

# One Country, 700+ Languages: NLP Challenges for Underrepresented Languages and Dialects in Indonesia

Anonymous ACL submission

## Abstract

NLP research is impeded by a lack of resources and awareness of the challenges presented by under-represented languages and dialects. Focusing on the languages spoken in Indonesia, the second most linguistically diverse and the fourth most populous nation of the world, we provide an overview of the current state of NLP research for Indonesia’s 700+ languages. We highlight challenges in Indonesian NLP and how these affect the performance of current NLP systems. Finally, we provide general recommendations to help develop NLP technology not only for languages of Indonesia, but also other underrepresented languages.

## 1 Introduction

Research in natural language processing (NLP) has traditionally focused on developing models for English and a small set of other languages with large amounts of data (see Figure 1, bottom right). While the lack of data is generally cited as the key reason for the lack of progress in NLP for underrepresented languages (Hu et al., 2020; Joshi et al., 2020), we argue that another factor relates to the diversity and the lack of understanding of the linguistic characteristics of such languages. Through the lens of the languages spoken in Indonesia, the world’s second-most linguistically diverse country, we seek to illustrate the challenges in applying NLP technology to a such diverse pool of languages.

Indonesia is the 4th most populous nation globally with 273 million people spread over 17,508 islands. There are more than 700 languages spoken in Indonesia, equal to 10% of the world’s languages, second only to Papua New Guinea (Eberhard et al., 2021). However, most of these languages are not well documented in the literature; many are not formally taught and no established standard exists across speakers (Novitasari et al., 2020). Many of them are decreasing in use, as Indonesian (*Bahasa Indonesia*), the national lan-

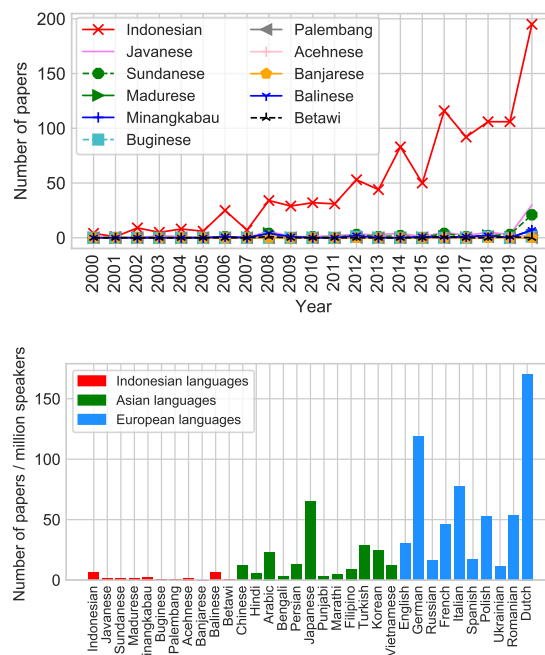


Figure 1: Following Joshi et al. (2020), we compile ACL Anthology (main conferences and workshops) to count the distribution of published works that mention Indonesian languages. **Top:** Distribution of papers in 20 years. **Bottom:** Distribution per million speakers, compared to other Asian and European languages.

guage, is more frequently used as the primary language across the country. This process may ultimately result in a monolingual society (Cohn and Ravindranath, 2014). One study finds that among 98 Indonesian local languages, 36 are considered safe, 51 are endangered, and 11 are already extinct (Anindyatri and Mufidah, 2020).

Table 1 shows the 10 Indonesian local languages with the most speakers, along with Indonesian for comparison (Eberhard et al., 2021). Javanese and Sundanese are at the top with 84M and 34M speakers respectively, while Madura, Minangkabau, and Buginese have around 6M speakers. Despite their large speaker populations, these local languages are poorly represented in the NLP literature. Compared to Indonesian, the number of research papers

Language	ISO	# Speakers
Indonesian	id	198 M
Javanese	jav	84 M
Sundanese / Sunda	su	34 M
Madurese / Madura	mad	7 M
Minangkabau	min	6 M
Buginese	bug	6 M
Betawi	bew	5 M
Acehnese / Aceh	ace	4 M
Banjar	bjn	4 M
Balinese	ban	3 M
Palembang Malay (musi)	mus	3 M

Table 1: The 10 most spoken Indonesian local languages according to Ethnologue (Eberhard et al., 2021).

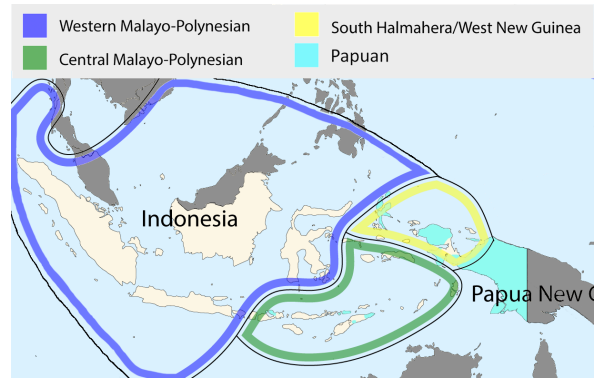


Figure 2: Map of Austronesian and Papuan languages in Indonesia

057 mentioning these languages has barely increased  
058 over the past 20 years (Figure 1, top). Furthermore,  
059 compared to their European counterparts, Indone-  
060 sian languages are drastically understudied (Figure  
061 1, bottom). This is true even for Indonesian, which  
062 has nearly 200M speakers.

063 Language technology should be accessible to  
064 everyone in their native languages (European Lan-  
065 guage Resources Association, 2019), including In-  
066 donesians. In the context of Indonesia, language  
067 technology research offers some benefits. First,  
068 language technology is one of the potential peace-  
069 maker tools in a multi-ethnic country, helping In-  
070 donesians understand each other better and avoid  
071 the ethnic conflicts of the past (Bertrand, 2004).  
072 On a larger scale, language technology promotes  
073 language use (European Language Resources As-  
074 sociation, 2019) and helps language preservation.  
075 Despite these benefits, following Bird (2020), we  
076 recommend a careful assessment of individual use  
077 scenarios of language technology so they are  
078 implemented for the good of the local population.

079 For language technology to be useful in the In-  
080 donesian context, it additionally has to account  
081 for the dialects of local languages. Dialects in  
082 Indonesia are influenced by the geographical lo-  
083 cation and regional culture of their speakers (Van-  
084 der Kloek, 2015) and thus often differ substantially  
085 in morphology and vocabulary, posing challenges  
086 for NLP systems. In this paper, we provide an  
087 overview of the current state of NLP for Indone-  
088 sian languages. We then discuss the challenges  
089 presented by those languages and demonstrate how  
090 they affect state-of-the-art systems in NLP. We fi-  
091 nally provide recommendations for developing bet-  
092 ter NLP technology not only for languages in In-  
093 donesia but also other under-represented languages.

## 2 Background and Related Work 094

### 2.1 History and Taxonomy 095

096 Indonesia is one of the richest countries in the  
097 world in terms of linguistic diversity. More  
098 than 400 of its languages belong to the Aus-  
099 tronesian language family, while the others are  
100 Papuan languages spoken in the eastern part of  
101 the country. As shown in Figure 2, the Aus-  
102 tronesian languages in Indonesia belong to three  
103 main groups: Western-Malayo-Polynesian (WMP),  
104 Central-Malayo-Polynesian (CMP), and South-  
105 Halmahera-West-New-Guinea (SHWNG) (Blust,  
106 1980). WMP languages are Malay, Indonesian, Ja-  
107 vanese, Sundanese, Balinese, and Minangkabau,  
108 among others. All languages mentioned in Table 1  
109 are in this group. Languages belonging to CMP  
110 are languages of the Lesser Sunda Islands from  
111 East Sumbawa (with Bimanese) onwards to the  
112 east, and languages of the central and southern  
113 Moluccas (including the Aru Islands and the Sula  
114 Archipelago). The SHWNG group consists of lan-  
115 guages of Halmahera and Cenderawasih Bay, as  
116 far as the Mamberamo River, and of the Raja Am-  
117 pat Islands. The Papuan languages, meanwhile,  
118 are mainly spoken in Papua, such as Dani, Asmat,  
119 Maybrat, and Sentani. Some Papuan languages  
120 are also spoken in Halmahera, Timor, and the Alor  
121 Archipelago (Palmer, 2018; Ross, 2005).

122 Most Austronesian linguists and archaeologists  
123 agree that the original ‘homeland’ of Austronesian  
124 languages must be sought in Taiwan and, prior to  
125 Taiwan, in coastal South China (Adelaar, 2005;  
126 Bellwood and Dizon, 2008; Bellwood et al., 2011).  
127 The Austronesian people moved from Taiwan to the  
128 Philippines in the second millennium CE. From the  
129 Philippines, they moved southward to Borneo and

130 Sulawesi. From Borneo, they migrated to Sumatra, 131 the Malay Peninsula, Java, and even to Madagascar. 132 From Sulawesi, they moved southward to the CMP 133 area and eastward to the SHWNG area. From there, 134 they migrated to Oceania and Polynesia, as far as 135 New Zealand, Easter Island, and Hawaii (Gray and 136 Jordan, 2000). The people that lived in insular 137 Southeast Asia, such as in the Philippines and In- 138 donesia, before the arrival of Austronesians were 139 Australo-Melanesians (Bellwood, 1997). Gradual 140 assimilation with Austronesians occurred although 141 some pre-Austronesian groups still survive such 142 as Melanesian people in eastern Indonesia (Ross, 143 2005; Coupe and Kratochvíl, 2020).

144 At the time of the arrival of the first Europeans, 145 Malay had become the major language (lingua 146 franca) of interethnic communication in Southeast 147 Asia and beyond (Steinhauer, 2005; Coupe and 148 Kratochvíl, 2020). It functioned as the language of 149 trade and the language of Islam because Muslim 150 merchants from India and the Middle East were 151 the first to introduce the religion into the harbor 152 towns of Indonesia. After the arrival of Europeans, 153 Malay was used by the Portuguese and Dutch 154 to spread Catholicism and Protestantism. When 155 the Dutch extended their rule over areas outside 156 Java in the nineteenth century, the importance of 157 Malay increased, and thus, the first standardiza- 158 tion of the spelling and grammar occurred in 1901, 159 based on Classical Malay (Abas, 1987; Sneddon, 160 2003). In 1928, the participants of the Second Na- 161 tional Youth Congress proclaimed Malay (hence- 162 forth called Indonesian) as the unifying language 163 of Indonesia. During World War II, the Japanese 164 occupying forces forbade all use of Dutch in favor 165 of Indonesian, which from then onward effectively 166 became the new national language. After indepen- 167 dence until the present, Indonesian has functioned 168 as the main language in education, mass media, and 169 government activities. Many local language speak- 170 ers are increasingly using Indonesian with their 171 children because they believe it will help them to- 172 ward a better education and career (Klamer, 2018).

