

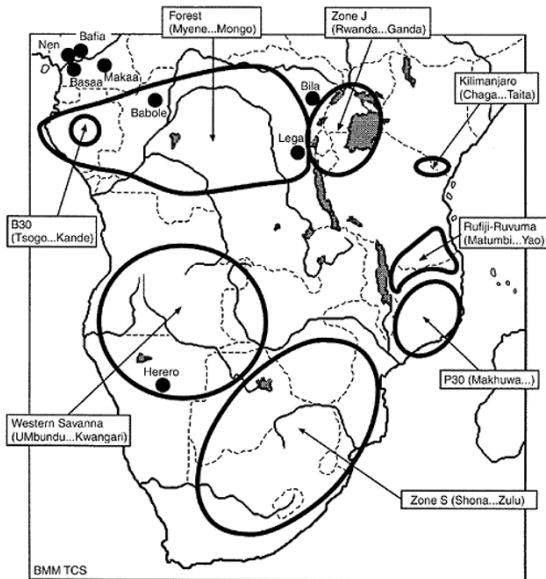
Computational Resources and Tools Developed for Bantu Languages

Joan Byamugisha

August 25, 2021



- 1 Background on Bantu Languages
- 2 Computational Resources
- 3 Computational Tools



Noun Class	Description of Associated Nouns
1 and 2	People and kinship
3 and 4	Plants, nature, and some parts of the body
5 and 6	Fruits, liquids, some parts of the body, and paired things
7 and 8	Inanimate objects
9 and 10	Tools and animals

- Stem: *-ntu*
 - Person: ***omuntu***
 - Thing: ***ekintu***
 - Place: ***ahantu***
- Stem: *-nyankore*
 - Person from Ankore: ***omunyangore***
 - Language of Ankore: ***orunyangore***
- The noun prefix is based on a class prefix

Noun Class (NC), Subject Concord (SC), Possessive Concord (PC), Adjective Concord (AC)

NC	SC	PC	AC
1. o-mu-	-a-	o-wa	o-mu-
2. a-ba-	-ba-	a-ba	a-ba-
3. o-mu-	-gu-	o-gwa	o-mu-
4. e-mi-	-gi-	e-ya	e-mi-
5. ei-/e-ri-	-ri-	e-rya	e-ri-
6. a-ma-	-ga-	a-ga	a-ma-

- Comprises bound morphemes, a verb-root, and extensions
- Morphemes preceding verb-root specify person, noun class, aspect, time, negation
- Extensions specify valency-changing categories
- <Initial> <Subject> <Negative> <Tense and/or Aspect>
<Object> <Root> <Extension> <Final>
- For tense and aspect:
 - Tense reference occurs before that of the aspect
 - Most languages have one, two, three, or four past tenses
 - Most languages have one, two, or three discrete future tenses

- Runyankore: *Titukakimureeterahoganu*.
- Morphemes: Ti-tu-ka-ki-mu-reet-er-a-ho-ga-nu
- English: We have never ever brought it to him/her.

- Generated a one million sentence general-purpose domain independent corpus
- Created 18,816 different ways of varying the sentence structure
- Generated corpora labeled for morphology and sentiment
- Generated text in seven tenses:

- $S \rightarrow IG FM LM IF VR EX FN$
- $IG \rightarrow PN IT SN$
- $PN \rightarrow ti \mid ni$
- $IT \rightarrow a \mid o \mid n \mid tu \mid mu \mid \dots$
- $SN \rightarrow ta$
- $FM \rightarrow za \mid ka \mid riku \mid rikuza$
- $LM \rightarrow ki$
- $IF \rightarrow mu \mid ba \mid \dots$
- $VR \rightarrow verbRoot$
- $EX \rightarrow w \mid er \mid erer \mid \dots$
- $FN \rightarrow a \mid e \mid ire$

- Pre-trained word embeddings from distributional context from one million sentences
- Pre-trained model on morphology
- Pre-trained sentence-level sentiment classifier

- Deals with determining the noun class when the class prefix is the same
- Runyankore
 - ***Omuti***
 - ***Omuntu***
 - ***Omwaka***
 - ***Omwana***
- Kinyarwanda
 - ***Umugore***
 - ***Umutima***
 - ***Umwembe***
 - ***Umwalimu***

Query Word	Results
<i>omuntu</i> (person)	<i>omugyesi</i> (reaper), <i>omutaahi</i> (companion), <i>omukoreesa</i> (overseer), <i>omushomesa</i> (teacher), <i>omukuru</i> (elder)
<i>omuti</i> (tree)	<i>omutumba</i> (banana tree), <i>omwani</i> (coffee tree), <i>omuzaabibu</i> (grape or grapevine), <i>omucungwa</i> (orange), <i>omugusha</i> (sorghum)
<i>omukono</i> (arm)	<i>omunwa</i> (mouth), <i>omutwe</i> (head), <i>eriino</i> (tooth), <i>enkokora</i> (elbow), <i>okuguru</i> (leg)
<i>embwa</i> (dog)	<i>embeba</i> (rat), <i>enkyende</i> (monkey), <i>empungu</i> (bird of prey), <i>enumi</i> (bull), <i>enyawaawa</i> (green ibis)

Noun Class	Description of Associated Nouns
1 and 2	People and kinship
3 and 4	Plants, nature, and some parts of the body
5 and 6	Fruits, liquids, some parts of the body, and paired things
7 and 8	Inanimate objects
9 and 10	Tools and animals

Approach	Runyankore	Luganda	Kinyarwanda
Morphological only	69.23	57.53	43.94
Semantic only	66.67	47.95	40.91
Morphology, Syntax, Semantics	87.18	73.97	63.64

- Word segmentation from pre-trained model on morphology
- Sentiment analysis from pre-trained sentiment classifier
- From Treebank
 - Performing machine translation
 - Automated evaluation of generated versus human authored text
 - Computing linguistic diversity