15</sup> African Leadership University, Rwanda  
<sup>16</sup> University of Lagos, Nigeria, <sup>17</sup> Max Planck Institute for Informatics, Germany.  
<sup>18</sup> LT Group, Universität Hamburg, <sup>19</sup> University of Chinese Academy of Science, China  
<sup>20</sup> Lancaster University, <sup>21</sup> University of Electronic Science and Technology of China, China.  
<sup>22</sup> United States International University - Africa (USIU-A), Kenya. <sup>23</sup> Niger-Volta LTI  
<sup>24</sup> Luleå University of Technology, <sup>25</sup> African University of Science and Technology, Abuja  
<sup>26</sup> University of Ibadan, Nigeria, <sup>27</sup> Namibia University of Science and Technology  
<sup>28</sup> Instadeep, <sup>29</sup> Jacobs University Bremen, Germany, <sup>30</sup> University of Waterloo.

## Abstract

We take a step towards addressing the under-representation of the African continent in NLP research by creating the first large publicly available high-quality dataset for named entity recognition (NER) in ten African languages, bringing together a variety of stakeholders. We detail characteristics of the languages to help researchers understand the challenges that these languages pose for NER. We analyze our datasets and conduct an extensive empirical evaluation of state-of-the-art methods across both supervised and transfer learning settings. We release the data, code, and models in order to inspire future research on African NLP<sup>1</sup>.

<sup>1</sup><https://github.com/masakhane-io/masakhane-ner>

## 1 Introduction

Africa has over 2,000 languages (Eberhard et al., 2020); however, these languages are scarcely represented in existing natural language processing (NLP) datasets, research, and tools (Martinus and Abbott, 2019). √ et al. (2020) investigate the reasons for these disparities by examining how NLP for low-resource languages is constrained by several societal factors. One of these factors is the geographical and language diversity of NLP researchers. For example, only 5 out of the 2695 affiliations of authors whose work was published at the five major NLP conferences in 2019 are from African institutions (Caines, 2019). Many NLP tasks such as machine translation, text classification, part-of-speech tagging, and named entity

recognition would benefit from the knowledge of native speakers who are involved in the development of datasets and models.

In this paper, we focus on named entity recognition (NER)—one of the most impactful tasks in NLP (Sang and De Meulder, 2003; Lample et al., 2016). NER is a core NLP task in information extraction and an essential component of numerous products including spell-checkers, localization of voice and dialogue systems, and conversational agents. It also enables identifying African names, places and organizations for information retrieval. African languages are under-represented in this crucial task due to lack in datasets, reproducible results, and researchers who understand the challenges that such languages pose for NER.

In this paper, we take an initial step towards improving representation for African languages in the NER task. More specifically, this paper makes the following contributions:

- (i) We bring together language speakers, dataset curators, NLP practitioners, and evaluation experts to address the constraints facing African NER. Based on the availability of online news corpora and language annotators, we develop NER datasets, models, and evaluation covering ten widely spoken African languages.
- (ii) We curate NER datasets from local sources to ensure relevance of future research for native speakers of the respective languages.
- (iii) We train and evaluate multiple NER models for all ten languages. Our experiments provide insights into the transfer across languages, and highlight open challenges.
- (iv) We release the data, code, and model outputs in order to facilitate future research on the specific challenges of African NER.

## 2 Related Work

**African NER datasets** NER is a well-studied sequence labeling task (Yadav and Bethard, 2018) and has been subject of many shared tasks in different languages (Tjong Kim Sang, 2002; Tjong Kim Sang and De Meulder, 2003; *ijc*, 2008; Shaalan, 2014; Benikova et al., 2014). Most of the available datasets are from high-resource languages. Although there have been efforts to create NER datasets for lower-resourced languages,

such as the WikiAnn corpus (Pan et al., 2017) covering 282 languages, such datasets consist of “silver-standard” labels created by transferring annotations from English to other languages through cross-lingual links in knowledge bases. As the data comes from Wikipedia, it includes some African languages, most with fewer than 10k tokens.

Other NER datasets for African languages are SADiLaR (Eiselen, 2016) for ten South African languages based on government data, and small corpora of less than 2K sentences for Yorùbá (Alabi et al., 2020) and Hausa (Hedderich et al., 2020). Additionally, LORELEI language packs (Strassel and Tracey, 2016) were created for some African languages including Yorùbá, Hausa, Amharic, Somali, Twi, Swahili, Wolof, Kinyarwanda, and Zulu, but are not publicly available.

**NER models** Popular sequence labeling models for NER are CRF (Lafferty et al., 2001), CNN-BiLSTM (Chiu and Nichols, 2016), BiLSTM-CRF (Huang et al., 2015), and CNN-BiLSTM-CRF (Ma and Hovy, 2016). The traditional CRF makes use of hand-crafted features like part-of-speech tags, context words and word capitalization. Neural network NER models are initialized with word embeddings like Word2Vec (Mikolov et al., 2013), GloVe (Pennington et al., 2014) and FastText (Bojanowski et al., 2016). More recently, pre-trained language models such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and LUKE (Yamada et al., 2020) produce state-of-the-art results for the task. Multilingual variants of these models like mBERT and XLM-RoBERTa (Conneau et al., 2020) make it possible to train NER models for several languages using transfer learning. Language-specific parameters and adaptation to unlabelled data of the target language have yielded further gains (Pfeiffer et al., 2020a,b).

## 3 Focus Languages

Table 1 gives an overview of the languages considered in this work, in terms of their language family, number of speakers and the regions in Africa where they are spoken. We chose to focus on these languages due to the availability of online news corpora, annotators, and most importantly because they are widely spoken African languages. Both region and language family might indicate a notion of proximity for NER, either because of linguistic features shared within that family, or because data sources cover a common set of locally rele-

Language	Family	Speakers	Region
Amharic	Afro-Asiatic-Ethio-Semitic	26M	East
Hausa	Afro-Asiatic-Chadic	63M	West
Igbo	Niger-Congo-Volta-Niger	27M	West
Kinyarwanda	Niger-Congo-Bantu	12M	East
Luganda	Niger-Congo-Bantu	7M	East
Luo	Nilo Saharan	4M	East
Nigerian-Pidgin	English Creole	75M	West
Swahili	Niger-Congo-Bantu	98M	Central & East
Wolof	Niger-Congo-Senegambia	5M	West & NW
Yorùbá	Niger-Congo-Volta-Niger	42M	West

Table 1: Language, family, number of speakers (Eberhard et al., 2020), and regions in Africa.

vant entities. We highlight language specifics for each language to illustrate the diversity of this selection of languages in Section 3.1, and then showcase the differences in named entities across these languages in Section 3.2.

### 3.1 Language Characteristics

**Amharic** (amh) uses the Fidel script consisting of 33 basic scripts (**ሀ** (hä) **ለ** (lä) **መ** (mä) **ሠ** (sä) ...), each of them with at least 7 vowel sequences (such as **ሀ** (hä) **ሁ** (hu) **ሂ** (hi) **ሃ** (ha) **ሄ** (hē) **ህ** (hi) **ሖ** (ho)). This results in more than 231 characters or Fidels. Numbers and punctuation marks are also represented uniquely with specific Fidels (**፩** (1), **፪** (2), ... and **፬** (.), **!** (!), **፣** (;),).

**Hausa** (hau) has 23-25 consonants, depending on the dialect and five short and five long vowels. Hausa has labialized phonemic consonants, as in /gw/ e.g. ‘agwagwa.’ As found in some African languages, implosive consonants also exist in Hausa, e.g. ‘b, ‘d, etc as in ‘barna’. Similarly, the Hausa approximant ‘r’ is realized in two distinct manners: roll and trill, as in ‘rai’ and ‘ra’ayi’, respectively.

**Igbo** (ibo) is an agglutinative language, with many frequent suffixes and prefixes (Emenanjo, 1978). A single stem can yield many word-forms by addition of affixes that extend its original meaning (Onyenwe and Hepple, 2016). Igbo is also tonal, with two distinctive tones (high and low) and a down-stepped high tone in some cases. The alphabet consists of 28 consonants and 8 vowels (A, E, I, I, O, O, U, U). In addition to the Latin letters

(except *c*), Igbo contains the following digraphs: (ch, gb, gh, gw, kp, kw, nw, ny, sh).