## 173 2.2 Efforts in Multilingual Research

174 Recently, pretrained multilingual language models 175 such as mBERT (Devlin et al., 2019), mBART (Liu 176 et al., 2020), and mT5 (Xue et al., 2021b) were pro- 177 posed. Their coverage, however, focuses on high- 178 resource languages. Among them, only mBERT 179 and mT5 include Indonesian local languages, i.e.,

Javanese, Sundanese, and Minangkabau, but with 180 comparatively little pretraining data. 181

182 Some multilingual datasets for question answer- 183 ing (TyDi QA; Clark et al., 2020), dia- 184 logue (XPersona; Lin et al., 2021), passage 185 ranking (mMARCO; Bonifacio et al., 2021), 186 cross-lingual visual question answering (xGQA; 187 Pfeiffer et al., 2021), common sense reason- 188 ing (XCOPA; Ponti et al., 2020), abstractive sum- 189 marization (Hasan et al., 2021), language and vi- 190 sion reasoning (MaRVL; Liu et al., 2021), and ma- 191 chine translation (FLORES-101; Goyal et al., 2021) 192 include Indonesian but most others do not, and very 193 few include Indonesian local languages. An excep- 194 tion is the weakly supervised named entity recog- 195 nition dataset, WikiAnn (Pan et al., 2017), which 196 covers several Indonesian local languages, namely 197 Acehnese, Javanese, Minangkabau, and Sundanese.

198 Parallel corpora including Indonesian local lan- 199 guages are i) CommonCrawl; ii) Wikipedia parallel 200 corpora like MediaWiki Translations<sup>1</sup> and Wiki- 201 Matrix (Schwenk et al., 2021); iii) the Leipzig 202 corpora (Goldhahn et al., 2012), which include 203 Indonesian, Javanese, Sundanese, Minangkabau, 204 Madurese, Acehnese, Buginese, Banjar, and Bali- 205 nese; and iv) JW-300 (Agić and Vulić, 2019), which 206 includes dozens of Indonesian local languages, e.g., 207 Batak language groups, Javanese, Dayak language 208 groups, and several languages in Nusa Tenggara.<sup>2</sup>

## 209 2.3 Progress in Indonesian NLP

210 Most NLP research on Indonesian has been done 211 across multiple topics, such as sentiment analy- 212 sis (Naradhipa and Purwarianti, 2011; Lunando and 213 Purwarianti, 2013), hate speech detection (Alfina 214 et al., 2017; Ibrohim and Budi, 2019; Sutejo and 215 Lestari, 2018), morphological analysis (Pisceldo 216 et al., 2008), POS tagging (Wicaksono and Purwari- 217 anti, 2010; Dinakaramani et al., 2014; Kurniawan 218 and Aji, 2018), named entity recognition (Budi 219 et al., 2005; Gunawan et al., 2018), question an- 220 swering (Mahendra et al., 2008; Fikri and Purwari- 221 anti, 2012), machine translation (Yulianti et al., 222 2011), and speech recognition (Lestari et al., 2006; 223 Baskoro and Adriani, 2008; Zahra et al., 2009). 224 However, many of these studies either kept the data 225 private or used non-standardized resources with 226 a lack of documentation and open-sourced code,

<sup>1</sup>[https://mediawiki.org/wiki/Content\\_translation](https://mediawiki.org/wiki/Content_translation)

<sup>2</sup>Recent studies (Caswell et al., 2021), however, have raised concerns regarding the quality of such multilingual corpora for under-represented languages.

which makes them extremely difficult to reproduce. Recently, Wilie et al. (2020), Koto et al. (2020b), and Cahyawijaya et al. (2021) collect Indonesian NLP resources as benchmark data. Others have also begun to create standardized labelled data for Indonesian NLP, e.g. the works of Kurniawan and Louvan (2018), Guntara et al. (2020), Mahendra et al. (2021), Koto et al. (2021), and Artari et al. (2021).

On the other hand, a handful of NLP research explore the local languages. Suryani et al. (2015) study machine translation in Sundanese by using prior POS tagging information, while Suryani et al. (2018) develop a word stemmer for Sundanese. Koto and Koto (2020) explore sentiment analysis and machine translation for Minangkabau. Safitri et al. (2016) work on spoken data language identification in three Indonesian local languages, i.e., Minangkabau, Sundanese and Javanese. Azizah et al. (2020) develop end-to-end neural text-to-speech model for Indonesian, Sundanese, and Javanese. Recently, Cahyawijaya et al. (2021) established a machine translation benchmark in Sundanese and Javanese using Bible data. Wibowo et al. (2021) studied a family of colloquial Indonesian, which is influenced by some local languages via morphological transformation, and Putri et al. (2021) worked on abusive language and hate speech detection on Twitter for five local languages, namely Javanese, Sundanese, Madurese, Minangkabau, and Musi.

### 3 Challenges for Indonesian NLP

#### 3.1 Limited Resources

**Monolingual Data** Unlabelled corpora are crucial for building large language models, such as GPT-2 (Radford et al., 2019) or BERT (Devlin et al., 2019). Available unlabelled corpora such as Indo4B (Wilie et al., 2020) and Indo4B-Plus (Cahyawijaya et al., 2021) mainly include data in Indonesian, with the latter containing  $\approx 10\%$  of data in Javanese and Sundanese (see Appendix C). In comparison, in multilingual corpora such as CC-100 (Conneau et al., 2020), Javanese and Sundanese data accounts for only 0.001% and 0.002% of the corpus size while in mC4 (Xue et al., 2021b), there are only 0.6M Javanese and 0.3M Sundanese tokens out of a total of 6.3T tokens.

In addition, we measure data availability in Wikipedia, compared to the number of speakers as in Figure 3.<sup>3</sup> Among highly spoken local lan-

<sup>3</sup>The number of speakers is collected from Wiki-

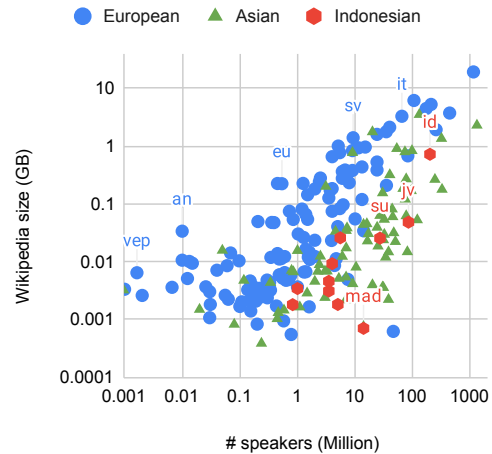


Figure 3: Wikipedia data size (in GB) compared to the number of speakers.

guages, much fewer data is available for Indonesian languages, compared to European languages with similar numbers of speakers. For example, Wikipedia contains more than 3 GB of Italian articles but less than 50 MB of Javanese articles, despite both languages having a comparable number of speakers. Similarly, Sundanese has less than 25 MB of articles, whereas languages of comparable speakers have more than 1.5 GB of articles. Similar trends hold for most other Asian languages.<sup>4</sup>

Beyond the most spoken local languages, other Indonesian local languages do not have Wikipedia instances, in contrast to European languages with few speakers. However, it is very difficult to find alternative sources for high-quality text data for other local languages of Indonesia (such as news websites), as most such sources are written in Indonesian. Resources in tail languages are even more lacking, due to a very low number of speakers. Moreover, most of the languages in the long tail are mainly used in a spoken context, making text data difficult to obtain.

These statistics demonstrate that collecting unlabelled corpora for Indonesian local languages is extremely difficult. This makes it impractical to develop strong pretrained language models for these languages, which have been the foundation for many recent state-of-the-art NLP systems.

**Labelled Data** Most work on Indonesian NLP (see §2) did not publicly release their data or mod-

data (Vrandečić and Krötzsch, 2014), from the number of speakers (P1098) property as of Nov 7th 2021, while the size is collected from the 20211101 Wikipedia dump.

<sup>4</sup>Other continents such as Africa are even more under-represented in terms of Wikipedia data (see Appendix D).

English	Mudung Laut	Dusun Teluk	Mersam	Suo Suo	Teluk Kual	Lubuk Telau	Bunga Tanjung	Pulau Aro
I/me	sayo	aku	awa?	sayo	kito, awa?	am <sup>b</sup> o	ambo	ambo
You	kau, kamu	kau	ka <sup>d</sup> n	kamu	kaan	kamu	aj, kau, kayo	ba?aj
he/she	dio?	dio?, jo	jo	kau	jo	jo	jo	ipo
if	kalu	jiko, kalu	kalu	bilao	kalu	jiko	ko?	kalu
one	satu	seko?	seko?	seko?	ciz?	seko?	seko?, so	seko?

Table 2: Lexical variation of Jambi Malay across different villages in Jambi, collected from [Anderbeck \(2008\)](#).

English	Context	Ngoko			Krama
		Western	Central	Eastern	Eastern
I/me	I like to eat fried rice.	inyong, enyong	aku	aku	kulo
You	Where will you go?	rika, kowe, ko	kowe, siro, sampeyan	koen, awakmu, sampeyan	panjenengan
How	How do I read this?	priwe	piye	yo’opo	pripun
Why	Why is this door broken?	ngapa	ngopo	opo’o	punapa
Will	Where will you go?	arep	arep	kate, ate	badhe
Not/no	The calculation is not correct.	ora	ora	gak	mboten

Table 3: Lexical variation of Javanese dialects and styles across different regions of Java island. Native speakers were asked to translate the words, given the context.

els, thus limiting reproducibility. Although recent Indonesian NLP benchmarks are addressing this issue, they mostly focus on the Indonesian language (see Appendix C). Some widely spoken local languages such as Javanese, Sundanese, or Minangkabau have extremely small labelled datasets compared to Indonesian, while others have barely any.

The lack of such datasets makes NLP development for the local languages difficult. However, constructing new labelled datasets is still challenging due to (1) the lack of speakers of some languages, (2) the vast continuum of dialectical variation (see §3.2.1), and (3) the absence of writing standard in most local languages (see §3.3).

### 3.2 Language Diversity

The diversity of Indonesian languages is not only due to the large number of local languages, but also the large number of dialects of these languages (§3.2.1). Speakers of local languages also often mix languages in conversation, which makes colloquial Indonesian more diverse (§3.2.2). In addition, some local languages are more commonly used in conversational contexts, so they do not have consistent writing forms in written media (§3.3).

