

Part-of-speech Tagset and Corpus Development for Igbo, an African Language

Ikechukwu E. Onyenwe Dept. of Computer Science, University of Sheffield Sheffield S1 4DP, UK i.onyenwe@shef.ac.uk	Dr. Chinedu Uchechukwu Dept. of Linguistics Nnamdi Azikiwe University Anambra State, Nigeria neduchi@yahoo.com	Dr. Mark Hepple Dept. of Computer Science, University of Sheffield Sheffield S1 4DP, UK m.r.hepple@shef.ac.uk
--	---	--

Abstract

This project aims to develop linguistic resources to support computational NLP research on the Igbo language. The starting point for this project is the development of a new part-of-speech tagging scheme based on the EAGLES tagset guidelines, adapted to incorporate additional language internal features. The tags are currently being used in a part-of-speech annotation task for the development of POS tagged Igbo corpus. The proposed tagset has 59 tags.

1 Introduction

Supervised machine learning methods in NLP require an adequate amount of training data. The first crucial step for a part-of-speech (POS) tagging system for a language is a well designed, consistent, and complete tagset (Bamba Dione et al., 2010) which must be preceded by a detailed study and analysis of the language. Our tagset was developed from scratch through the study of linguistics and electronic texts in Igbo, using the EAGLES recommendations.

This initial manual annotation is important. Firstly, information dealing with challenging phenomena in a language is expressed in the tagging guideline; secondly, computational POS taggers require annotated text as training data. Even in unsupervised methods, some annotated texts are still required as a benchmark in evaluation. With this in mind, our tagset design follows three main goals: to determine the tagset size, since a smaller granularity provides higher accuracy and less ambiguity (de Pauwy et al., 2012); to use a sizeable scheme to capture the grammatical distinctions at a word level suited for further grammatical analysis, such as parsing; and to deliver good accuracy for automatic tagging, using the manually tagged data. We discuss the development of the tagset and corpus for Igbo. This work is, to the best of our knowledge, the first published work attempting to develop statistical NLP resources for Igbo.

2 Some Grammatical Features of the Igbo Language

2.1 Language family and speakers

The Igbo language has been classified as a Benue-Congo language of the Kwa sub-group of the Niger-Congo family¹ and is one of the three major languages in Nigeria, spoken in the eastern part of Nigeria, with about 36 million speakers². Nigeria is a multilingual country having around 510 living languages¹, but English serves as the official language.

2.2 Phonology

Standard Igbo has eight vowels and thirty consonants. The 8 vowels are divided into two harmony groups that are distinguished on the basis of the Advanced Tongue Root (ATR) phenomenon. They are -ATR: i [ɪ], ɯ [ʊ], a [ɑ], ɔ [ɔ] and +ATR: i [i], u [u], e [e], o [o] (Uchechukwu, 2008). Many Igbo words select their vowels from the same harmony group. Also, Igbo is a tonal language. There are three distinct tones

This work is licensed under a Creative Commons Attribution 4.0 International Licence. Page numbers and proceedings footer are added by the organisers. Licence details: <http://creativecommons.org/licenses/by/4.0/>

¹<http://nigerianwiki.com/wiki/Languages>

²http://en.wikipedia.org/wiki/Igbo_people

recognized in the language viz; High, Low, and Downstep. The tones are represented as *High* [H] = [´], *Low* [L] = [˘], *downstep* = [˘˘] (Emenanjo, 1978; Ikekeonwu, 1999) and are placed above the tone bearing units (TBU) of the language.

There are two tone marking systems, either: all high tones are left unmarked and all low tones and downsteps are marked (Green and Igwe, 1963; Emenanjo, 1978), or only contrastive tones are marked (Welmers and Welmers, 1968; Nwachukwu, 1995). We used the first system to illustrate the importance of tonal feature in the language’s lexical or grammatical structure. For example, at the lexical level the word *akwa* without a tone mark can be given the equivalent of ‘bed/bridge’, ‘cry’, ‘cloth’, or ‘egg’. But these equivalents can be properly distinguished when tone marked, as follows: *akwa* “cry”, *akwà* “cloth”, *àkwà* “bed or brigde”, *àkwa* “egg”. At the grammatical level, an interrogative sentence can be distinguished from a declarative sentence through a change in tone of the person pronouns from a high tone (e.g. *Ọ nà-àbịa* “He is coming”) to a low tone (e.g. *Ò nà-àbịa* “Is he coming?”). Also, there are syllabic nasal consonants, which are tone bearing units in the language. The nasal consonants always occur before a consonant. For example: *ndo* ‘Sorry’ or explicitly tone marked as *ndó*.