**Kinyarwanda** (kin) makes use of 24 Latin characters with 5 vowels similar to English and 19 consonants excluding *q* and *x*. Moreover, Kinyarwanda has 74 additional complex consonants (such as *mb*, *mpw*, and *njyw*).<sup>2</sup> It is a tonal language with three tones: low (no diacritic), high (signaled by “/”) and falling (signaled by “^”). The default word order is Subject-Verb-Object.

**Luganda** (lug) is a tonal language with subject-verb-object word order. The Luganda alphabet is composed of 24 letters that include 17 consonants (*p*, *v*, *f*, *m*, *d*, *t*, *l*, *r*, *n*, *z*, *s*, *j*, *c*, *g*), 5 vowel sounds represented in the five alphabetical symbols (*a*, *e*, *i*, *o*, *u*) and 2 semi-vowels (*w*, *y*). It also has a special consonant *ng’*.

**Luo** (luo) is a tonal language with 4 tones (high, low, falling, rising) although the tonality is not marked in orthography. It has 26 Latin consonants without Latin letters (*c*, *q*, *v*, *x* and *z*) and additional consonants (*ch*, *dh*, *mb*, *nd*, *ng’*, *ng*, *ny*, *nj*, *th*, *sh*). There are nine vowels (*a*, *e*, *i*, *o*, *u*, *ɛ*, *ɛ*, *ɔ*, *ɔ*) which are distinguished primarily by advanced tongue root (ATR) harmony (De Pauw et al., 2007).

**Nigerian-Pidgin** (pcm) is a largely oral, national lingua franca with a distinct phonology from English, its lexifier language. Portuguese, French, and especially indigenous languages form the substrate of lexical, phonological, syntactic, and semantic influence on Nigerian-Pidgin (NP). English lexical items absorbed by NP are commonly phonologically closer to indigenous Nigerian languages, notably in the realization of vowels. As a rapidly evolving language, the orthography of NP is undergoing codification and indigenization (Ofiong Mensah, 2012; Onovbiona, 2012; Ojarikre, 2013).

**Swahili** (swa) is the most widely spoken language on the African continent. It has 30 letters including 24 Latin letters without characters (*q* and *x*) and six additional consonants (*ch*, *dh*, *gh*, *ng’*, *sh*, *th*) unique to Swahili pronunciation.

<sup>2</sup><https://gazettes.africa/archive/rw/2014/rw-government-gazette-dated-2014-10-13-no-41%20bis.pdf>

**Wolof** (wo1) has an alphabet similar to that of French. It consists of 29 characters, including all letters of the French alphabet except H, V and Z. It also includes the characters ɗ (“ng”, lowercase: ŋ) and Ñ (“gn” as in Spanish). Accents are present, but limited in number (À, É, È, Ó). However, unlike many other Niger-Congo languages, Wolof is not a tonal language.

**Yorùbá** (yɔɾ) has 25 Latin letters without the Latin characters (c, q, v, x and z) and with additional letters (ẹ, gb, ẹ, ọ). Yorùbá is a tonal language with three tones: low (“˘”), middle (“ˉ”, optional) and high (“ˊ”). The tonal marks and underdots are referred to as diacritics and they are needed for the correct pronunciation of a word. Yorùbá is a highly isolating language and the sentence structure follows Subject-Verb-Object.

### 3.2 Named Entities

As most work on for NER has centered around English, it is not clear how well existing systems will generalize to languages that realize NE differently, e.g., in terms of the surrounding sentence structure or surface forms. Hu et al. (2020) analyzed generalization from English to other languages for NER but only evaluated on two African languages. Among all languages, they found that models particularly struggled to generalize to named entities in Swahili.

In order to highlight the differences across our focus languages and in comparison to English, Table 2 shows an English example sentence containing color-coded PER, LOC, and DATE entities and their corresponding translations. The following characteristics could pose challenges for NER systems developed for English:

- Amharic shares no lexical overlap with the English source sentence.
- While “Zhang” is identical across all Latin-script languages, “Kano” features accents in Wolof and Yorùbá due to its localization.
- The Fidel script has no capitalization, which could hinder transfer from other languages.
- Igbo, Wolof, and Yorùbá all use diacritics, which are not present in the English alphabet.
- The sentence structure and surface form of NE is the same in English and Nigerian-Pidgin, but there exist lexical differences (e.g. in terms of how time is realized).

- Between the 10 African languages, “Nigeria” is spelled in 5 different ways.
- Numerical “18”: Igbo, Wolof and Yorùbá write out their numbers, resulting in different numbers of tokens for the entity span.

## 4 Data and Annotation Methodology

Our data was obtained from local news sources, in order to ensure relevance of the dataset for native speakers from those regions. The data was annotated using the ELISA tool (Lin et al., 2018). Our annotators were native speakers of the languages who volunteered through the *Masakhane* community.<sup>4</sup> The annotators come from the same regions as the news sources and were trained on how to perform NER annotation. We annotated four entity types: Personal name (PER), Location (LOC), Organization (ORG), and date & time (DATE). The annotated entities were inspired by the English CoNLL-2003 Corpus (Tjong Kim Sang, 2002). We replaced the MISC tag with the DATE tag following Alabi et al. (2020) as the MISC tag may be ill-defined and cause disagreement among non-expert annotators. We report the number of annotators as well as general statistics of the datasets in Table 3.

A key objective of our annotation procedure was to create high-quality datasets by ensuring a high annotator agreement. To achieve high agreement scores, we ran collaborative workshops for each language, which allowed annotators to discuss any disagreements. ELISA provides an entity-level F1-score and also an interface for annotators to correct their mistakes, making it easy to achieve inter-annotator agreement scores between 0.96 and 1.0 for all languages.

We report inter-annotator agreement scores in Table 4 using Fleiss’ Kappa (Fleiss, 1971) at both the token and entity level. The latter considers each span an annotator proposed as an entity. As a result of our workshops, all our datasets have exceptionally high inter-annotator agreement. For Kinyarwanda, Luo, Swahili, and Wolof, we report perfect inter-annotator agreement scores ( $\kappa = 1$ ). For each of these languages, we had two annotators per token, who were instructed to discuss and resolve conflicts among themselves. The Appendix provides a detailed confusion matrix.

<sup>4</sup><https://www.masakhane.io>

Language	Sentence
English	The Emir of <b>Kano</b> turbaned <b>Zhang</b> who has spent <b>18 years</b> in <b>Nigeria</b>
Amharic	<b>የካኖ ኢምር በናይጄርያ ጌጅ ዓመት ያሳለፈውን ዛንግን ዋና መሪ አደረጉት</b>
Hausa	Sarkin <b>Kano</b> yayi wa <b>Zhang</b> wanda yayi <b>shekara 18</b> a <b>Nigeria</b> sarauta
Igbo	Onye Emir nke <b>Kano</b> kpubere <b>Zhang</b> okpu onye nke nọgoro <b>afọ iri na asato</b> na <b>Naijiria</b>
Kinyarwanda	Emir w'i <b>Kano</b> yimitse <b>Zhang</b> wari umaze <b>imyaka 18</b> muri <b>Nijeriya</b>
Luganda	Emir w'e <b>Kano</b> yatikkidde <b>Zhang</b> amaze <b>emyaka 18</b> mu <b>Nigeria</b>
Luo	Emir mar <b>Kano</b> ne orwakone turban <b>Zhang</b> ma osedak <b>Nigeria</b> kwuom <b>higni 18</b>
Nigerian-Pidgin	Emir of <b>Kano</b> turban <b>Zhang</b> wey don spend <b>18 years</b> for <b>Nigeria</b>
Swahili	Emir wa <b>Kano</b> alimvisha kilemba <b>Zhang</b> ambaye alikaa <b>miaka 18</b> nchini <b>Nigeria</b>
Wolof	Emiiru <b>Kanó</b> dafa kaala kii di <b>Zhang</b> mii def <b>Nigeria</b> <b>fukki at ak juróom ñett</b>
Yorùbá	Èmíà ilú <b>Kánò</b> wé lówàní lé orí <b>Zhang</b> èní tí ó tí lo <b>òdún méjìdínlógún</b> ní orílẹ̀-èdè <b>Nàìjíríà</b>