#### 3.2.1 Regional Dialects and Style Differences

Indonesian local languages often have multiple dialects, depending on the geographical location. Local languages of Indonesian spoken in different locations might be different (have some lexical variation) to one another, despite still being categorized as the same language ([Fauzi and Puspitorini, 2018](#)). For example, [Anderbeck \(2008\)](#) shows that

villages across the Jambi province use different dialects of Jambi Malay. Similarly, [Kartikasari et al. \(2018\)](#) show that Javanese between different cities in central and eastern Java can have more than 50% lexical variation, while [Purwaningsih \(2017\)](#) shows that Javanese in different districts in the Lamongan Regency has up to 13% lexical variation. Similar studies have been conducted on other languages, such as Balinese ([Maharani and Candra, 2018](#)) and Sasak ([Sarwadi et al., 2019](#)).

Moreover, Indonesian and its local languages have multiple styles, even within the same dialect. One factor that affects style is the level of politeness and formality—similar to Japanese and other Asian languages ([Bond and Baldwin, 2016](#)). More polite language is used when speaking to a person with a higher social position, especially to elders, seniors, and sometimes strangers. Different politeness levels manifest in the use of different honorifics and even different lexical terms.

To illustrate the distinctions between regional dialects and styles, we highlight common words and utterances across dialects and styles in Jambi Malay and Javanese in Tables 2 and 3. For Jambi Malay, we sample the result from [Anderbeck \(2008\)](#). For Javanese, we ask native speakers to translate basic words across three regional dialects: Western, Central, and Eastern Javanese, and two different styles: *Ngoko* (standard, daily-use Javanese) and *Krama* (polite Javanese, used to communicate to elders and those with higher social status). However, since contemporary *Krama* Javanese is not very different between regions, we only consider *Krama* from Eastern speaker’s perspective.

Style	Region	Model				
		Langid.py		FastText		CLD3
		Top-1	Top-3	Top-1	Top-3	Top-1
<i>Ngoko</i>	Western	0.241	0.621	0.069	0.379	0.759
<i>Ngoko</i>	Central	0.345	0.690	0.379	0.724	0.828
<i>Ngoko</i>	Eastern	0.276	0.552	0.103	0.379	0.552
<i>Krama</i>	Eastern	0.345	0.759	0.379	0.586	0.897

Table 4: Language identification accuracy based on different Javanese dialects and styles. Systems do not perform equally well across dialects and styles.

Jambi Malay is not widely spoken (1M speakers), but has many dialects across villages. As shown in Table 2, many common words are spoken differently across dialects and styles. Similarly, Javanese is also different across regions. Not every Javanese speaker understands *Krama*, since its usage is very limited. Moreover, the number of Javanese speakers who can use *Krama* is declining (Cohn and Ravindranath, 2014).<sup>5</sup> Examples from other languages are shown in Appendix E.

### Case Study in Javanese

Dialectical and style differences pose a challenge to NLP systems. To explore the extent of this challenge, we conduct an experiment to test the robustness of NLP systems to variations in Javanese dialects. We ask native speakers<sup>6</sup> to translate 29 simple sentences into Javanese according to the specified dialect and style. We then evaluate several language identification systems on those instances. Language identification is a core part of multilingual NLP and a necessary step for collecting textual data in a language. Despite its importance, it is an open research area, particularly for under-represented languages (Caswell et al., 2020).

We compare Langid.py (Lui and Baldwin, 2012), FastText (Joulin et al., 2017), and CLD3.<sup>7</sup> The results can be seen in Table 4. In general, the language identification systems are more accurate in detecting Javanese texts in the *Ngoko*-Central dialect, or *Krama*, since the systems were trained on Javanese Wikipedia data, which is written in either

<sup>5</sup>*Krama* is used to speak formally (e.g., with older or respected people). Nowadays, however, people prefer to use Indonesian more in formal situation. People who move from sub-urban areas to bigger cities tend to continue to use *Ngoko* and thus also pass *Ngoko* on to their children.

<sup>6</sup>Our annotators are based in Banyumas for Western Javanese, Jogjakarta for Central Javanese, and Jember for Eastern Javanese. Using dialects from different cities might result in a slightly different result.

<sup>7</sup><https://github.com/google/cld3>

Colloquial Indonesian	Translation
Ada yang <b>ngetag</b> foto <b>lawas</b> di FB	Someone is tagging old photos in FB
<b>Quotenya</b> Andrew Ng ini relevan <b>banget</b>	This Andrew Ng quote is very relevant
<b>Bilo</b> kita pergi main lagi?	When will we go play again?
Ini <b>teh</b> aksara jawa kenapa susah <b>banget</b> ?	Why is this Javanese script very difficult?

Table 5: Colloquial Indonesian code-mixing examples from social media. Color code: **English**, **Betawinese**, **Javanese**, **Minangkabau**, **Sundanese**, Indonesian.

the *Ngoko*-Central or *Krama* dialects and styles. If an NLP system can only detect certain dialects, then this information should be conveyed explicitly. Problems arise if we assume that the model works equally well across dialects. For example, in the case of language identification, if we use the model to collect datasets automatically, then Javanese datasets with poor-performing dialects will be under represented in the data.

### 3.2.2 Code-Mixing

Code-mixing is an occurrence where a person speaks alternately in two or more languages in a conversation (Poplack, 1980; Winata et al., 2019, 2021a). This phenomenon is common in Indonesian conversations (Barik et al., 2019; Wibowo et al., 2020, 2021). In a conversational context, people sometimes mix their local languages with standard Indonesian, resulting in colloquial Indonesian (Siregar et al., 2014). This colloquial-style Indonesian is used daily in speech and conversation, and is common on social media (Sutrisno and Ariesta, 2019). Some frequently used code-mixed words (especially on social media) are even intelligible to people that do not speak the original local languages. Interestingly, code-mixing can also occur in border areas where people are exposed to multiple languages, therefore mixing them together. For example, people in Jember (a regency district in East Java) combine Javanese and Madurese in their daily conversation (Haryono, 2012).

Indonesian code-mixing not only occurs at the word level but also at the morpheme level (Winata, 2021). For example, *quotenya* (‘his/her quote’, see Table 5) combines the English word ‘quote’ and the Indonesian suffix *-nya*, which denotes possession; similarly, *ngetag* combines the Betawinese prefix *nge-* and the English word ‘tag’. More examples can be found in Table 5.

Language	Meaning	Written Variation	IPA
Javanese (Eastern– <i>Ngoko</i> )	what	apa / opo	/ɔpɔ/
	there is	ana / ono / onok	/ɔnɔʔ/
	you	kon / koen	/kɔn/
Balinese (Alus– <i>Singgih</i> )	yes	inggih / nggih	/ʔɪŋgih/
	I / me	tiang / tyang	/tiaŋ/
	<greeting>	swastyastu / swastiastu	/swastiastu/
Sundanese (Badui– <i>Loma</i> )	please / sorry	punten / punteun	/puntən/
	red	beureum / berem	/bɔrim/
	salivating	ngacai / ngacay	/ŋacai/

Table 6: Written form variations in several local languages, confirmed by native speakers.

### 3.3 Orthography Variation

Many Indonesian local languages are mainly used in spoken settings and have no established standard orthography system. Some local languages do originally have their own archaic writing systems that derive from the Jawi alphabet or Kawi script, and even though standard transliteration into the Roman alphabet exists for some (e.g., Javanese and Sundanese), they are not widely known and practiced (Soeparno, 2015). Hence, some words have multiple romanized writings that are mutually intelligible by speakers, as they are pronounced the same. Some examples can be seen in Table 6. Such variety of the written form is common in many local languages in Indonesia. This variation leads to a significantly larger vocabulary size, especially for NLP systems that use word-based representations, and results in words being spelled differently despite referring to the same word, a challenge for subword-based models.

### 3.4 Societal Challenges

Language evolves together with the speakers. A more widely used language may have a larger digital presence, which fosters a more written form of communication while languages that are used only within small communities may emphasize the spoken form. There are also languages that are declining, where the speakers prefer to use Indonesian rather than their local language. In contrast, there are isolated residents that use the local language daily and are less proficient in Indonesian (Nurjanah et al., 2018; Jahang and Meirina, 2021). These variations give rise to different requirements and there is no single solution for all.

Technology and education is not well-distributed within the nation. Internet penetration in Indonesia is 73.7% in 2020, but is mainly concentrated on the

Java island. Among the non-Internet users, 39% explain that they do not understand the technology, while 15% state that they do not have the device to access the internet.<sup>8</sup> In some areas where Internet is not seen as a basic need, imposing NLP technology on them may not necessarily be relevant. At the same time, general NLP development within the nation faces difficulties due to the lack of funding especially in universities outside of Java. GPU servers are still scarce, even on Java.<sup>9</sup>

The dynamics of population movement in Indonesia also need to be taken into consideration. For example, there are urban communities who transmigrate to remote areas for social purposes, such as teaching or becoming doctors for underdeveloped villages. Each of these situations might call for various new NLP technologies to be developed to facilitate better communication.

## 4 Opportunities

Based on the challenges for Indonesian NLP highlighted in the previous section, we formulate proposals for improving the state of Indonesian NLP research, as well as of other under-represented languages. Our proposals cover several aspects including metadata documentation; potential research directions; and engagement with communities.

### 4.1 Better Documentation

In line with studies promoting proper data documentation for NLP research (Bender and Friedman, 2018; Rogers et al., 2021; Alyafeai et al., 2021), we recommend the following considerations.

**Regional Dialect Metadata** We have shown that the same languages can have a large variation depending on region and dialect. Therefore, we suggest adding regional dialect metadata to NLP datasets and models, not only for Indonesian but for other languages as well. This is particularly important for languages with large dialectal differences. Regional dialect metadata is also important to clearly communicate NLP capabilities to stakeholders and end users as it will help set an expectation of what types of dialects systems can handle. Additionally, regional metadata can indirectly inform the topics of the data, especially for crawled data sources.

<sup>8</sup>The Indonesian Internet Providers Association (APJII) survey: <https://apjii.or.id/survei2019x>

<sup>9</sup>For instance, we estimate the whole computer science faculty of the nation’s top university owns 8 V100 GPUs.

**Style and Register Metadata** Similarly, we also suggest adding style and register metadata. This metadata can capture the politeness level of the text, not only for Indonesian but also other languages. In addition, this metadata can be used to document the formality level of the text, so may be useful for research on modeling style or style transfer.