2.3 Writing System

The Igbo orthography is based on the Standard Igbo by the Ọnwụ Committee (Ọnwụ Committee, 1961). There are 28 consonants: *b gb ch d f g gh gw h j k kw kp l m n nw ny ñ p r s sh t v w y z*, and 8 vowels (see phonology section). Nine of the consonants are digraphs: *ch, gb, gh, gw, kp, kw, nw, ny, sh*.

Igbo is an agglutinative language in which its lexical categories undergo affixation, especially the verbs, to form a lexical unit. For example, the word form *ericharịrị* is a verbal structure with four morphemes: verbal vowel prefix *e-*, verb root *-ri-*, extensional suffix *-cha-*, and a second extensional suffix *-rịrị*. Its occurrence in the sentence “*Obi must eat up that food*” is *Obi ga-ericharịrị nri ahụ*, that is, *Obi aux-eat.completely.must food DET*. Igbo word order is Subject-Verb-Object (SVO), with a complement to the right of the head.

2.4 Grammatical Classes

Generally, Emenanjo (1978) identified the following broad word classes for Igbo: verbal, nominal, nominal modifier, conjunction, preposition, suffixes, and enclitics. The verbal is made up of verbs, auxiliaries and participles, while the nominal is made up of nouns, numerals, pronouns and interrogatives. Nouns are further classified into five lexical classes, viz; proper, common, qualificative, adverbial and ideophones. However, we identified extra five in the tagset design phase (see the appendix). Nominal modifiers occur in a noun phrase. Its four classes are adjectives, demonstratives, quantifiers and pronominal modifiers. Conjunctions link words or sentences together, while prepositions are found preceding nominals and verbals and cannot be found in isolation. Suffixes and enclitics are the only bound elements in the language. Suffixes are primarily affixed to verbals only, while enclitics are used with both verbals and other word classes. Suffixes are found in verb phrase slots and enclitics can be found in both verb phrase and noun phrase slots. The language does not have a grammatical gender system.

3 Language Resources

The development of NLP resources for any language is based on the linguistics resources available for the language. This includes appropriate fonts and text processing software as well as the available electronic texts for the work. The font and software problems of the language have been addressed through the Unicode development (Uchechukwu, 2005; Uchechukwu, 2006). The next is the availability of Igbo texts.

Any effort towards the Igbo corpus development is a non-trivial task. There are basic issues connected with the nature of the language. The first major surprise is that Igbo texts ‘by native speakers’ written ‘for native speakers’ vary in forms due to dialectal difference and are usually not tone-marked. Indeed, the tone marking used in the sections above are usually found in academic articles. It would be strange to find an Igbo text (literary work) that is fully tone marked and no effort has been made to undertake a tone marking of existing Igbo texts. Such an effort looks impossible as more Igbo texts are written and

published. Such is the situation that confronts any effort to develop an Igbo corpus. Hence, developing NLP resources for the language has to start with the available resources; otherwise, such an endeavour would have to first take a backward step of tone marking all the texts to be added to its corpus and normalizing the dialectal differences. This is a no mean task.

It is for this reason that we chose the New World Translation (NWT) Bible version for Igbo corpus with its English parallel text³. The NWT Bible does not adopt a particular tone marking system, neither is there a consistent use of tone marks for all the sentences in the Bible. Instead, there is narrow use of tone marks in specific and restricted circumstances throughout the book. An example is when there is a need to disambiguate a particular word. For instance, *ihé* without tone mark could mean ‘thing’ or ‘light’. These two are always tone marked in the Bible to avoid confusion; hence *ihè* ‘light’ and *íhé* ‘thing’. The same applies to many other lexical items. Another instance is the placement of a low tone on the person pronouns to indicate the onset of an interrogative sentence, which otherwise would be read as a declarative sentence. This particular example has already been cited as one of the uses of tone mark in the language. Apart from such instances, the sentences in the Bible are not tone marked. As such, one cannot rely on such restricted use of tone marks for any major conclusions on the grammar of the language. With regard to corpus work in general, the Bible has been described as consistent in its orthography, most easily accessible, carefully translated (most translators believe it is the word of God), and well structured (books, chapters, verses), etc. (Resnik et al., 1999; Kanungo and Resnik, 1999; Chew et al., 2006). The NWT Bible is generally written in standard Igbo.