Table 2: Example of named entities in different languages. **PER**, **LOC**, and **DATE** are in colours purple, orange, and green respectively. The original sentence is from BBC Pidgin.<sup>3</sup>

Language	Data Source	Train/ dev/ test	# Anno.	PER	ORG	LOC	DATE	% of Entities in Tokens	# Tokens
Amharic	DW & BBC	1750/ 250/ 500	4	730	403	1,420	580	15.13	37,032
Hausa	VOA Hausa	1903/ 272/ 545	3	1,490	766	2,779	922	12.17	80,152
Igbo	BBC Igbo	2233/ 319/ 638	6	1,603	1,292	1,677	690	13.15	61,668
Kinyarwanda	IGIHE news	2110/ 301/ 604	2	1,366	1,038	2096	792	12.85	68,819
Luganda	BUKEDDE news	2003/ 200/ 401	3	1,868	838	943	574	14.81	46,615
Luo	Ramogi FM news	644/ 92/ 185	2	557	286	666	343	14.95	26,303
Nigerian-Pidgin	BBC Pidgin	2100/ 300/ 600	5	2,602	1,042	1,317	1,242	13.25	76,063
Swahili	VOA Swahili	2104/ 300/ 602	6	1,702	960	2,842	940	12.48	79,272
Wolof	Lu Defu Waxu & Saabal	1,871/ 267/ 536	2	731	245	836	206	6.02	52,872
Yorùbá	GV & VON news	2124/ 303/ 608	5	1,039	835	1,627	853	11.57	83,285

Table 3: Statistics of our datasets including their source, number of sentences in each split, number of annotators, number of entities of each label type, percentage of tokens that are named entities, and total number of tokens.

## 5 Experimental Setup

### 5.1 NER baseline models

To evaluate baseline performance on our dataset, we experiment with three popular NER models: CNN-BiLSTM-CRF, multilingual BERT (mBERT), and XLM-RoBERTa (XLM-R). The latter two models are implemented using the HuggingFace transformers toolkit (Wolf et al., 2019). For each language, we train the model on the in-language training data and evaluate on its test data.

**CNN-BiLSTM-CRF** This architecture was proposed for NER by Ma and Hovy (2016). For each input sequence, we first compute the vector representation for each word by concatenating character-level encodings from a CNN and vector embeddings for each word. Following Rijhwani et al. (2020), we use randomly initialized word embeddings since we do not have high-quality pre-

trained embeddings for all the languages in our dataset. Our model is implemented using the DyNet toolkit (Neubig et al., 2017).

**mBERT** We fine-tune multilingual BERT (Devlin et al., 2019) on our NER corpus by adding a linear classification layer to the pre-trained transformer model, and training it end-to-end. mBERT was trained on 104 languages including only two African languages: Swahili and Yorùbá. We use the mBERT-base cased model with 12-layer Transformer blocks consisting of 768-hidden size and 110M parameters.

**XLM-R** XLM-R (Conneau et al., 2020) was trained on 100 languages including Amharic, Hausa, and Swahili. The major differences between XLM-R and mBERT are (1) XLM-R was trained on Common Crawl while mBERT was trained on Wikipedia; (2) XLM-R is based on

Dataset	Token		Entity		Disagreement from Type
	Fleiss' Kappa	Fleiss' Kappa	Fleiss' Kappa	Fleiss' Kappa	
amh	0.987	0.959	0.987	0.959	0.044
hau	0.988	0.962	0.988	0.962	0.097
ibo	0.995	0.983	0.995	0.983	0.071
kin	1.0	1.0	1.0	1.0	0.0
lug	0.997	0.99	0.997	0.99	0.023
luo	1.0	1.0	1.0	1.0	0.0
pcm	0.989	0.966	0.989	0.966	0.048
swa	1.0	1.0	1.0	1.0	0.0
wol	1.0	1.0	1.0	1.0	0.0
yor	0.99	0.964	0.99	0.964	0.079

Table 4: Inter-annotator agreement for our datasets calculated using Fleiss’ kappa  $\kappa$  at the token and entity level. Disagreement from type refers to the proportion of all entity-level disagreements, which are due only to type mismatch.

RoBERTa, which is trained with a masked language model (MLM) objective while mBERT was additionally trained with a next sentence prediction objective. For all the experiments, we use the XLM-R-base cased model consisting of 12 layers, with a hidden size of 768 and 270M parameters.

**MeanE-BiLSTM** This is a simple BiLSTM model with an additional linear classifier. For each input sequence, we first extract a sentence embedding from mBERT or XLM-R LM before passing it into the BiLSTM model. Following, Reimers and Gurevych (2019), we make use of the mean of the 12-layer output embeddings of the LM (i.e. *MeanE*). This has been shown to provide better sentence representations than the embedding of the [CLS] token used for fine-tuning mBERT and XLM-R.

## 5.2 Improving the Baseline Models

In this section, we consider techniques to improve the baseline models such as utilizing gazetteers, transfer learning from another domain and language, and aggregating NER datasets by regions. For these experiments, we focus on the PER, ORG, and LOC categories, because the gazetteers from Wikipedia do not contain DATE entities and some source domains and languages that we transfer from do not have DATE annotation. We make use of the XLM-R model because it generally outperforms mBERT in our experiments (see Section 6).

### 5.2.1 Gazetteers for NER

Gazetteers are lists of named entities collected from manually crafted resources such as GeoN-

ames or Wikipedia. Before the widespread adoption of neural networks, NER methods used gazetteers-based features to improve performance (Ratinov and Roth, 2009). These features are created for each  $n$ -gram in the dataset and are typically binary-valued, indicating whether the  $n$ -gram is present in the gazetteer.

Recently, Rijhwani et al. (2020) showed that augmenting the neural CNN-BiLSTM-CRF model with gazetteer features can improve NER performance for low-resource languages. We conduct similar experiments on the languages in our dataset, using entity lists from Wikipedia as gazetteers. For Luo and Nigerian-Pidgin, which do not have their own Wikipedia, we use entity lists from English Wikipedia.

### 5.2.2 Transfer Learning

Here, we focus on cross-domain transfer from Wikipedia to the news domain, and cross-lingual transfer from English and Swahili NER datasets to the other languages in our dataset.

**Domain Adaptation from WikiAnn** We make use of the WikiAnn corpus (Pan et al., 2017), which is available for five of the languages in our dataset: Amharic, Igbo, Kinyarwanda, Swahili and Yorùbá. For each language, the corpus contains 100 sentences in each of the training, development and test splits except for Swahili, which contains 1K sentences in each split. For each language, we train on the corresponding WikiAnn training set and either zero-shot transfer to our respective test set or additionally fine-tune on our training data.

**Cross-lingual transfer** For training the cross-lingual transfer models, we use the CoNLL-2003 NER dataset in English with over 14K training sentences and our annotated corpus. We make use of the languages that are supported by the XLM-R model and are widely spoken in East and West Africa like Swahili and Hausa. The English corpus has been shown to transfer very well to low resource languages (Hedderich et al., 2020; Lauscher et al., 2020). We first train on either the English CoNLL-2003 data or our training data in Swahili, Hausa, or Nigerian-Pidgin before testing on the African languages.