## 4.2 Potential Research Direction

In Indonesia, there are only few widely spoken languages that have been investigated in NLP, while the rest remain unstudied. Mitigating this limitation, we suggest future research to focus more on under-represented and unexplored languages.

**Data-Efficient NLP** Pretrained language models, which have taken NLP world by storm, require an abundant amount of monolingual data. However, data collection has been a long-standing problem for low-resource languages. Therefore, we recommend more exploration into designing data-efficient approaches such as adaptation methods (Artetxe et al., 2020; Aji et al., 2020; Gururangan et al., 2020; Koto et al., 2021), few-shot learning (Winata et al., 2021b; Madotto et al., 2021; Le Scao and Rush, 2021), and learning from related languages (Khanuja et al., 2021; Khemchandani et al., 2021). The goal of these methods is effective resource utilization, that is, to minimize the financial costs for computation and data collection as advocated by Schwartz et al. (2020), Cahyawijaya (2021), and Nityasya et al. (2021).

**Data Generation** Data collection efforts need to be commenced as soon as possible, despite all the challenges (§3.1). Here, we suggest collecting parallel data between Indonesian and each of the local languages due to several reasons. First, a lot of Indonesians are bilingual (Koto and Koto, 2020), that is, they speak both Indonesian and their local language, which facilitates data collection. Moreover, the fact that the local languages have some vocabulary overlap with Indonesian (See Table 7 in Appendix) might help building translation systems using relatively fewer parallel data (Nguyen and Chiang, 2017). Finally, having such parallel data, we can build translation systems for synthetic data generation. In line with this approach, the effectiveness of models trained on synthetic translated dataset can be explored.

**Robustness to Code-mixing and Non-Standard Orthography** Languages in Indonesia are prone

to variations due to code-mixing and non-standard orthography, which occurs on the morpheme or even grapheme level. Models that are applied to Indonesian code-mixed data need to be able to learn morphologically faithful representations. Therefore, we recommend more exploration on methods derived from subword tokenization (Gage, 1994; Kudo, 2018) and token-free models (Gillick et al., 2016; Tay et al., 2021; Xue et al., 2021a) to deal with this problem.

**NLP Beyond Text** For many Indonesian local languages that are rarely if ever written, speech is a more natural communication format. We thus recommend more attention on less text-focused research, such as spoken language understanding (SLU) (Chung et al., 2021; Serdyuk et al., 2018), speech recognition (Besacier et al., 2014; Winata et al., 2020), and multimodality (Dai et al., 2020, 2021) in order to progress NLP in such languages.

## 4.3 Engage with Communities

As discussed in §3.4, it is difficult to generalize a solution across local languages. We thus encourage the NLP community to work more closely with native speakers and local communities (Nekoto et al., 2020). This is necessary to provide solutions and resources that support use cases benefiting the native speakers and communities of under-represented languages. We advise the involvement of linguists, for example to aid the language documentation process (Anastasopoulos et al., 2020). As GPU access can be a challenge for Indonesian research institutions, we suggest to engage with academic communities. We support open-science movements such as BigScience<sup>10</sup> or ICLR CoSubmitting Summer<sup>11</sup>, which help to start collaborations and to reduce the entry barrier to NLP research.

## 5 Conclusion

In this paper, we have highlighted challenges in Indonesian NLP. Indonesia is one of the most populous country and the second-most linguistically diverse, with over 700 local languages, yet Indonesian NLP is under represented and under explored. Based on the observed challenges, we have also presented recommendations to improve the situation, not only for Indonesian, but for other under represented languages as well.

<sup>10</sup><https://bigscience.huggingface.co/>

<sup>11</sup><https://blog.iclr.cc/2021/08/10/broadening-our-call-for-participation-to-iclr-2022/>



## References

- 619 Husen Abas. 1987. *Indonesian as a unifying language*  
620 *of wider communication: a historical and sociolin-*  
621 *guistic perspective*. Number 73 in Pacific Linguistics  
622 Series D. Canberra.
- 623 Alexander Adelaar. 2005. The Austronesian languages  
624 of Asia and Madagascar: A historical perspective.  
625 In Alexander Adelaar and Nikolaus P. Himmel-  
626 mann, editors, *The Austronesian Languages of Asia*  
627 *and Madagascar*, chapter 1, pages 1–42. Routledge,  
628 Oxon.
- 629 Željko Agić and Ivan Vulić. 2019. [JW300: A wide-](#)  
630 [coverage parallel corpus for low-resource languages](#).  
631 In *Proceedings of the 57th Annual Meeting of the*  
632 *Association for Computational Linguistics*, pages  
633 3204–3210, Florence, Italy. Association for Computa-  
634 tional Linguistics.
- 635 Alham Fikri Aji, Nikolay Bogoychev, Kenneth  
636 Heafield, and Rico Sennrich. 2020. In neural ma-  
637 chine translation, what does transfer learning trans-  
638 fer? In *Proceedings of the 58th Annual Meeting*  
639 *of the Association for Computational Linguistics*,  
640 pages 7701–7710.
- 641 Ika Alfina, Ruli Manurung, and Mohamad Ivan Fanany.  
642 2016. [DBpedia entities expansion in automatically](#)  
643 [building dataset for Indonesian NER](#). In *2016 In-*  
644 *ternational Conference on Advanced Computer Sci-*  
645 *ence and Information Systems (ICACIS)*, pages  
646 335–340.
- 647 Ika Alfina, Rio Mulia, Mohamad Ivan Fanany, and  
648 Yudo Ekanata. 2017. Hate speech detection in the  
649 indonesian language: A dataset and preliminary  
650 study. In *2017 International Conference on Ad-*  
651 *vanced Computer Science and Information Systems*  
652 *(ICACIS)*, pages 233–238. IEEE.
- 653 Zaid Alyafeai, Maraim Masoud, Mustafa Ghaleb, and  
654 Maged S Al-shaibani. 2021. Masader: Metadata  
655 sourcing for arabic text and speech data resources.  
656 *arXiv preprint arXiv:2110.06744*.
- 657 Antonios Anastasopoulos, Christopher Cox, Graham  
658 Neubig, and Hilaria Cruz. 2020. [Endangered lan-](#)  
659 [guages meet Modern NLP](#). In *Proceedings of*  
660 *the 28th International Conference on Computa-*  
661 *tional Linguistics: Tutorial Abstracts*, pages 39–45,  
662 Barcelona, Spain (Online). International Committee  
663 for Computational Linguistics.
- 664 Karl Ronald Anderbeck. 2008. *Malay dialects of the*  
665 *Batanghari river basin (Jambi, Sumatra)*. SIL Inter-  
666 national.
- 667 Anisya O. Anindyatri and Imarotul Mufidah. 2020.  
668 *Gambaran Kondisi Vitalitas Bahasa Daerah di In-*  
669 *donesia*. Kementerian Pendidikan dan Kebudayaan  
670 Pusat Data dan Teknologi Informasi, Tangerang Se-  
671 latan, Indonesia.
- Valentina Kania Prameswara Artari, Rahmad Ma- 672  
hendra, Meganingrum Arista Jiwangi, Adityo 673  
Anggraito, and Indra Budi. 2021. Multi-pass sieve 674  
coreference resolution for Indonesian. In *Proceed-* 675  
*ings of Recent Advances in Natural Language Pro-* 676  
*cessing*, pages 84–90. 677
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 678  
2020. On the cross-lingual transferability of mono- 679  
lingual representations. In *Proceedings of the 58th* 680  
*Annual Meeting of the Association for Computa-* 681  
*tional Linguistics*, pages 4623–4637. 682
- Jessica Naraiswari Arwidarasti, Ika Alfina, and 683  
Adila Alfa Krisnadhi. 2019. [Converting an Indone-](#) 684  
[sian constituency treebank to the Penn treebank for-](#) 685  
[mat](#). In *2019 International Conference on Asian* 686  
*Language Processing (IALP)*, pages 331–336. 687
- Annisa Nurul Azhar, Masayu Leylia Khodra, and 688  
Arie Pratama Sutiono. 2019. Multi-label aspect cat- 689  
egorization with convolutional neural networks and 690  
extreme gradient boosting. In *2019 International* 691  
*Conference on Electrical Engineering and Informat-* 692  
*ics (ICEEI)*, pages 35–40. IEEE. 693
- Kurniawati Azizah, Mirna Adriani, and Wisnu Jat- 694  
miko. 2020. [Hierarchical transfer learning for](#) 695  
[multilingual, multi-speaker, and style transfer dnn-](#) 696  
[based tts on low-resource languages](#). *IEEE Access*, 697  
8:179798–179812. 698
- Anab Maulana Barik, Rahmad Mahendra, and Mirna 699  
Adriani. 2019. Normalization of indonesian-english 700  
code-mixed twitter data. In *Proceedings of the 5th* 701  
*Workshop on Noisy User-generated Text (W-NUT* 702  
*2019)*, pages 417–424. 703
- Sadar Baskoro and Mirna Adriani. 2008. Developing 704  
an indonesian speech recognition system. In *Second* 705  
*MALINDO Workshop. Selangor, Malaysia*. 706
- Peter Bellwood. 1997. *Prehistory of the Indo-* 707  
*Malaysian Archipelago*. University of Hawaii Press, 708  
Honolulu. 709
- Peter Bellwood, Geoffrey Chambers, Malcolm Ross, 710  
and Hsiao-chun Hung. 2011. Are ‘cultures’ inher- 711  
ited? Multidisciplinary perspectives on the origins 712  
and migrations of Austronesian-speaking peoples 713  
prior to 1000 BC. In *Investigating archaeological* 714  
*cultures*, pages 321–354. Springer. 715
- Peter Bellwood and Eusebio Dizon. 2008. Austrone- 716  
sian cultural origins : out of Taiwan, via the Batanes 717  
Islands, and onwards to western Polynesia. In Ali- 718  
cia Sanchez-Mazas, editor, *Past human migrations* 719  
*in East Asia: matching archaeology, linguistics and* 720  
*genetics*. Routledge, London. 721
- Emily M Bender and Batya Friedman. 2018. Data 722  
statements for natural language processing: Toward 723  
mitigating system bias and enabling better science. 724  
*Transactions of the Association for Computational* 725  
*Linguistics*, 6:587–604. 726