4 Tokenization

We outline here the method we used in the tokenization of the text. For the sake of a start-up, we tokenized based on the whitespace. The Igbo language uses whitespace to represent lexical boundaries; we used the following regex:

Separate characters if the string matches:

- “ga-” or “n” or “N” or “na-” or “Na-” or “ana-” or “ina-”; for example, the following samples *n’elu*, *na-erughari*, *ina-akwa*, *ana-egbu* in the Bible will be separated into *n’*, *elu*, *na-*, *erughari*, *ina-*, *akwa*, *ana-*, *egbu* tokens.
- Any non-zero length sequence consisting of a–z, A–Z, 0–9, combining grave accent (`), combining acute accent (´), combining dot below (˙); for example, these words *ihè*, *ahú*, *ájá* in the corpus will be separated as tokens with their diacritics.
- Any single character from: left double-quotation mark (“), right double-quotation mark (”), comma (,), colon (:), semicolon (;), exclamation (!), question (?), dot (.).
- Any single non-whitespace character.

In place of sentence splitting, we use verses since all 66 books of the Bible is written in verse level. Our major aim is to use this Igbo corpus to implement our new tagset, which will capture all the inflected and non-inflected tokens in the corpus. For lack of space, issues with tokenization with respect to morphemes, manual annotation implementations and platform used will not be discussed in this paper.

5 Tagset Design

We adopt the (Leech, 1997) definition of a POS tagset as a set of word categories to be applied to the tokens of a text. We designed our tagset following the standard EAGLES guidelines, diverging where necessary (e.g. EAGLES, which favours European languages, specifies *articles* at the obligatory level, but this category does not apply for Igbo). A crucial question in tagset design is the extent of fine-grained distinctions to encode within the tagset. A too coarsely grained tagset may fail to capture distinctions that would be valuable for subsequent analysis, e.g. syntactic parsing; too fine-grained may make automatic (and manual) POS tagging difficult, resulting in errors that lead to different problems for later processing. In what follows, we introduce a sizeable tagset granularity with the intention of providing a basis for practical POS tagging.

³Obtained from jw.org.

NNM	Number marking nouns	NNT	Instrumental nouns
NNQ	Qualificative nouns	VrV	– <i>rV</i> implies suffix
NND	Adverbial nouns	VCJ	Conjunctive verbs
NNH	Inherent complement nouns	α _XS	any POS tag with affixes
NNA	Agentive nouns		

Table 1: Selected distinctive tags from the tagset scheme

The tagset is intended to strike an appropriate balance for practical purposes regarding granularity, capturing what we believe will be the key lexico-grammatical distinctions of value for subsequent processing, such as parsing. Further subcategorization of the grammatical classes, as described in section 2.4, results in 59 tags which apply to whole tokens (produced by the tokenisation stage described above). An important challenge comes from the complex morphological behaviour of Igbo. Thus, a verb such as *bja*, which we assign the tag VSI (a verb in its simple or base form), can combine with extensional suffixes, such as *ghi* and *kwa*, to produce variants such as *bjaghi*, *bjakwa* and *bjaghikwa*, which exhibit similar grammatical behaviour to the base form. As such, we might have assigned these variants the VSI tag also, but have instead chosen to assign VSI_XS, which serves to indicate both the core grammatical behaviour and the presence of extensional suffixes. In *abjakwa*, we find the same base form *bja*, plus a verbal vowel prefix *a*, resulting in the verb being a participle, which we assign the tag VPP_XS. For the benefit of cross-lingual training and other NLP tasks, a smaller tagset that captures only the grammatical distinctions between major classes is required. The present 59 tags can easily be simplified to a coarse-grained tagset of 15 tags, which will principally preserve just the core distinctions between word classes, such as nouns, verb, adjective, etc.