### 5.3 Aggregating Languages by Regions

As previously illustrated in Table 2, several entities have the same form in different languages while some entities may be more common in the

Language	In mBERT?	In XLM-R?	Percent OOV Test Entities	CNN-BiLSTM CRF	mBERT-base MeanE / FTune	XLM-R-base MeanE / FTune
amh	✗	✓	72.94	52.89	0.0 / 0.0	63.57 / <b>70.96</b>
hau	✗	✓	33.40	83.70	81.15 / 87.34	86.10 / <b>89.44</b>
ibo	✗	✗	46.56	78.48	76.45 / <b>85.11</b>	73.77 / 84.51
kin	✗	✗	57.85	64.61	65.77 / 70.98	63.07 / <b>73.93</b>
lug	✗	✗	61.12	74.31	70.46 / <b>80.56</b>	67.75 / <b>80.71</b>
luo	✗	✗	65.18	66.42	56.76 / 72.65	52.60 / <b>75.14</b>
pcm	✗	✗	61.26	66.43	81.16 / <b>87.78</b>	82.12 / <b>87.39</b>
swa	✓	✓	40.97	79.26	83.13 / 86.37	84.46 / <b>87.55</b>
wol	✗	✗	69.73	60.43	57.16 / <b>66.10</b>	54.38 / 64.38
yor	✓	✗	65.99	67.07	74.35 / <b>78.64</b>	66.67 / 77.58
avg	–	–	57.50	69.36	64.64 / 71.55	69.45 / <b>79.16</b>
avg (excl. amh)	–	–	55.78	71.19	71.82 / 79.50	70.10 / <b>80.07</b>

Table 5: NER model comparison, showing F1-score on the test sets after 50 epochs averaged over 5 runs. This result is for all 4 tags in the dataset: PER, ORG, LOC, DATE. **Bold** marks the top score (tied if within the range of SE). mBERT and XLM-R are trained in two ways (1) MeanE: mean output embeddings of the 12 LM layers are used to initialize BiLSTM + Linear classifier, and (2) FTune: LM fine-tuned end-to-end with a linear classifier

region where the language is spoken. To study the performance of NER models across geographical areas, we combine languages based on the region of Africa that they are spoken in Table 1: (1) East region including Kinyarwanda, Luganda, Luo, and Swahili; (2) West Region with Hausa, Igbo, Nigerian-Pidgin, Wolof, and Yorùbá languages, (3) East and West regions i.e all languages except Amharic because of the writing script.

## 6 Results

### 6.1 Baseline Models

Table 5 shows the F1-score obtained by CNN-BiLSTM-CRF, mBERT and XLM-R models on the test set of the ten African language corpus when training on our in-language data. We additionally indicate whether the language is supported by the pre-trained language models (✓). The number of out of vocab (OOV) entities in the test set is also reported alongside results of the baseline models. In general, the datasets with greater numbers of OOV entities have lower performance with the CNN-BiLSTM-CRF model, while those with lower OOV rates (Hausa, Igbo, Swahili) have higher performance. Similar to the findings of Devlin et al. (2019), the CNN-BiLSTM-CRF model performs worse than mBERT and XLM-R. We expect performance to be better (e.g., for Amharic and Nigerian-Pidgin with over 18 F1 point difference) when using pre-trained word embeddings for the initialization of the BiLSTM model rather than random initialization (we leave this for future work as discussed in Section 7).

Interestingly, the pre-trained language models (PLMs) have reasonable performance even on languages they were not trained on such as Igbo, Kinyarwanda, Luganda, Luo, and Wolof. However, languages supported by the PLM tend to have better performance overall. We observe that XLM-R models have significantly better performance on five languages; two of the languages (Amharic and Swahili) are supported by the pre-trained XLM-R. Similarly, mBERT has better performance for Yorùbá since the language is part of the PLM’s training corpus. Although mBERT is trained on Swahili, XLM-R shows better performance. This observation is consistent with (Hu et al., 2020) and could be because XLM-R is trained on more Swahili texts (Common Crawl with 275M tokens (Conneau et al., 2020)) whereas mBERT is trained on smaller texts from Wikipedia (6M tokens<sup>5</sup>).

Another observation is that mBERT tends to have better performance for the non-Bantu Niger-Congo languages i.e., Igbo, Wolof, and Yorùbá. On the other hand, XLM-R works better for Afro-Asiatic languages (i.e., Amharic and Hausa), Nilo-Saharan (i.e., Luo) and Bantu languages like Kinyarwanda, Luganda, and Swahili. We also note that the writing script is one of the primary factors for the transfer of knowledge in PLMs for the languages they were not trained on. For example, mBERT achieves an F1-score of 0.0 on Amharic because of the script.

We performed further analysis on the transfer

<sup>5</sup><https://github.com/mayhewsw/multilingual-data-stats>

Method	amh	hau	ibo	kin	lug	luo	pcm	swa	wol	yor	avg
CNN-BiLSTM-CRF	50.05	85.26	81.28	62.08	75.51	67.47	63.20	77.96	63.42	66.45	69.26
+ Gazetteers	49.54	85.46	81.45	<b>65.55</b>	75.08	67.03	<b>68.14</b>	<b>80.15</b>	61.17	<b>67.37</b>	<b>70.09</b>

Table 6: Improving NER models using Gazetteers. The result is only for 3 Tags: PER, ORG & LOC. Models trained for 50 epochs. Result is an average over 5 runs.

Method	amh	hau	ibo	kin	lug	luo	pcm	swa	wol	yor	avg
XLM-R-base	69.98	90.08	86.10	73.73	80.84	76.28	86.57	<b>88.77</b>	68.82	78.40	77.36
WikiAnn zero-shot	29.57	–	20.39	8.21	–	–	–	35.36	–	6.41	–
eng-CoNLL zero-shot	–	66.39	45.48	36.11	34.82	31.91	68.38	74.47	21.57	33.37	33.87
pcm zero-shot	–	69.69	44.58	43.02	44.28	32.72	–	76.33	25.30	34.49	37.40
swa zero-shot	–	83.12*	53.80	57.58	57.27*	34.74	70.48*	–	34.40	47.96	47.63
hau zero-shot	–	–	59.99*	59.88*	57.00	43.94*	69.70	84.38*	43.11*	60.17*	54.02*
WikiAnn + finetune	<b>70.82</b>	–	85.57	73.44	–	–	–	87.66	–	76.31	–
eng-CoNLL + finetune	–	89.56	86.59	71.82	77.82	71.12	84.24	87.18	67.45	75.12	74.99
pcm + finetune	–	90.31	87.10	74.65	79.94	75.50	–	87.81	64.50	78.31	77.33
swa + finetune	–	91.56	86.18	74.23	80.21	75.22	86.94	–	70.77	80.21	77.80
hau + finetune	–	–	86.73	74.49	80.45	75.06	88.21	87.59	<b>70.98</b>	79.59	77.88
combined East Langs.	–	–	–	<b>75.65</b>	81.31	77.54	–	88.29	–	–	–
combined West Langs.	–	90.89	87.16	–	–	–	87.24	–	69.73	80.10	–
combined 9 Langs.	–	<b>91.68</b>	<b>87.78</b>	75.16	<b>81.45</b>	<b>78.36</b>	<b>88.31</b>	88.10	69.86	<b>80.30</b>	<b>78.81</b>

Table 7: Transfer Learning Result (i.e. F1-score). 3 Tags: PER, ORG & LOC. WikiAnn, eng-CoNLL, and the annotated datasets are trained for 50 epochs. Fine-tuning is only for 10 epochs. Results are averaged over 5 runs and the total average (avg) is computed over ibo, kin, lug, luo, wol, yor languages. The overall highest F1-score is in **bold**, and the best F1-score in the zero-shot settings has asterisk (\*) beside it.

Source Language	PER	ORG	LOC
eng-CoNLL	45.2	26.2	29.8
pcm	34.8	46.8	32.8
swa	52.0	51.5	38.0
hau	60.0	51.7	47.7

Table 8: Average Per-named entity F1-score for the zero-shot NER using the XLM-R model. Average computed over ibo, kin, lug, luo, wol, yor languages.

abilities of mBERT and XLM-R by extracting sentence embeddings from LM to train a BiLSTM model (*MeanE-BiLSTM*) instead of fine-tuning them end-to-end. Table 5 shows that languages that are not supported by mBERT or XLM-R generally perform worse than CNN-BiLSTM-CRF model (despite being randomly initialized) except *kin*. Also, sentence embeddings extracted from mBERT often lead to better performance than XLM-R for languages they both do not support (like *ibo*, *kin*, *lug*, *luo*, and *wol*).

## 6.2 Evaluation of Gazetteer Features

Table 6 shows the performance of the CNN-BiLSTM-CRF model with the addition of gazetteer features as described in Section 5.2.1.