727	Jacques Bertrand. 2004. <i>Nationalism and ethnic conflict in Indonesia</i> . Cambridge University Press.	<i>Human Language Technologies</i> , pages 1897–1907, Online. Association for Computational Linguistics.	782
728			783
729	Laurent Besacier, Etienne Barnard, Alexey Karpov, and Tanja Schultz. 2014. Automatic speech recognition for under-resourced languages: A survey. <i>Speech communication</i> , 56:85–100.	Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. <i>TyDi QA: A benchmark for information-seeking question answering in typologically diverse languages</i> . <i>Transactions of the Association for Computational Linguistics</i> , 8:454–470.	784
730			785
731			786
732			787
733	Steven Bird. 2020. <i>Decolonising speech and language technology</i> . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 3504–3519, Barcelona, Spain (Online). International Committee on Computational Linguistics.		788
734			789
735			790
736		Abigail C Cohn and Maya Ravindranath. 2014. Local languages in indonesia: Language maintenance or language shift. <i>Linguistik Indonesia</i> , 32(2):131–148.	791
737			792
738	Robert Blust. 1980. Austronesian etymologies. <i>Oceanic Linguistics</i> , 19:1–181.		793
739			794
740	Francis Bond and Timothy Baldwin. 2016. <i>Introduction to Japanese Computational Linguistics</i> .	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. <i>Unsupervised cross-lingual representation learning at scale</i> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451, Online. Association for Computational Linguistics.	795
741			796
742	Luiz Henrique Bonifacio, Israel Campiotti, Roberto Lotufo, and Rodrigo Nogueira. 2021. <i>mmarco: A multilingual version of ms marco passage ranking dataset</i> .		797
743			798
744			799
745			800
746	Indra Budi, Stéphane Bressan, Gatot Wahyudi, Zainal A. Hasibuan, and Bobby A. A. Nazief. 2005. Named entity recognition for the Indonesian language: Combining contextual, morphological and part-of-speech features into a knowledge engineering approach. In <i>Discovery Science</i> , pages 57–69, Berlin, Heidelberg. Springer Berlin Heidelberg.		801
747			802
748		Alexander R. Coupe and František Kratochvíl. 2020. <i>Asia before English</i> , pages 15–48. John Wiley & Sons, Inc.	804
749			805
750			806
751		Wenliang Dai, Samuel Cahyawijaya, Zihan Liu, and Pascale Fung. 2021. Multimodal end-to-end sparse model for emotion recognition. In <i>NAACL</i> .	807
752			808
753	Samuel Cahyawijaya. 2021. <i>Greenformers: Improving computation and memory efficiency in transformer models via low-rank approximation</i> .		809
754			
755		Wenliang Dai, Zihan Liu, Tiezheng Yu, and Pascale Fung. 2020. <i>Modality-transferable emotion embeddings for low-resource multimodal emotion recognition</i> . In <i>Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing</i> , pages 269–280, Suzhou, China. Association for Computational Linguistics.	810
756	Samuel Cahyawijaya, Genta Indra Winata, Bryan Wilie, Karissa Vincentio, Xiaohong Li, Adhiguna Kuncoro, Sebastian Ruder, Zhi Yuan Lim, Syafri Bahar, Masayu Leylia Khodra, et al. 2021. <i>Indonlg: Benchmark and resources for evaluating Indonesian natural language generation</i> . <i>arXiv preprint arXiv:2104.08200</i> .		811
757			812
758			813
759			814
760			815
761			816
762			817
763	Isaac Caswell, Theresa Breiner, Daan van Esch, and Ankur Bapna. 2020. <i>Language ID in the wild: Unexpected challenges on the path to a thousand-language web text corpus</i> . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 6588–6608, Barcelona, Spain (Online). International Committee on Computational Linguistics.		818
764			
765		Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. <i>BERT: Pre-training of deep bidirectional transformers for language understanding</i> . In <i>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)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	819
766			820
767			821
768			822
769			823
770			824
771	Isaac Caswell, Julia Kreutzer, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, et al. 2021. <i>Quality at a glance: An audit of web-crawled multilingual datasets</i> . <i>arXiv preprint arXiv:2103.12028</i> .		825
772			826
773			827
774		Arawinda Dinakaramani, Fam Rashel, Andry Luthfi, and Ruli Manurung. 2014. <i>Designing an Indonesian part of speech tagset and manually tagged Indonesian corpus</i> . In <i>2014 International Conference on Asian Language Processing (IALP)</i> , pages 66–69.	828
775			829
776			830
777	Yu-An Chung, Chenguang Zhu, and Michael Zeng. 2021. <i>SPLAT: Speech-language joint pre-training for spoken language understanding</i> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics:</i>		831
778			832
779		David M. Eberhard, Gary F. Simons, and Charles D. Fennig. 2021. <i>Ethnologue: Languages of the World. Twenty-fourth edition</i> . Dallas, Texas: SIL International.	833
780			834
781			835
			836

837	European Language Resources Association. 2019.	<i>the 13th Workshop on Building and Using Comparable Corpora</i> , pages 35–43, Marseille, France. European Language Resources Association.	892
838	BLT4All: Language Technologies for All. <a href="https://lt4all.elra.info/en/">https://lt4all.elra.info/en/</a> . [Online; accessed Dec. 2019.].		893
839			894
840	Muhammad Fachri. 2014. Named entity recognition for Indonesian text using hidden markov model. Undergraduate Thesis.	Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8342–8360.	895
841			896
842			897
843	Andri Imam Fauzi and Dwi Puspitorini. 2018. Dialect and identity: A case study of javanese use in whatsapp and line. In <i>IOP Conference Series: Earth and Environmental Science</i> , volume 175, page 012111. IOP Publishing.		898
844			899
845			900
846			901
847			902
848	Abdurrisyad Fikri and Ayu Purwarianti. 2012. Case based Indonesian closed domain question answering system with real world questions. In <i>2012 7th International Conference on Telecommunication Systems, Services, and Applications (TSSA)</i> , pages 181–186.	Akhmad Haryono. 2012. <i>Perubahan dan perkembangan bahasa: Tinjauan historis dan sosiolinguistik</i> . Ph.D. thesis, Udayana University.	903
849			904
850		Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. <i>XL-sum: Large-scale multilingual abstractive summarization for 44 languages</i> . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4693–4703, Online. Association for Computational Linguistics.	905
851			906
852			907
853	Philip Gage. 1994. A new algorithm for data compression. <i>C Users J.</i> , 12(2):23–38.		908
854			909
855	Dan Gillick, Cliff Brunk, Oriol Vinyals, and Amarnag Subramanya. 2016. <i>Multilingual language processing from bytes</i> . In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1296–1306, San Diego, California. Association for Computational Linguistics.		910
856			911
857			912
858			913
859			914
860			915
861			916
862			917
863	Dirk Goldhahn, Thomas Eckart, and Uwe Quasthoff. 2012. <i>Building large monolingual dictionaries at the Leipzig corpora collection: From 100 to 200 languages</i> . In <i>Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)</i> , pages 759–765, Istanbul, Turkey. European Language Resources Association (ELRA).	Devin Hoesen and Ayu Purwarianti. 2018. Investigating bi-lstm and crf with pos tag embedding for Indonesian named entity tagger. In <i>2018 International Conference on Asian Language Processing (IALP)</i> , pages 35–38. IEEE.	918
864			919
865			920
866			921
867			922
868			923
869			924
870	Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2021. The flores-101 evaluation benchmark for low-resource and multilingual machine translation.	Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. <i>XTREME: A Massively Multilingual Multi-task Benchmark for Evaluating Cross-lingual Generalization</i> . In <i>Proceedings of ICML 2020</i> .	925
871			926
872			927
873			928
874			929
875			930
876	Russell D. Gray and Fiona M. Jordan. 2000. Language trees support the express-train sequence of Austronesian expansion. <i>Nature</i> , 405:1052–1055.	Muhammad Okky Ibrohim and Indra Budi. 2019. <i>Multi-label hate speech and abusive language detection in Indonesian Twitter</i> . In <i>Proceedings of the Third Workshop on Abusive Language Online</i> , pages 46–57, Florence, Italy. Association for Computational Linguistics.	931
877			932
878			933
879	Yohanes Gultom and Wahyu Catur Wibowo. 2017. Automatic open domain information extraction from Indonesian text. In <i>2017 International Workshop on Big Data and Information Security (IW BIS)</i> , pages 23–30. IEEE.	Arfinda Ilmania, Abdurrahman, Samuel Cahyawijaya, and Ayu Purwarianti. 2018. <i>Aspect detection and sentiment classification using deep neural network for Indonesian aspect-based sentiment analysis</i> . In <i>2018 International Conference on Asian Language Processing (IALP)</i> , pages 62–67.	934
880			935
881			936
882			937
883			938
884	William Gunawan, Derwin Suhartono, Fredy Purnomo, and Andrew Ongko. 2018. Named-entity recognition for indonesian language using bidirectional lstm-cnns. <i>Procedia Computer Science</i> , 135:425–432.	Benediktus Sridin Sulu Jahang and Zita Meirina. 2021. <i>1,3 juta anak di ntt belum bisa berbahasa indonesia</i> . Last accessed on 05/10/2021.	939
885			940
886			941
887			942
888			943
889	Tri Wahyu Guntara, Alham Fikri Aji, and Radityo Eko Prasojo. 2020. <i>Benchmarking multidomain English-Indonesian machine translation</i> . In <i>Proceedings of</i>	Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. <i>The State and Fate of Linguistic Diversity and Inclusion in the NLP World</i> . In <i>Proceedings of ACL 2020</i> .	944
890			945
891			946
		Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. Bag of tricks for efficient text classification. In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers</i> , pages 427–431. Association for Computational Linguistics.	947
			948