Although Emenanjo (1978) classified **ideophones** as a form of noun, we have assigned them a separate tag **IDEO**, as these items can be found performing many grammatical functions. For instance, the **ideophone** *koi*, “to say that someone walks *koi koi*” has no nominal meaning, rather its function here is adverbial. A full enumeration of this scheme is given in the appendix.

5.1 The development of an POS tagged Igbo Corpus

Here we analyse the manual POS tagging process that is ongoing based on the tagset scheme. The Bible books were allocated randomly to six groups, producing six corpora portions of approximately 45,000 tokens each. Our plan was for each human annotator to tag at least 1000 tokens per day, resulting in complete POS tagging in 45 days. The overall corpus size allocated is 264,795 tokens of the new testament Bible. There are six human annotators, who are students of the Department of Linguistics at Nnamdi Azikiwe University, Awka, supervised by a senior lecturer in the same department; giving an effective total of seven human annotators. Additionally, a common portion of the corpus (38,093 tokens) was given to all the annotators, as a basis for calculating inter-annotator agreement.

6 Conclusions

We have outlined our current progress in the development of a POS tagging scheme for Igbo from scratch. Our project aims to build linguistic computational resources to support research in natural language processing (NLP) for Igbo. It is important to note that these tags are applicable on unmarked, not fully marked, and fully tone marked Igbo texts, since the fully tone marked tokens play the same grammatical roles as in the none tone marked texts, written by native speakers for fellow native speakers.

Our method of tagset design could be used for other African or under-resourced languages. African languages are morphologically rich, and of around 2000 languages in the continent, only a small number have featured in NLP research.

Acknowledgements

We acknowledge the support of Tertiary Education Trust Fund (TETFund) in Nigeria, and would like to thank Mark Tice for his useful comments and help in preparing this paper.

References

- Cheikh M. Bamba Dione, Jonas Kuhn, and Sina Zarriß. 2010. Design and Development of Part-of-Speech-Tagging Resources for Wolof (Niger-Congo, spoken in Senegal). In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*. ELRA).
- Peter A. Chew, Steve J. Verzi, Travis L. Bauer, and Jonathan T. McClain. 2006. Evaluation of the Bible as a Resource for Cross-Language Information Retrieval. In *Proceedings of the Workshop on Multilingual Language Resources and Interoperability*. Association for Computational Linguistics.
- E. Nọlue Emenanjo. 1978. *Elements of Modern Igbo Grammar: A Descriptive Approach*. Ibadan Ox. Uni. Press.
- Margaret M. Green and G. Egemba Igwe. 1963. *A descriptive grammar of Igbo*. London: Oxford University Press and Berlin: Akademie-Verlag.
- Clara Ikekeonwu. 1999. "Igbo", *Handbook of the International Phonetic Association*. C. U. Press.
- Tapas Kanungo and Philip Resnik. 1999. The Bible, Truth, and Multilingual OCR Evaluation. In *Proceedings of SPIE Conf. on Document Recognition and Retrieval*, pages 86–96.
- Geoffrey Leech. 1997. *Introducing Corpus Annotation*. Longman, London.
- P. Akujuobi Nwachukwu. 1995. Tone in Igbo Syntax. Technical report, Nsukka: Igbo Language Association.
- Guy de Pauwy, Gilles-Maurice de Schryver, and Janneke van de Loo. 2012. Resource-Light Bantu Part-of-Speech Tagging. In *Proceedings of the Workshop on Language Technology for Normalisation of Less-Resourced Languages*, pages 85–92.
- Philip Resnik, Mari Broman Olsen, and Mona Diab. 1999. The Bible as a Parallel Corpus: Annotating the 'Book of 2000 Tongues'. *Computers and the Humanities*, 33.
- Chinedu Uchechukwu. 2005. The Representation of Igbo with the Appropriate Keyboard. In Clara Ikekeonwu and Inno Nwadike, editors, *Igbo Lang. Dev.: The Metalanguage Perspective*, pages 26–38. CIDJAP Enugu.
- Chinedu Uchechukwu. 2006. Igbo Language and Computer Linguistics: Problems and Prospects. In *Proceedings of the Lesser Used Languages and Computer Linguistics Conference*. European Academy (EURAC).
- Chinedu Uchechukwu. 2008. African Language Data Processing: The Example of the Igbo Language. In *10th International pragmatics conference, Data processing in African languages*.
- Beatrice F. Welmers and William E. Welmers. 1968. Igbo: A Learner's Manual. Published by authors.
- Ọnwụ Committee. 1961. The Official Igbo Orthography.