On average, the model that uses gazetteer features performs better than the baseline. In general, languages with larger gazetteers, such as Swahili and Yorùbá (16K and 10K entities in their gazetteers respectively), have more improvement in performance than those with fewer gazetteer entries, such as Amharic and Luganda (2K and 500 gazetteer entities respectively). This indicates that having high-coverage gazetteers is important for the model to efficiently learn from gazetteer features.

## 6.3 Transfer Learning Experiments

Table 7 shows the result for the different transfer learning approaches, which we discuss individually in the following sections.

### 6.3.1 Cross-domain Transfer

We evaluate cross-domain transfer from Wikipedia to the news domain for the five languages that are available in the WikiAnn (Pan et al., 2017) dataset. In the zero-shot setting, the NER F1-score is low i.e., less than 40 F1-score for all languages, with Kinyarwanda and Yorùbá having less than 10 F1-score. This is likely due to the number of training sentences present in WikiAnn: there are only 100 sentences in the datasets of Amharic, Igbo, Kin-

yarwanda and Yorùbá. Although the Swahili corpus has 1,000 sentences, the 35 F1-score shows that transfer is not very effective. In general, cross-domain transfer is a challenging problem, and is even harder when the number of training examples from the source domain is small. Fine-tuning on the in-domain news NER data does not improve over the baseline (XLM-R-base).

### 6.3.2 Cross-Lingual Transfer

**Zero-shot** In the zero-shot setting, we evaluated NER models trained on English *eng-CoNLL03* dataset, Nigerian-Pidgin (*pcm*), Swahili (*swa*), and Hausa (*hau*) annotated corpus. We excluded the MISC entity in the *eng-CoNLL03* corpus because it is absent in our target dataset. Table 7 shows the result for the (zero-shot) transfer performance. We observe that the closer the source and target languages are geographically, the better the performance. The *pcm* model (trained on only 2K sentences) transfers better than the *eng-CoNLL03* model (trained on 14K sentences) with an average improvement of 3.5 F1 and individual language improvements ranging from 1–9 F1 (apart from Igbo). *swa* performs better than *pcm* with an improvement of over 10 F1 on average. For some languages like Luo, the improvement is quite small (2%), probably because they do not belong to the same family. We found that, on average, transferring from Hausa gave the best F1, with improvement of over 20%, 16%, and 6% than using *eng-CoNLL*, *pcm*, and *swa* respectively. Per-entity analysis in Table 8 shows that the largest improvements are obtained for ORG and LOC entities. In general, zero-shot transfer is most effective when transferring from Hausa and Swahili.

**Fine-tuning** We use the target language corpus to fine-tune the NER models previously trained on *eng-CoNLL*, *pcm*, and *swa*. On average, there is only a small improvement as compared to the XLM-R base model. In particular, we see significant improvement for Hausa, Nigerian-Pidgin, Wolof, and Yorùbá using either *swa* or *hau* as the source NER model.

## 6.4 Regional Influence on NER

We evaluate whether combining different language training datasets by region influences the performance on individual language test sets. Table 7 shows that all languages spoken in West Africa (*hau*, *ibo*, *wol*, *pcm*, *yor*) have slightly better per-

formance (0.6–1.7 F1) when we train on their combined training data. However, for the East-African languages, the F1 score only improved for three languages (*kin*, *lug*, *luo*). Training NER model on all nine languages leads to better performance on all languages except Swahili. On average over six languages (*ibo*, *kin*, *lug*, *luo*, *wol*, *yor*), the performance improves by 1.5 F1.

## 7 Conclusion and Future Work

We address the NER task for African languages by bringing together a variety of stakeholders to create a high-quality NER dataset for ten African languages. We evaluate multiple state-of-the-art NER models and establish strong baselines. We have released our best model that can recognize named entities for ten African languages on HuggingFace Model Hub<sup>6</sup> (see Appendix for details). We also investigate cross-domain transfer with experiments on five languages with the WikiAnn dataset, along with cross-lingual transfer for low-resource NER using the English CoNLL-2003 dataset and other languages supported by XLM-R. In the future, we plan to use pretrained word embeddings such as GloVe (Pennington et al., 2014) and fastText (Bojanowski et al., 2017) instead of random initialization for the CNN-BiLSTM-CRF, investigate more NER models such as XLM-R-large, increase the number of annotated sentences per language, and expand the dataset to more African languages.

### Acknowledgements

We would like to thank Heng Ji and Ying Lin for providing the ELISA NER tool used for annotation. We also thank the Spoken Language Systems Chair, Dietrich Klakow at Saarland University for providing GPU resources to train the models. Finally, we thank Adhi Kuncoro for useful feedback on a draft of this paper.

### References

- CLAD centre de linguistique appliquée de dakar. <http://clad.ucad.sn/>. Accessed: 2021-01-20.
- 2008. *Proceedings of the IJCNLP-08 Workshop on Named Entity Recognition for South and South East Asian Languages*.
- MustGo MustGo et al. 2015. *Igbo Language - Structure, Writing & Alphabet - MustGo*.

<sup>6</sup><https://huggingface.co/Davlan/xlm-roberta-large-masakhaner>

- Jesujoba Alabi, Kwabena Amponsah-Kaakyire, David Adelani, and Cristina España-Bonet. 2020. **Masive vs. curated embeddings for low-resourced languages: the case of Yorùbá and Twi**. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 2754–2762, Marseille, France. European Language Resources Association.
- F. O. Asahiah, O. A. Odejobi, and E. R. Adagunodo. 2017. Restoring tone-marks in standard Yorùbá electronic text: Improved model. *Computer Science*, 18(3):301–315.
- BBC. 2016. Pidgin - West African lingua franca. <https://www.bbc.com/news/world-africa-38000387>. Accessed: 2021-01-19.
- Mary E. Beckman and Janet B Pierrehumbert. 1986. *Intonational Structure in English and Japanese*. Phonology Yearbook.
- Darina Benikova, Chris Biemann, and Marc Reznicek. 2014. **NoSta-D named entity annotation for German: Guidelines and dataset**. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 2524–2531, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2016. Enriching word vectors with subword information. *arXiv preprint arXiv:1607.04606*.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Andrew Caines. 2019. **The geographic diversity of NLP conferences**.
- Jason P.C. Chiu and Eric Nichols. 2016. **Named entity recognition with bidirectional LSTM-CNNs**. *Transactions of the Association for Computational Linguistics*, 4:357–370.
- Bernard Comrie, editor. 2009. *The World's Major Languages*. Routledge, London.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. **Unsupervised cross-lingual representation learning at scale**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Guy De Pauw, Peter W Wagacha, and Dorothy Atieno Abade. 2007. **Unsupervised induction of Dholuo word classes using maximum entropy learning**, page 8.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. **BERT: Pre-training of deep bidirectional transformers for language understanding**. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, Minneapolis, Minnesota. Association for Computational Linguistics.
- David M. Eberhard, Gary F. Simons, and Charles D. Fennig (eds.). 2020. **Ethnologue: Languages of the world**, twenty-third edition.
- Roald Eiselen. 2016. **Government domain named entity recognition for South African languages**. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 3344–3348, Portorož, Slovenia. European Language Resources Association (ELRA).
- Nolue Emenanjo. 1978. *Elements of Modern Igbo Grammar - a descriptive approach*. Oxford University Press, Ibadan, Nigeria.
- Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.
- ∇, Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa Matsila, Timi Fasubaa, Taiwo Fagbohunge, Solomon Oluwole Akinola, Shamsuddeen Muhammad, Salomon Kabongo Kabenamualu, Salomey Osei, Freshia Sackey, Rubungo Andre Niyongabo, Ricky Macharm, Perez Ogayo, Orevaoghene Ahia, Musie Meressa Berhe, Mofetoluwa Adeyemi, Masabata Mokgesi-Seling, Lawrence Okegbemi, Laura Martinus, Kolawole Tajudeen, Kevin Degila, Kelechi Ogueji, Kathleen Siminyu, Julia Kreutzer, Jason Webster, Jamiil Toure Ali, Jade Abbott, Iroro Orife, Ignatius Ezeani, Idris Abdulkadir Dangana, Herman Kamper, Hady Elshahar, Goodness Duru, Ghollah Kioko, Murhabazi Espoir, Elan van Biljon, Daniel Whitenack, Christopher Onyefuluchi, Chris Chinenye Emezue, Bonaventure F. P. Dossou, Blessing Sibanda, Blessing Bassey, Ayodele Olabiyi, Arshath Ramkilowan, Alp Öktem, Adewale Akinfaderin, and Abdallah Bashir. 2020. **Participatory research for low-resourced machine translation: A case study in African languages**. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, Online.
- Elodie Gauthier, Laurent Besacier, Sylvie Voisin, Michael Melese, and Uriel Pascal Elingui. 2016. **Collecting resources in sub-Saharan African languages for automatic speech recognition: a case study of Wolof**. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 3863–3867, Portorož, Slovenia. European Language Resources Association (ELRA).
- Michael A. Hedderich, David Adelani, Dawei Zhu, Jesujoba Alabi, Udia Markus, and Dietrich Klakow.