949	Erlin Kartikasari, Kisyani Laksono, Dian Savitri	<i>Meeting of the Association for Computational Lin-</i>	1006
950	Agusniar, and Diah Yovita Suryarini. 2018. A	<i>guistics (Volume 1: Long Papers)</i> , pages 66–75, Mel-	1007
951	study of dialectology on Javanese "Ngoko" in	bourne, Australia. Association for Computational	1008
952	Banyuwangi, Surabaya, Magetan, and Solo. <i>Human-</i>	Linguistics.	1009
953	<i>iora</i> , 30(2):128.		
954	Simran Khanuja, Diksha Bansal, Sarvesh Mehtani,	Kemal Kurniawan and Alham Fikri Aji. 2018. Toward	1010
955	Savya Khosla, Atreyee Dey, Balaji Gopalan,	a standardized and more accurate Indonesian part-of-	1011
956	Dilip Kumar Margam, Pooja Aggarwal, Rajiv Teja	speech tagging. In <i>2018 International Conference</i>	1012
957	Nagipogu, Shachi Dave, Shruti Gupta, Subhash	<i>on Asian Language Processing (IALP)</i> , pages 303–	1013
958	Chandra Bose Gali, Vish Subramanian, and Partha	307. IEEE.	1014
959	Talukdar. 2021. <b>MuRIL: Multilingual Representa-</b>	Kemal Kurniawan and Samuel Louvan. 2018. Indo-	1015
960	<b>sum: A new benchmark dataset for Indonesian text</b>	summarization. In <i>2018 International Conference</i>	1016
961	<b>tations for Indian Languages.</b> <i>arXiv preprint</i>	<i>on Asian Language Processing (IALP)</i> , pages 215–	1017
	<i>arXiv:2103.10730.</i>	220. IEEE.	1018
962	Yash Khemchandani, Sarvesh Mehtani, Vaidehi Patil,		1019
963	Abhijeet Awasthi, Partha Talukdar, and Sunita	Teven Le Scao and Alexander M Rush. 2021. How	1020
964	Sarawagi. 2021. <b>Exploiting Language Relatedness</b>	many data points is a prompt worth? In <i>Proceedings</i>	1021
965	<b>for Low Web-Resource Language Model Adapta-</b>	<i>of the 2021 Conference of the North American Chap-</i>	1022
966	<b>tion: An Indic Languages Study.</b> In <i>Proceedings</i>	<i>ter of the Association for Computational Linguistics:</i>	1023
967	<i>of ACL 2021.</i>	<i>Human Language Technologies</i> , pages 2627–2636.	1024
968	Marian Klamer. 2018. Documenting the linguistic di-	Dessi Puji Lestari, Koji Iwano, and Sadaoki Furui.	1025
969	versity of indonesia: Time is running out. In <i>Pro-</i>	2006. A large vocabulary continuous speech recog-	1026
970	<i>ceedings of ‘Revitalization of local languages as the</i>	nition system for indonesian language. In <i>15th In-</i>	1027
971	<i>pillar of pluralism’</i> , pages 1–10. APBL (Asosiasi	<i>ndonesian Scientific Conference in Japan Proceed-</i>	1028
972	Peneliti Bahasa-bahasa Lokal) and Nusa Cendana	<i>ings</i> , pages 17–22.	1029
973	University, Kupang, Satya Wacana Press.		
974	Fajri Koto and Ikhwan Koto. 2020. <b>Towards compu-</b>	Zhaojiang Lin, Zihan Liu, Genta Indra Winata, Samuel	1030
975	<b>tational linguistics in Minangkabau language: Stud-</b>	Cahyawijaya, Andrea Madotto, Yejin Bang, Etsuko	1031
976	<b>ies on sentiment analysis and machine translation.</b>	Ishii, and Pascale Fung. 2021. <b>XPersona: Evaluat-</b>	1032
977	In <i>Proceedings of the 34th Pacific Asia Conference</i>	<b>ing multilingual personalized chatbot.</b> In <i>Proceed-</i>	1033
978	<i>on Language, Information and Computation</i> , pages	<i>ings of the 3rd Workshop on Natural Language Pro-</i>	1034
979	138–148, Hanoi, Vietnam. Association for Computa-	<i>cessing for Conversational AI</i> , pages 102–112, On-	1035
980	tional Linguistics.	line. Association for Computational Linguistics.	1036
981	Fajri Koto, Jey Han Lau, and Timothy Baldwin. 2020a.	Fangyu Liu, Emanuele Bugliarello, Edoardo Maria	1037
982	<b>Liputan6: A large-scale Indonesian dataset for text</b>	Ponti, Siva Reddy, Nigel Collier, and Desmond El-	1038
983	<b>summarization.</b> In <i>Proceedings of the 1st Confer-</i>	liott. 2021. <b>Visually grounded reasoning across lan-</b>	1039
984	<i>ence of the Asia-Pacific Chapter of the Association</i>	<b>guages and cultures.</b> In <i>Proceedings of the 2021</i>	1040
985	<i>for Computational Linguistics and the 10th Interna-</i>	<i>Conference on Empirical Methods in Natural Lan-</i>	1041
986	<i>tional Joint Conference on Natural Language Pro-</i>	<i>guage Processing</i> , pages 10467–10485, Online and	1042
987	<i>cessing</i> , pages 598–608, Suzhou, China. Associa-	Punta Cana, Dominican Republic. Association for	1043
988	tion for Computational Linguistics.	Computational Linguistics.	1044
989	Fajri Koto, Jey Han Lau, and Timothy Baldwin. 2021.	Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey	1045
990	<b>IndoBERTweet: A pretrained language model for In-</b>	Edunov, Marjan Ghazvininejad, Mike Lewis, and	1046
991	<b>donesian twitter with effective domain-specific vo-</b>	Luke Zettlemoyer. 2020. <b>Multilingual denoising</b>	1047
992	<b>cabulary initialization.</b> In <i>Proceedings of the 2021</i>	<b>pre-training for neural machine translation.</b> <i>Transac-</i>	1048
993	<i>Conference on Empirical Methods in Natural Lan-</i>	<i>tions of the Association for Computational Linguis-</i>	1049
994	<i>guage Processing.</i>	<i>tics</i> , 8:726–742.	1050
995	Fajri Koto, Afshin Rahimi, Jey Han Lau, and Timo-	Marco Lui and Timothy Baldwin. 2012. <b>langid.py:</b>	1051
996	thy Baldwin. 2020b. <b>IndoLEM and IndoBERT: A</b>	<b>An off-the-shelf language identification tool.</b> In <i>Pro-</i>	1052
997	<b>benchmark dataset and pre-trained language model</b>	<i>ceedings of the ACL 2012 system demonstrations</i> ,	1053
998	<b>for Indonesian NLP.</b> In <i>Proceedings of the 28th In-</i>	pages 25–30.	1054
999	<i>ternational Conference on Computational Linguis-</i>	Edwin Lunando and Ayu Purwarianti. 2013. <b>Indone-</b>	1055
1000	<i>tics</i> , pages 757–770, Barcelona, Spain (Online). In-	<b>isian social media sentiment analysis with sarcasm</b>	1056
1001	ternational Committee on Computational Linguis-	<b>detection.</b> In <i>2013 International Conference on Ad-</i>	1057
1002	tics.	<i>vanced Computer Science and Information Systems</i>	1058
1003	Taku Kudo. 2018. <b>Subword regularization: Improving</b>	<i>(ICACSYS)</i> , pages 195–198.	1059
1004	<b>neural network translation models with multiple sub-</b>		
1005	<b>word candidates.</b> In <i>Proceedings of the 56th Annual</i>		

1060	Andrea Madotto, Zhaojiang Lin, Genta Indra Winata, and Pascale Fung. 2021. Few-shot bot: Prompt-based learning for dialogue systems. <i>arXiv preprint arXiv:2110.08118</i> .	1117
1061		1118
1062		1119
1063		1120
1064	Putu Devi Maharani and Komang Dian Puspita Candra. 2018. Variasi leksikal bahasa Bali dialek kuta selatan. <i>Mudra Jurnal Seni Budaya</i> , 33(1):76–84.	1121
1065		1122
1066		1123
1067	Rahmad Mahendra, Alham Fikri Aji, Samuel Louvan, Fahrurrozi Rahman, and Clara Vania. 2021. IndoNLI: A natural language inference dataset for Indonesian. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	1124
1068		1125
1069		1126
1070		1127
1071		1128
1072		1129
1073	Rahmad Mahendra, Septina Dian Larasati, and Ruli Manurung. 2008. Extending an Indonesian semantic analysis-based question answering system with linguistic and world knowledge axioms. In <i>Proceedings of the 22nd Pacific Asia Conference on Language, Information and Computation</i> , pages 262–271, The University of the Philippines Visayas Cebu College, Cebu City, Philippines. De La Salle University, Manila, Philippines.	1130
1074		1131
1075		1132
1076		1133
1077		1134
1078		1135
1079		1136
1080		1137
1081		1138
1082	Rahmad Mahendra, Heninggar Septiantri, Haryo Akbarianto Wibowo, Ruli Manurung, and Mirna Adriani. 2018. Cross-lingual and supervised learning approach for Indonesian word sense disambiguation task. In <i>Proceedings of the 9th Global Wordnet Conference</i> , pages 245–250, Nanyang Technological University (NTU), Singapore. Global Wordnet Association.	1139
1083		1140
1084		1141
1085		1142
1086		1143
1087		1144
1088		1145
1089		1146
1090	Miftahul Mahfuzh, Sidik Soleman, and Ayu Purwarianti. 2019. Improving joint layer rnn based keyphrase extraction by using syntactical features. In <i>2019 International Conference of Advanced Informatics: Concepts, Theory and Applications (ICAICTA)</i> , pages 1–6. IEEE.	1147
1091		1148
1092		1149
1093		1150
1094		1151
1095		1152
1096	Ryan McDonald, Joakim Nivre, Yvonne Quirnbach-Brundage, Yoav Goldberg, Dipanjan Das, Kuzman Ganchev, Keith Hall, Slav Petrov, Hao Zhang, Oscar Täckström, Claudia Bedini, Núria Bertomeu Castelló, and Jungmee Lee. 2013. Universal Dependency annotation for multilingual parsing. In <i>Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 92–97, Sofia, Bulgaria. Association for Computational Linguistics.	1153
1097		1154
1098		1155
1099		1156
1100		1157
1101		1158
1102		1159
1103		1160
1104		1161
1105		1162
1106	Aqsath Rasyid Naradhipa and Ayu Purwarianti. 2011. Sentiment classification for Indonesian message in social media. In <i>Proceedings of the 2011 International Conference on Electrical Engineering and Informatics</i> , pages 1–4.	1163
1107		1164
1108		1165
1109		1166
1110		1167
1111	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 Elsahar, 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 <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2144–2160, Online. Association for Computational Linguistics.	1168
1112		1169
1113		1170
1114		1171
1115		1172
1116		1173
		1174
	Toan Q Nguyen and David Chiang. 2017. Transfer learning across low-resource, related languages for neural machine translation. In <i>Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 296–301.	1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
		1311
		1312
		1313
		1314
		1315
		1316
		1317
		1318
		1319
		1320
		1321
		1322
		1323
		1324
		1325
		1326
		1327
		1328
		1329
		1330
		1331
		1332
		1333
		1334
		1335
		1336
		1337
		1338
		1339
		1340
		1341
		1342
		1343
		1344
		1345
		1346
		1347
		1348
		1349
		1350
		1351
		1352
		1353
		1354
		1355
		1356
		1357
		1358
		1359
		1360
		1361
		1362
		1363
		1364
		1365
		1366
		1367
		1368
		1369
		1370
		1371
		1372
		1373
		1374
		1375
		1376
		1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400