A A Tagset Design for the Igbo Language

Noun Class	
Tag	Description/Example
NNP	Noun Proper. <i>Chineke</i> 'God', Onyeka, Okonkwo, Osita.
NNC	Noun Common. <i>Oku</i> 'fire', <i>uwa</i> 'earth', <i>osisi</i> 'tree, stick', <i>ala</i> 'ground', <i>eluigwe</i> 'sky, heaven'
NNM	Number Marking Noun. <i>Ndi</i> 'people', <i>nwa</i> 'child', <i>umu</i> 'children'. <i>ndi</i> is classified as a common noun with an attached phrase of "thing/person associated with" (Emenanjo, 1978). <i>ndi</i> preceding a noun marks plurality of that noun, <i>nwa</i> marks it singular (e.g. <i>nwa agboghọ</i> 'a maiden'), and <i>umu</i> also indicate plurality (e.g. <i>umu agboghọ</i> 'maidens').
NNQ	Qualificative noun. Nouns that are inherently semantically descriptive. E.g. <i>ogologo</i> [height, long, tall]
NND	Adverbial noun. This lexical class function to modify verbals, e.g. <i>O ji nwayọọ eri nri ya</i>
NNH	Inherent Complement. Igbo verb has a [verb + NP/PP] structure. NP/PP are the verb complement. They cooccur with the verb, at times quite distant from the verb, e.g. (1) <i>igu egwu</i> 'to sing', (2) <i>iti igba</i> 'to drum', (3) <i>igwu ji</i> 'harvest yam'.
NNA	Agentive Noun. Nouns are formed through verbs nominalization. Compare (1) with <i>ogo egwu</i> 'singer' and (2) with <i>oti igba</i> 'drummer'. For links <i>NNAV</i> ... <i>NNAC</i> .
NNT	Instrumental Noun. Refer to instruments and are formed via nominalization. Compare (3) with <i>ngwu ji</i> 'digger'. For links <i>NNTV</i> ... <i>NNTC</i> .
NOTE: We introduced link indicators in NNA and NNT, V and C, Where V and C stand for verbal and Complementary respectively. So, <i>NNAV</i> indicates derivation from the verbal component of the inherent complement verb and <i>NNAC</i> is the inherent complement of the whole verbal complex. E.g., <i>ogu/NNAV egwu/NNAC</i> . Also, <i>NNTV</i> and <i>NNTC</i> , where <i>NNTV</i> is derived from the verbal component of the inherent complement verb and <i>NNTC</i> is the inherent complement of the whole verbal complex. E.g., <i>ngwu/NNTV ji/NNTC</i>	