2020. [Transfer learning and distant supervision for multilingual transformer models: A study on African languages](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2580–2591, Online. Association for Computational Linguistics.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. [XTREME: A Massively Multilingual Multi-task Benchmark for Evaluating Cross-lingual Generalization](#). In *Proceedings of ICML 2020*.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. [Bidirectional lstm-crf models for sequence tagging](#).
- Alexandre Kimenyi. 1980. *A Relational Grammar of Kinyarwanda*, volume 91 of *University of California Publications in Linguistics*. University of California Press, Berkeley.
- John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. [Conditional random fields: Probabilistic models for segmenting and labeling sequence data](#). In *Proceedings of the Eighteenth International Conference on Machine Learning, ICML '01*, pages 282–289, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. [Neural Architectures for Named Entity Recognition](#). In *Proceedings of NAACL-HLT 2016*.
- Anne Lauscher, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. 2020. [From zero to hero: On the limitations of zero-shot language transfer with multilingual Transformers](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4483–4499, Online. Association for Computational Linguistics.
- Ying Lin, Cash Costello, Boliang Zhang, Di Lu, Heng Ji, James Mayfield, and Paul McNamee. 2018. [Platforms for non-speakers annotating names in any language](#). In *Proceedings of ACL 2018, System Demonstrations*, pages 1–6, Melbourne, Australia. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#).
- Xuezhe Ma and Eduard Hovy. 2016. [End-to-end sequence labeling via bi-directional LSTM-CNNs-CRF](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1064–1074, Berlin, Germany. Association for Computational Linguistics.
- Laura Martinus and Jade Z Abbott. 2019. [A focus on neural machine translation for African languages](#). *arXiv preprint arXiv:1906.05685*.
- Marilyn Merritt and Mohamed H Abdulaziz. 1985. [Swahili as a National Language in East Africa](#).
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. [Distributed representations of words and phrases and their compositionality](#). In *Advances in Neural Information Processing Systems*, volume 26, pages 3111–3119. Curran Associates, Inc.
- Jackson Muhirwe. 2009. [Morphological analysis of tone marked Kinyarwanda text](#). In *International Workshop on Finite-State Methods and Natural Language Processing*, pages 48–55. Springer.
- Judith Nakayiza. 2013. *The sociolinguistics of multilingualism in Uganda: A case study of the official and non-official language policy, planning and management of Luruuri-lunyara and Luganda*. Ph.D. thesis, SOAS, University of London.
- Graham Neubig, Chris Dyer, Y. Goldberg, A. Matthews, Waleed Ammar, Antonios Anastasopoulos, Miguel Ballesteros, David Chiang, Daniel Clothiaux, Trevor Cohn, Kevin Duh, Manaal Faruqui, Cynthia Gan, Dan Garrette, Yangfeng Ji, Lingpeng Kong, Adhiguna Kuncoro, Manish Kumar, Chaitanya Malaviya, Paul Michel, Y. Oda, M. Richardson, Naomi Saphra, Swabha Swayamdipta, and Pengcheng Yin. 2017. [DyNet: The dynamic neural network toolkit](#). *ArXiv*, abs/1701.03980.
- Rubungo Andre Niyongabo, Qu Hong, Julia Kreutzer, and Li Huang. 2020. [KINNEWS and KIRNEWS: Benchmarking cross-lingual text classification for Kinyarwanda and Kirundi](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5507–5521, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Richard Nordquist. 2020. [West African Pidgin English \(WAPE\)](#). [thoughtco.com/west-african-pidgin-english-wape-1692496](https://thoughtco.com/west-african-pidgin-english-wape-1692496). Accessed: 2021-01-19.
- Antoine Nzeyimana. 2020. [Morphological disambiguation from stemming data](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4649–4660, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Eyo Offiong Mensah. 2012. [Grammaticalization in Nigerian Pidgin](#). *Íkala, revista de lenguaje y cultura*, 17(2):167–179.
- Chinyere Ohiri-Aniche. 2007. [Stemming the tide of centrifugal forces in Igbo orthography](#). *Dialectical Anthropology*, 31(4):423–436.
- Anthony Ojarikre. 2013. [Perspectives and problems of codifying nigerian pidgin english orthography](#). *Perspectives*, 3(12).

- Fagbolu Olutola et al. 2019. Model for translation of English language noun phrases to Luganda. *London Journal of Research in Computer Science and Technology*.
- Ijite Blessing Onovbiona. 2012. Serial verb construction in Nigerian Pidgin.
- Ikechukwu E. Onyenwe and Mark Hepple. 2016. Predicting morphologically-complex unknown words in igbo. In *Text, Speech, and Dialogue*, pages 206–214, Cham. Springer International Publishing.
- 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)*, pages 1946–1958, Vancouver, Canada. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Lothar Peter and Hans-Georg Wolf. 2007. A comparison of the varieties of West African pidgin English. *World Englishes*, 26(1):3–21.
- Jonas Pfeiffer, Ivan Vuli, Iryna Gurevych, and Sebastian Ruder. 2020a. [MAD-X: An Adapter-based Framework for Multi-task Cross-lingual Transfer](#). In *Proceedings of EMNLP 2020*.
- Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2020b. [Unks everywhere: Adapting multilingual language models to new scripts](#). *arXiv preprint arXiv:2012.15562*.
- Lev Ratinov and Dan Roth. 2009. [Design challenges and misconceptions in named entity recognition](#). In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning (CoNLL-2009)*, pages 147–155, Boulder, Colorado. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Shruti Rijhwani, Shuyan Zhou, Graham Neubig, and Jaime Carbonell. 2020. [Soft gazetteers for low-resource named entity recognition](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8118–8123, Online. Association for Computational Linguistics.
- Ibrahim T. Sabiu, Fakhru A. Zainol, and Mohammed S. Abdullahi. 2018. Hausa people of northern Nigeria and their development. *Asian People Journal (APJ)*, 1(1):179–189.
- Erik F Sang and Fien De Meulder. 2003. Introduction to the conll-2003 shared task: Language-independent named entity recognition. In *Proceedings of CoNLL 2003*.
- Mark Sebba. 1997. *Contact languages: Pidgins and creoles*. Macmillan International Higher Education.
- K. Shaalan. 2014. A survey of arabic named entity recognition and classification. *Computational Linguistics*, 40:469–510.
- Stephanie Strassel and Jennifer Tracey. 2016. [LORELEI language packs: Data, tools, and resources for technology development in low resource languages](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 3273–3280, Portorož, Slovenia. European Language Resources Association (ELRA).
- Erik F. Tjong Kim Sang. 2002. [Introduction to the CoNLL-2002 shared task: Language-independent named entity recognition](#). In *COLING-02: The 6th Conference on Natural Language Learning 2002 (CoNLL-2002)*.
- Erik F. Tjong Kim Sang and Fien De Meulder. 2003. [Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition](#). In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pages 142–147.
- R. Vanamali. 1993. *Pidgin and Creole linguistics*. University of Calabar, Calabar, Nigeria.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R’emi Louf, Morgan Funtowicz, and Jamie Brew. 2019. [Huggingface’s transformers: State-of-the-art natural language processing](#). *ArXiv*, abs/1910.03771.
- Vikas Yadav and Steven Bethard. 2018. [A survey on recent advances in named entity recognition from deep learning models](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2145–2158, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. 2020. [LUKE: Deep contextualized entity representations with entity-aware self-attention](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6442–6454, Online. Association for Computational Linguistics.
- Seid Muhie Yimam, Abinew Ali Ayele, Gopalakrishnan Venkatesh, and Chris Biemann. 2020. [Introducing various semantic models for Amharic: Experimentation and evaluation with multiple tasks and datasets](#). *arXiv:2011.01154*.