1175	Gurevych. 2021. xgqa: Cross-lingual visual question answering. <i>ArXiv</i> , abs/2109.06082.	
1176		
1177	Femphy Pisceldo, Rahmad Mahendra, Ruli Manurung, and I Wayan Arka. 2008. A two-level morphological analyser for the Indonesian language. In <i>Proceedings of the Australasian Language Technology Association Workshop 2008</i> , pages 142–150.	
1178		
1179		
1180		
1181		
1182	Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. <a href="#">XCOPA: A multilingual dataset for causal commonsense reasoning</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2362–2376, Online. Association for Computational Linguistics.	
1183		
1184		
1185		
1186		
1187		
1188		
1189	Shana Poplack. 1980. Sometimes i’ll start a sentence in Spanish y termino en español: toward a typology of code-switching I.	
1190		
1191		
1192	Apriyani Purwaningsih. 2017. Geografi dialek bahasa Jawa pesisiran di desa paciran kabupaten Lamongan. In <i>Proceeding of International Conference on Art, Language, and Culture</i> , pages 594–605.	
1193		
1194		
1195		
1196	Ayu Purwarianti and Ida Ayu Putu Ari Crisdayanti. 2019. Improving bi-lstm performance for Indonesian sentiment analysis using paragraph vector. In <i>2019 International Conference of Advanced Informatics: Concepts, Theory and Applications (ICAICTA)</i> , pages 1–5. IEEE.	
1197		
1198		
1199		
1200		
1201		
1202	Ayu Purwarianti, Masatoshi Tsuchiya, and Seiichi Nakagawa. 2007. A machine learning approach for Indonesian question answering system. In <i>Artificial Intelligence and Applications</i> , pages 573–578.	
1203		
1204		
1205		
1206	Shofianina Dwi Ananda Putri, Muhammad Okky Ibrahim, and Indra Budi. 2021. Abusive language and hate speech detection for Indonesian-local language in social media text. In <i>Recent Advances in Information and Communication Technology 2021</i> , pages 88–98, Cham. Springer International Publishing.	
1207		
1208		
1209		
1210		
1211		
1212	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	
1213		
1214		
1215	Anna Rogers, Tim Baldwin, and Kobi Leins. 2021. Just what do you think you’re doing, dave? a checklist for responsible data use in nlp. <i>arXiv preprint arXiv:2109.06598</i> .	
1216		
1217		
1218		
1219	Malcolm Ross. 2005. Pronouns as a preliminary diagnostic for grouping Papuan languages. <i>Papuan pasts: Cultural, linguistic and biological histories of Papuan-speaking peoples</i> , 572:15–65.	
1220		
1221		
1222		
1223	Nur Endah Safitri, Amalia Zahra, and Mirna Adriani. 2016. Spoken language identification with phonotactics methods on minangkabau, Sundanese, and Javanese languages. <i>Procedia Computer Science</i> , 81:182–187.	
1224		
1225		
1226		
1227		
	Mei Silviana Saputri, Rahmad Mahendra, and Mirna Adriani. 2018. Emotion classification on Indonesian twitter dataset. In <i>2018 International Conference on Asian Language Processing (IALP)</i> , pages 90–95. IEEE.	1228 1229 1230 1231 1232
	Gita Sarwadi, Mahsun Mahsun, and Burhanuddin Burhanuddin. 2019. Lexical variation of Sasak Kuto-Kute dialect in North Lombok district. <i>Jurnal Kata: Penelitian tentang Ilmu Bahasa dan Sastra</i> , 3(1):155–169.	1233 1234 1235 1236 1237
	Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. 2020. Green ai. <i>Communications of the ACM</i> , 63(12):54–63.	1238 1239 1240
	Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021. <a href="#">WikiMatrix: Mining 135M parallel sentences in 1620 language pairs from Wikipedia</a> . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 1351–1361, Online. Association for Computational Linguistics.	1241 1242 1243 1244 1245 1246 1247 1248
	Dmitriy Serdyuk, Yongqiang Wang, Christian Fuegen, Anuj Kumar, Baiyang Liu, and Yoshua Bengio. 2018. <a href="#">Towards end-to-end spoken language understanding</a> . IEEE.	1249 1250 1251 1252
	Ken Nabila Setya and Rahmad Mahendra. 2018. <a href="#">Semi-supervised textual entailment on Indonesian wikipedia data</a> . In <i>2018 International Conference on Computational Linguistics and Intelligent Text Processing (CICLing)</i> .	1253 1254 1255 1256 1257
	Masitowarni Siregar, Syamsul Bahri, Dedi Sanjaya, et al. 2014. Code switching and code mixing in Indonesia: Study in sociolinguistics. <i>English Language and Literature Studies</i> , 4(1):77–92.	1258 1259 1260 1261
	James Neil Sneddon. 2003. <i>The Indonesian language: Its history and role in modern society</i> . UNSW Press, Sydney.	1262 1263 1264
	Soeparno. 2015. <a href="#">Kerancuan fono-ortografis dan ortofologis bahasa Indonesia ragam lisan dan tulis</a> . <i>Diksi</i> , 12(2).	1265 1266 1267
	Hein Steinhauer. 2005. Colonial history and language policy in insular Southeast Asia and Madagascar. In Alexander Adelaar and Nikolaus P. Himmelmann, editors, <i>The Austronesian Languages of Asia and Madagascar</i> , chapter 3, pages 65–86. Routledge, Oxon.	1268 1269 1270 1271 1272 1273
	Arie Ardiyanti Suryani, Dwi Hendratmo Widyantoro, Ayu Purwarianti, and Yayat Sudaryat. 2015. <a href="#">Experiment on a phrase-based statistical machine translation using pos tag information for Sundanese into Indonesian</a> . In <i>2015 International Conference on Information Technology Systems and Innovation (IC-ITSI)</i> , pages 1–6.	1274 1275 1276 1277 1278 1279 1280

1281	Arie Ardiyanti Suryani, Dwi Hendratmo Widyantoro, Ayu Purwarianti, and Yayat Sudaryat. 2018. <a href="#">The rule-based sundanese stemmer</a> . <i>ACM Trans. Asian Low-Resour. Lang. Inf. Process.</i> , 17(4).	1336
1282		1337
1283		1338
1284		
1285	Taufic Leonardo Sutejo and Dessi Puji Lestari. 2018. Indonesia hate speech detection using deep learning. In <i>2018 International Conference on Asian Language Processing (IALP)</i> , pages 39–43. IEEE.	1339
1286		1340
1287		1341
1288		1342
1289	Bejo Sutrisno and Yessika Ariesta. 2019. Beyond the use of code mixing by social media influencers in instagram. <i>Advances in Language and Literary Studies</i> , 10(6):143–151.	1343
1290		1344
1291		1345
1292		
1293	Yi Tay, Vinh Q Tran, Sebastian Ruder, Jai Gupta, Hyung Won Chung, Dara Bahri, Zhen Qin, Simon Baumgartner, Cong Yu, and Donald Metzler. 2021. Charformer: Fast character transformers via gradient-based subword tokenization. <i>arXiv preprint arXiv:2106.12672</i> .	1346
1294		1347
1295		1348
1296		1349
1297		1350
1298		1351
1299	Jozina Vander Klok. 2015. The dichotomy of auxiliaries in javanese: Evidence from two dialects. <i>Australian Journal of Linguistics</i> , 35(2):142–167.	1352
1300		1353
1301		1354
1302	Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. <i>Communications of the ACM</i> , 57(10):78–85.	1355
1303		1356
1304		1361
1305	Haryo Akbarianto Wibowo, Made Nindyatama Nityasya, Afra Feyza Akyürek, Suci Fitriany, Alham Fikri Aji, Radityo Eko Prasajo, and Derry Tanti Wijaya. 2021. <a href="#">IndoCollex: A testbed for morphological transformation of Indonesian word colloquialism</a> . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3170–3183, Online. Association for Computational Linguistics.	1362
1306		1363
1307		1364
1308		1365
1309		1366
1310		
1311		1367
1312		1368
1313		1369
1314	Haryo Akbarianto Wibowo, Tatag Aziz Prawiro, Muhammad Ihsan, Alham Fikri Aji, Radityo Eko Prasajo, Rahmad Mahendra, and Suci Fitriany. 2020. Semi-supervised low-resource style transfer of Indonesian informal to formal language with iterative forward-translation. In <i>2020 International Conference on Asian Language Processing (IALP)</i> , pages 310–315. IEEE.	1370
1315		1371
1316		1372
1317		1373
1318		1374
1319		1375
1320		
1321		
1322	Alfan Farizki Wicaksono and Ayu Purwarianti. 2010. Hmm based part-of-speech tagger for bahasa indonesia. In <i>4th International MALINDO (Malaysian-Indonesian Language) Workshop</i> .	1376
1323		1377
1324		1378
1325		1379
1326	Bryan Wilie, Karissa Vincentio, Genta Indra Winata, Samuel Cahyawijaya, Xiaohong Li, Zhi Yuan Lim, Sidik Soleman, Rahmad Mahendra, Pascale Fung, Syafri Bahar, et al. 2020. Indonlu: Benchmark and resources for evaluating Indonesian natural language understanding. In <i>Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing</i> , pages 843–857.	1380
1327		1381
1328		1382
1329		1383
1330		1384
1331		1385
1332		
1333		1386
1334		1387
1335		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
		1406
		1407
		1408
		1409
		1410
		1411
		1412
		1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
		1421
		1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
		1431
		1432
		1433
		1434
		1435
		1436
		1437
		1438
		1439
		1440
		1441
		1442
		1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
		1458
		1459
		1460
		1461
		1462
		1463
		1464
		1465
		1466
		1467
		1468
		1469
		1470
		1471
		1472
		1473
		1474
		1475
		1476
		1477
		1478
		1479
		1480
		1481
		1482
		1483
		1484
		1485
		1486
		1487
		1488
		1489
		1490
		1491
		1492
		1493
		1494
		1495
		1496
		1497
		1498
		1499
		1500

Task: *Multilingual Parsing from Raw Text to Universal Dependencies*, pages 1–21, Brussels, Belgium. Association for Computational Linguistics.

## A Language Statistics

In Figure 4, we contrast the distribution of publication with Indonesian language compared to European languages. Despite number of Indonesian speakers is much larger compared to some European languages, number of published research in Indonesian is still comparatively lower.

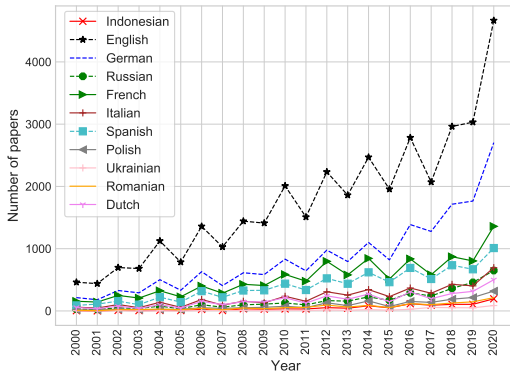


Figure 4: Number of published research works per year in Indonesian and European languages from 2000 to 2020.