Verb Class	
VIF	Infinitive. Marked through the addition of the vowel [i] or [i] to the verb root.
VSI	Simple verb. Has only one verb root.
VCO	Compound Verb. Involves a combination of two verb roots.
VIC	Inherent Complement Verb (ICV). Involves the combination of a simple or compound verb with a noun phrase or a prepositional phrase. It gives rise to the structures (1) V + NP, or (2) V + PP
VMO	Modal Verb. Its formed by inherent complement verbs and simple verbs. [See the section on suffixes]
VAX	Auxiliary Verb. ga [Future marking], na [progressive]
VPP	Participle. Always occurs after the auxiliary, and prefixed e/a to the verb root using vowel harmony.
V CJ	Conjunctive Verb. A verb that has a conjunctive meaning, especially in narratives: <i>wee</i>
VBC	Bound Verb Compliment or Bound Cognate Noun. Its formed by harmonizing prefix <i>a/e</i> to the verb root. It
(BVC)	looks like the participle but occurs after the participle in same sentence as the verb. It can be formed from every verb.
VGD	Gerund. Reduplication of the verb root plus harmonizing vowel o/ọ. Also, internal vowel changes can occur. E.g. <i>ba</i> ‘enter’ [ọ + bụ + ba]= <i>ọbuba</i> ‘the entering’
Inflectional Class	
VrV	– <i>rV</i> (e.g. <i>-ra</i>). If attached to an active verb, it means simple past; but a stative meaning with a stative verb.
VPERF	Perfect (e.g. <i>-la/-le, -go</i>). Describes the ‘perfect tense’. <i>-la/-le</i> obeys vowel harmony and the variant <i>-go</i> does not.
Other part-of-speech tags	
ADJ	Adjective. The traditional part of speech ‘adjective’ that qualifies a noun. Igbo has very few of them.
PRN	Pronoun. The 3 persons are 1st (sing + pl), 2nd (sing + pl), and 3rd (sing + pl) person Pronouns.
PRNREF	Reflexive Pronoun. Formed by combination of the personal pronouns with the noun <i>onwe</i> ‘self’.
PRNEMP	Emphatic pronoun. This involves the structure [pronoun+onwe+pronoun].
ADV	Adverb. Changes or simplifies the meaning of a verb. They are few in Igbo.
CJN	Conjunction. There are complex and simple conjunctions distinguish based on grammatical functions viz; co-ordinators, sub-ordinators and correlatives. Link indicators <i>CJN1...CJN2</i> are for “correlative CJN”. E.g. <i>ma/CJN1...ma/CJN2</i> .
PREP	Preposition. The preposition <i>na</i> is realised as <i>n’</i> if the modified word begins with a vowel.
WH	Interrogative. Questions that return useful data through explanation. <i>Ọnye, ọ́nị, olee, ...</i>
PRNYNQ	Pronoun question. Questions that return YES or NO answer. E.g. <i>m, à, hà, ò, `ọ, ...</i>
IDEO	Ideophone. This is used for sound-symbolic realization of various lexico-grammatical function. E.g. <i>niganiga, murii, koi, etc.</i>
QTF	Quantifier. This can be found after their nominals in the NP structure. E.g. <i>dum, naabo, nille.</i>
DEM	Demonstrative. This is made up of only two deictics and always used after their nominals. E.g. <i>a, ahụ.</i>
INTJ	Interjection. <i>Ee</i>
FW	Borrowed word. <i>amen.</i>
SYM	Punctuation. It includes all symbols.
CD	Number. This includes all digits 1,2,3, ... and <i>otu, mbi, abua, atọ, ...</i>
DIGR	Digraph. All combined graphemes that represent a character in Igbo, which occur in the text. <i>gb, gw, kp, nw, ...</i>
TTL	Title . Includes foreign and Igbo titles. E.g. <i>Maazi.</i>
CURN	Currency.
ABBR	Abbreviation.
Any type of suffixes	
α_XS	any POS tag with affixes. for $\alpha \in \{VIF, VSI, VCO, VPP, VGD, VAX, CJN, WH, VPERF, VrV, PREP, DEM, QTF, ADJ, ADV\}$. See verb, other POS, inflectional classes.
NOTE: Tags with affixes identify inflected token forms in the corpus for use in further analysis, e.g. morphology. For practical POS tagging, such tags may be simplified, i.e. $\alpha_XS \Rightarrow \alpha$.	
Any type of Enclitics	
ENC	Collective. <i>cha, si nu, ko</i> – means all, totality forming a whole or aggregate. Negative Interrogative. <i>di, ri, du</i> – indicates scorn or disrespect and are mainly used in Rhetorical Interrogatives. Adverbial ‘Immediate present and past’. <i>fo/hu</i> – it indicates action that is just/has just taking/taken place. <i>rii</i> – indicates that an action/event has long taken place Adverbial ‘Additive’. <i>kwa (kwọ), kwu</i> – mean ‘also’, ‘in addition to’, ‘denoting’, ‘repetition or emphasis’. Adverbial ‘Confirmative’. <i>noo (noo; nnoo)</i> – this means really or quite.

B The Major Classes of the Tagset

ADJ	adjective	FW	foreign word	QTF	quantifier	ADV	adverb	NNC	common noun
INTJ	interjection	SYM	symbol	CJN	conjunction	NNP	proper noun	PREP	preposition
WH	interrogative	PRN	pronoun	V	verb	CD	number	DEM	demonstration
There is no article in the language.									