## A Language Characteristics

**Amharic** (amh) is an Ethio-Semitic and Afro-Asiatic language. It is mainly spoken in Ethiopia, Israel, and America (Eberhard et al., 2020). It has about 57 million speakers, where 32 million of them are native speakers, and uses Ge'ez or Fidel script for writing. The Fidel script consists of 33 basic scripts (ሀ (hä) ለ (lä) ጠ (mä) ሠ (šä) ...) each of them with at least seven vowel sequences (such as ... ሀ (hä) ሁ (hu) ሂ (hī) ሃ (ha) ሄ (hē) ህ (hi) ሆ (ho)) which results in more than 231 characters or Fidels. Numbers and punctuation marks are also represented uniquely with a specific Fidels (፩ (1), ፪ (2), ፫ (3), ... or # (.), !(!), ፤ (;), ...). Named entity recognition in Amharic is particularly challenging as 1) orthographic features such as initial capitals are not helpful 2) pre-processing components such as part-of-speech tagging and lemmatization are not publicly available (Yimam et al., 2020), 3) named entities are difficult to identify unless contextual information is provided, for example, the token አበበ (abebe) can be a verb or a person name. While there is some effort for Amharic NER, the only publicly available dataset is from the SAY project at New Mexico State University's Computing Research Laboratory<sup>7</sup>. However, the annotation guideline and strategies of the SAY dataset are not available with the dataset and most of the annotations are inconsistent (non-named entities are marked as named entities, a lot of named entities are missed from each document, etc.). The annotation and generation of benchmark datasets, building statistical models, and evaluation of results are the main contributions of our work for Amharic NER. While Preparing the dataset we have found difficulties that come from simple white-space tokenization like prepositions like ከ which is closer to the word "in" being a part of a word.

**Hausa** (hau): The Hausa language is widely spoken in the West and Central African Sahel belt. In West Africa, it is the most spoken indigenous, native language and the third most spoken language, after English and French. The Hausa language belongs to the West Chadic branch of the Afro-Asiatic phylum (Eberhard et al., 2020; Sabiu et al., 2018). It is spoken by about 150 million people and most of the speakers reside in Nigeria, Niger, Chad, Ghana, Benin, Burkina Faso and

Cameroon - all resulting in Arabic, Anglophone and Francophone influences. It has five alphabetic vowel symbols (a, e, i, o, u) when written in Romanized Hausa writing system. However, these symbols can have either short or long variants in pronunciation, leading to ten monophthongs. In addition, there are four vowel diphthongs (ai, au, iu, ui), making it a total of fourteen vowel phonemes. There are about about 23-25 consonant phonemes in the Hausa language depending on the speaker. Hausa is a tonal language with three tones: high, low, and falling. The high tone is always left unmarked, the low tone is represented by a grave accent (̀), and the falling tone combines the high and low tones and it is represented by a circumflex (ˆ).

**Igbo** (Ásùsù Ìgbò) (ibo) belongs to the Benue-Congo group of the Niger-Congo language family and is spoken by over 27 million people (Eberhard et al., 2020). It is the primary language of the Igbo people, an ethnic group of southeastern Nigeria, but also spoken in some parts of Equatorial Guinea and Cameroon. There are approximately 30 Igbo dialects, some of which are not mutually intelligible: their major differences are lexical and phonological (MustGo et al, 2015). This large number of dialects led to the development of a standardized spoken and written Igbo language, known as 'Standard Igbo' in 1962. In the early 1960s, the Nigerian government committee created the Ọnwụ orthography as the official Igbo orthography (Ohiri-Aniche, 2007). Ọnwụ consists of 28 consonants and 8 vowels. Igbo has some interesting acoustic properties. The Standard Igbo consists of eight vowels, thirty consonants and has two syllable types: consonant + vowel (the most common syllable type), vowel or syllabic nasal. Every syllable in Igbo has a tone. There are two tones: high and low. High tone is marked with an acute accent, e.g., á, while low tone is marked with a grave accent, e.g., à. These are not normally represented in the orthography. An interesting feature of Igbo phonology is tonal downstep (Beckman and Pierrehumbert, 1986). For example, two adjacent high-tone syllables will normally be produced with the same pitch. However, if a low-tone syllable occurs between two high tones, then the second high tone will be produced with a lower pitch than the first one, e.g., áá will become ááá.

The **Kinyarwanda language** (kin) belongs to the Niger-Congo family of languages, it is one

<sup>7</sup><https://github.com/geezorg/data/tree/master/amharic/tagged/nmsu-say>

of the four official languages spoken in Rwanda (along with English, French and Swahili) with at least 30 million speakers of which 12 million are natives. Other Kinyarwanda speakers are from the three countries neighboring Rwanda, i.e., Democratic Republic of the Congo, Uganda, and Tanzania (Niyongabo et al., 2020). Based on its morphological rich and tonal system, it is considered to be a more generic prototype of the larger group of Bantu languages, such as Kinyamurenge, Kinyabwisha, Rufumbira, and Ha (Kimenyi, 1980; Muhirwe, 2009; Nzeyimana, 2020). Moreover, it is mutually intelligible with Kirundi, an official language of Burundi.

**Luganda** (lug) is a Bantu language that is widely spoken in the African Great Lakes region. It is one of the major languages in Uganda primarily spoken in the south eastern Buganda region mainly along the shores of Lake Victoria and up north towards the Lake Kyoga shores (Nakayiza, 2013; Olutola et al., 2019). Luganda is spoken by more than six million people principally in central Uganda and other parts. It is the most widely spoken indigenous language and the most widely spoken second language after English and Swahili which are official languages in Uganda (Merritt and Abdulaziz, 1985). It belongs to the Bantu branch of the Niger-Congo language family. Typologically, it is a highly-agglutinating, tonal language with subject-verb-object, word order, and nominative-accusative morphosyntactic alignment (Olutola et al., 2019).

**Dholuo** (luo): Dholuo or Nilotic Kavirondo is a dialect of the Luo group of Nilotic languages. It is spoken by about 5.2 million Luo people of Kenya and Tanzania who mainly occupy the eastern and southern shores of lake Victoria. It is also spoken as a second, third or fourth language by the neighbouring communities of Suba, Kisii, Kuria and Luhya. Suba people who speak Dholuo as their first language are called Luo-Abasuba.

Dholuo is a tonal language with 4 tones (high, low, falling, rising) although the tonality is not marked in orthography. It uses the Latin alphabet without extra diacritics. There are 26 consonants without Latin letters (c, q, v, x and z) and additional (ch, dh, mb, nd, ng', ng, ny, nj, th, sh). There are nine vowels (a, e, i, o, u, e, ε, ɔ, ɔ) which are distinguished primarily by [ATR] harmony (De Pauw et al., 2007). Dholuo sentences follow Subject-Verb-Object word order.

Dholuo is mutually intelligible with Alur, Lango, Acholi and Adhola of Uganda spoken in Uganda. Kenyan Dholuo has two dialects, namely the Trans-Yala dialect spoken in Ugenya, Alego, Yimbo and parts of Gem; and the South Nyanza dialect spoken in South Nyanza, Siaya and Kisumu (De Pauw et al., 2007).