## B Wikipedia Vocabulary Overlap

Lang	# Vocab		
	All (k)	Top 1% $\cap$ KBBI (%)	Top 100 $\cap$ KBBI (%)
id	2023	59.3	<b>96</b>
jav	435	46.8	43
su	286	44.3	47
min	252	30.3	41
bug	23	35.7	27
map-bms	14	76.7	79
gor	12	40.5	49
ace	12	37.6	46
ban	10	43.3	46
bjn	4	62.9	69
nia	1	25.9	30
mad	1	26.9	24

Table 7: Vocabulary of Indonesian languages in Wikipedia, filtered with KBBI third edition <sup>12</sup>

In Table 7, we present vocabulary statistics of Indonesian languages in Wikipedia. Due to the noisy nature of Wikipedia, we use “*Kamus Besar Bahasa Indonesia*” (KBBI) third edition,<sup>13</sup> the official dictionary for the Indonesian language to filter the top 1% and top-100 most frequent words. As

<sup>12</sup>KBBI is the official Indonesian dictionary.

<sup>13</sup><https://github.com/geovedi/indonesian-wordlist>

expected, the top 1% words are less reliable, with only 59.3% of the vocabulary overlap between id and KBBI. In the top-100 words, there is a 96% word overlap with KBBI, making this set more reliable. Previous work on Minangkabau by Koto and Koto (2020) also showed that id-min words have a 55% overlap in a manually curated bilingual dictionary, closer to the top-100 value for min in Table 7.

## C Indonesian NLP Resources

On Table 8, we list statistics of Indonesian language corpora for different tasks, including sentiment analysis, part-of-speech tagging, summarization, NLI, and discourse. Although most datasets are in Indonesian and only a few are in Minangkabau (Koto and Koto, 2020), Javanese and Sundanese (Cahyawijaya et al., 2021), these resource collections are arguably beneficial for constructing resources in other local languages. This is because 1) Indonesian can be used as a pivot language with regard to local languages due to the large vocabulary overlap (see Table 7), and 2) most Indonesians are bilingual, speaking both Indonesian and their local language (Koto and Koto, 2020).

## D Wikipedia Availability

In Figure 5 we compare Wikipedia size (in GB file size) compared to the number of speakers across various languages. We show that some African languages are even more under-resourced.

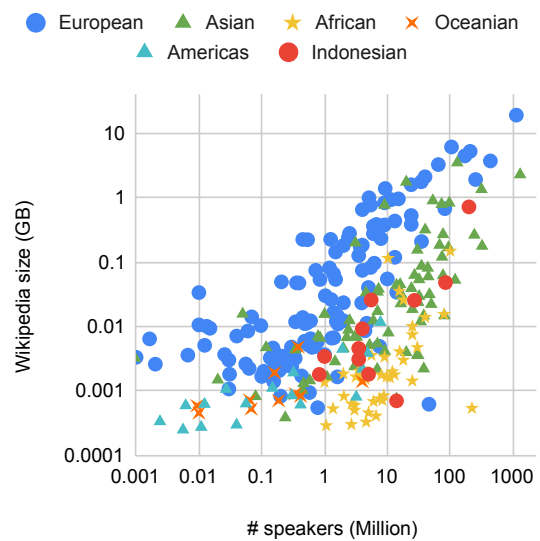


Figure 5: Wikipedia data size (in GB) compared to the number of speakers.



Name	Task Type	Size	% Indo	% Local
Labelled Datasets				
POSP (Hoesen and Purwarianti, 2018)	PoS Tagging	8k	100%	0%
BaPOS (Dinakaramani et al., 2014)	PoS Tagging	10k	100%	0%
NERGrit (Wilie et al., 2020)	Named Entity Recognition	2k	100%	0%
NERP (Hoesen and Purwarianti, 2018)	Named Entity Recognition	8k	100%	0%
(Gultom and Wibowo, 2017)	Named Entity Recognition	2k	100%	0%
Singgalang (Alfina et al., 2016)	Named Entity Recognition	48K	100%	0%
(Fachri, 2014)	Named Entity Recognition	2K	100%	0%
KEPS (Mahfuzh et al., 2019)	Keyphrase Extraction	1k	100%	0%
FacQA (Purwarianti et al., 2007)	Question Answering	3k	100%	0%
WReTE (Setya and Mahendra, 2018)	Natural Language Inference	0.5k	100%	0%
IndoNLI (Mahendra et al., 2021)	Natural Language Inference	18k	100%	0%
CASA (Ilmania et al., 2018)	Sentiment Analysis	1k	100%	0%
SmSA (Purwarianti and Crisdayanti, 2019)	Sentiment Analysis	13k	100%	0%
HoASA (Azhar et al., 2019)	Sentiment Analysis	3k	100%	0%
SA in IndoLEM (Koto et al., 2020b)	Sentiment Analysis	5k	100%	0%
SA in MinangNLP (Koto and Koto, 2020)	Sentiment Analysis	5K	0%	100%
(Saputri et al., 2018)	Emotion Classification	4k	100%	0%
(Ibrohim and Budi, 2019)	Hate Speech Detection	13k	100%	0%
TED En-Id (Guntara et al., 2020)	Machine Translation	93k	100%	0%
News En-Id (Guntara et al., 2020)	Machine Translation	42k	100%	0%
Religion En-Id (Guntara et al., 2020)	Machine Translation	590k	100%	0%
MT in MinangNLP (Koto and Koto, 2020)	Machine Translation	11K	0%	100%
EN↔ID MT (Cahyawijaya et al., 2021)	Machine Translation	31K	100%	0%
SU↔ID MT (Cahyawijaya et al., 2021)	Machine Translation	16K	0%	100%
JV↔ID MT (Cahyawijaya et al., 2021)	Machine Translation	16K	0%	100%
IndoSum (Kurniawan and Louvan, 2018)	Summarization	20k	100%	0%
Liputan6 (Koto et al., 2020a)	Summarization	215k	100%	0%
Kethu (Arwidarasti et al., 2019)	Constituency Parsing	1k	100%	0%
UD-Id GSD (McDonald et al., 2013)	Dependency Parsing	5k	100%	0%
UD-Id PUD (Zeman et al., 2018)	Dependency Parsing	1k	100%	0%
(Mahendra et al., 2018)	Word Sense Disambiguation	2k	100%	0%
IndoCoref (Artari et al., 2021)	Coreference Resolution	0.2k	100%	0%
NTP and Tweet Ordering (Koto et al., 2020b)	Discourse	7k	100%	0%
Pretraining Corpora				
Indo4B (Wilie et al., 2020)	-	3.6B words	100%	0%
Indo4B-Plus (Cahyawijaya et al., 2021)	-	4.0B words	89.64%	10.36%

Table 8: Statistics of publicly available datasets, most datasets are covered on the existing Indonesian languages NLP benchmarks (Wilie et al., 2020; Koto et al., 2020b; Cahyawijaya et al., 2021).

## E Dialect Differences

In this section, we present more examples of lexical variation of other local languages. Maharani and Candra (2018) and Sarwadi et al. (2019) show lexical variation of Balinese and Sasak, respectively, where they ask locals to translate general/common words. Then, they compare the vocabulary across different locations (in this case, villages) to each other. Some of the examples can be seen in Table 9 and 10. Unfortunately, they did not provide quantitative results. Pamolango (2012) conducted a similar experiment in the Banggai district in South Sulawesi across 31 observation points for the Saluan language. While Pamolango (2012) did not provide full examples, they reported up to 23.5% lexical variation among 200 basic vocabulary items.

## F Local Language Classification

As shown in Table 11, some of the Javanese texts are misidentified as Indonesian, English, and Malaysian. Javanese and Indonesian (which is similar to Malaysian) share some words. We believe English mis-classification is due to the data size bias.

English	Kedonganan	Jimbaran	Unggasan
I/me	Tyang	Tyang	Aku
You	Béné	Béné	Éngko
Umbrella	Pajéng	Pajéng	Pajong
Hat	Capil	Topong	Cecapil, Tetopong
How	Engken	Engken	Kengen
Where	Dijé	Dijé	Di joho
All	Konyangan	Onyé	Konyangan, onyang
Swallow (vb)	Gélék, ngélék	Gélék, ngélék	Ngélokang
Scratch (vb)	Gagas	Gagas	Gauk
Cough (vb)	Kokoan	Dékah	Kohkohan
Dawn	Plimunan	Plimunan	Sémongan
Afternoon	Sanjé	Sanjé	Sanjano

Table 9: Lexical variation of Balinese across different villages in South Kuta district, Bali (Maharani and Candra, 2018)

English	Pemenang Timur	Jenggala	Genggelang	Kayangan	Akar-Akar
Here	Ite	ite	ite	ite	tinI
There	Ito	ito	ito	ito	tinO
You	di?	sita	di?	sita	di?
Husband	kurənan	sawa	sawa	sawa	sawa
No	de?	de?	de?	de?	sora?
Paddle	bose	bose	dayung	dayung	bose
Spear	tər	cinəkan	tər	tombak	tombak
Black	birəŋ	birəŋ	birəŋ	birəŋ	pisak
Red	bəŋəŋ	bəŋəŋ	bəŋəŋ	bəŋəŋ	abaŋ
White	putz?	putz?	putz?	putz?	pətak
Worm	gumbər	loŋa	gumbər	gumbər	gumbər

Table 10: Lexical variation of Sasak across different villages in North Lombok district (Sarwadi et al., 2019)

Dialect/ Style	Method	classified as			
		jv	id	en	ms
Western- <i>Ngoko</i>	Langid	0.241	0.103	0.172	0.069
	FastText	0.069	0.276	0.276	0.069
	CLD3	0.759	0.000	0.000	0.034
Central- <i>Ngoko</i>	Langid	0.345	0.138	0.069	0.069
	FastText	0.379	0.310	0.069	0.069
	CLD3	0.828	0.000	0.000	0.034
Eastern- <i>Ngoko</i>	Langid	0.276	0.103	0.069	0.138
	FastText	0.103	0.310	0.103	0.034
	CLD3	0.552	0.103	0.000	0.000
Eastern- <i>Krama</i>	Langid	0.345	0.241	0.034	0.172
	FastText	0.379	0.310	0.069	0.034
	CLD3	0.897	0.000	0.000	0.000

Table 11: Language identification mis-classification rate.