**Nigerian-Pidgin** (pcm): Pidgin English is broadly spoken across West Africa, with some 75 million speakers in Nigeria and approximately 5 million speakers in Ghana, as of 2016. Although it is a second (L2) language for a majority of the speakers, there are some 3-5 million speakers for whom it is the first language (BBC, 2016). Emerging as a lingua franca between European and African during the Atlantic Slave Trade of the 15th to 19th centuries, West African Pidgin evolved in a linguistically complex continuum where hundreds of languages are spoken. However, it maintains a consistency across the region, in part due to its use among Africans speaking mutually unintelligible languages (Peter and Wolf, 2007; Onovbiona, 2012). The Nigerian-Pidgin (NP) that has been the focus of this work, has been recorded as early as the 15th century along the coastal areas of Badagry, Sapele, Warri, Port-Harcourt and Calabar. While NP has English as the lexifier language, Portuguese, French and especially indigenous languages form the substrate of lexical, phonological, syntactic and semantic influence (Offiong Mensah, 2012; Onovbiona, 2012). For example, NP is rhythmically dissimilar to English, as its cadence and intonation are tone based (Ojarikre, 2013). Moreover, early contact with the Portuguese produced such common words in the NP lexicon: pikin (*child*) is derived from the Portuguese word pequenino; palaver (*problem*) from palavra; dash (*a gift*) from doação and sabi (*to know*) from saber (Vanamali, 1993). Finally, like many pidgins, NP has few prepositions, while tense and aspect are non-inflectional (Sebba, 1997; Nordquist, 2020).

**Swahili** (swa): Swahili, also known by its native name Kiswahili, is one of the most spoken languages in Africa. It is spoken by 100–150 million people across East Africa. Swahili is spoken by countries such as Tanzania, Kenya, Uganda, Rwanda, and Burundi, some parts of Malawi, Somalia, Zambia, Mozambique and the Democratic Republic of the Congo (DRC). Swahili is also one of the working languages of the African Union and

officially recognized as a lingua franca of the East African Community. In 2018, South Africa legalized the teaching of Swahili in South African schools as an optional subject to begin in 2020. The Southern African Development Community (SADC) officially recognized the Swahili as their official language. The Swahili alphabet consists of five vowels and nineteen consonants, almost like the English alphabet. The vowels, as we mentioned before, include a, e, i, o, u. The consonants include b, c, d, f, g, h, j, k, l, m, n, p, r, s, t, v, w, y, z. The consonants that are missing are q and x. The additions of ch, dh, gh, ng', sh, th and others are considered consonants because of the unique Swahili pronunciation.

**Wolof** (wo1): As a West-Atlantic language mainly spoken in Senegal and Gambia, Wolof (Ouolof, Volof, Walaf, Waro-Waro, Yallof) is also used in the Southern part of Mauritania. It belongs to the Atlantic group of the Niger-Congo language family and over seven million people spreading across three West African states is currently speaking Wolof. While only about 40% of the Senegalese population are Wolof, about 90% of the people speak the language as either their first, second or third language (Comrie, 2009). There are two major geographical varieties of Wolof: one spoken in Senegal, and the other spoken in The Gambia (Eberhard et al., 2020). Even if people who speaks Wolof understand each other, the Senegalese Wolof and the Gambian Wolof are two distinct languages: both own their ISO 639-3 language code (respectively “WOL” and “WOF”) (Gauthier et al., 2016). Within Senegal, (Eberhard et al., 2020) distinguishes five dialects: Baol, Cayor, Dylof, Lebou, and Jander. Although it has a long tradition of writing using the Arabic script known as Ajami or Wolofal, it has also been adapted to Roman script. The Wolof alphabet is quite close to the French one, we can find all the letters of its alphabet except H, V and Z. It also includes the characters ð (“ng”, lowercase: ŋ) and ñ (“gn”, as in Spanish). Accents are present, but in limited number (À, É, È, Ó). Twenty nine (29) Roman-based characters are used from the Latin script and most of them are involved in digraphs standing for geminate and prenasalized stops. Unlike many other Niger-Congo languages, Wolof does not have tones. Nevertheless, Wolof syllables differ in intensity, e.g., long vowels are pronounced with more intensity than short ones.

Length is represented by double vowel letters in writing and most Wolof consonants can be also geminated (doubled). However, Wolof is not a standardized language (and some sources exclude the “H” from the alphabet) since no single variety has ever been accepted as the norm. Nonetheless, the Center of Applied Linguistics of Dakar (CLAD), coordinates the orthographic standardization of the Wolof language (cla).

The **Yorùbá language** (yor) is the third most spoken language in Africa, and is native to the south-western Nigeria and the Republic of Benin. It is one of the national languages in Nigeria, Benin and Togo, and it is also spoken in other countries like Ghana, Côte d’Ivoire, Sierra Leone, Cuba, and Brazil. The language belongs to the Niger-Congo family, and is spoken by over 40 million native speakers (Eberhard et al., 2020). Yorùbá has several dialects but the written language has been standardized by the 1974 Joint Consultative Committee on Education (Asahiah et al., 2017), it has 25 Latin letters without the Latin characters (c, q, v, x and z) and with additional letters (e, gb, s, o). Yorùbá is a tonal language with three tones: low, middle and high. These tones are represented by the grave (“`”), optional macron (“—”) and acute (“’”) accents respectively. These tones are applied on vowels and syllabic nasals, but the mid tone is usually ignored in writings. Yorùbá is a highly isolating language and the sentences structure follows the Subject-Verb-Object.

## B Annotator Agreement

To shed more light on the few cases where annotators disagreed, we provide entity-level confusion matrices across all ten languages in Table 9. The most common disagreement is between organizations and locations.

	DATE	LOC	ORG	PER
DATE	32,978	-	-	-
LOC	10	70,610	-	-
ORG	0	52	35,336	-
PER	2	48	12	64,216

Table 9: Entity-level confusion matrix between annotators, calculated over all ten languages.

Lang.	WikiANN (zero-shot)			eng-CoNLL (zero-shot)			pcm (zero-shot)			swa (zero-shot)			hau (zero-shot)		
	PER	ORG	LOC	PER	ORG	LOC	PER	ORG	LOC	PER	ORG	LOC	PER	ORG	LOC
amh	41	15	29	–	–	–	–	–	–	–	–	–	–	–	–
hau	–	–	–	64	31	81	68	45	77	82	63	89	–	–	–
ibo	32	10	16	54	47	30	41	67	24	55	67	35	70	69	38
kin	27	3	2	35	29	50	29	59	55	59	60	56	63	62	57
lug	–	–	–	46	22	37	42	55	45	56	64	55	58	59	55
luo	–	–	–	53	20	16	40	18	25	44	31	22	64	33	26
pcm	–	–	–	70	64	71	–	–	–	77	68	66	81	63	65
swa	48	13	57	78	60	80	82	59	81	–	–	–	92	73	85
wol	–	–	–	31	18	9	29	30	9	41	41	19	43	42	44
yor	27	3	7	52	21	37	28	52	39	57	46	41	62	45	66
Ave.				45.2	26.2	29.8	34.8	46.8	32.8	52	51.5	38	60	51.7	47.7

Table 10: Per-named entity F1-score for the zero-shot NER using XLM-R model. Average computed over ibo, kin, lug, luo, wol, yor languages

Language	F1-score
amh	75.76
hau	91.75
ibo	86.26
kin	76.38
lug	84.64
luo	80.65
pcm	89.55
swa	89.48
wol	70.70
yor	82.05

Table 11: Evaluation of *xlm-roberta-large* NER model on different language test set

## C Zero-shot NER – Per-Named Entity F1-score

Table 10 shows the per-named entity F1-score for the zero-shot NER models when performing transfer learning from WikiANN, English, Nigerian-Pidgin, Swahili and Hausa languages.

## D Model Hyper-parameters for Reproducibility

For fine-tuning mBERT and XLM-R, we used the base models with maximum sequence length of 164, batch size of 32, learning rate of 5e-5, and number of epochs 50. For the MeanE-BiLSTM model, the hyper-parameters are similar to fine-tuning the LM except for the learning rate that we set to be 5e-4, the BiLSTM hyper-parameters are: input dimension is 768 (since the embedding size from mBERT and XLM-R is 768) in each direction of LSTM, one hidden layer, hidden layer size of 64, and drop-out probability of 0.3 before the last linear layer. All the experiments were performed on a single GPU (Nvidia V100).

## E XLM-RoBERTa-Large Model for all languages

We released a single NER model trained on the aggregation of all the ten African languages. The model is trained by fine-tuning the *xlm-roberta-large* model. Table 11 shows the performance across different languages.