A Brief Analysis of Amharic NLP: From POS Tagging to Question Answering

Seid Muhie Yimam Language Technology Lab TU Darmstadt 6 September 2016

My profile

- Doctoral researcher at Language Technology lab
- My Topic: Narrowing the loop: integration of resources and linguistic dataset development with interactive machine learning
 - Adaptive annotation in WebAnno
 - Semantic writing aid using contextual paraphrasing

Ecosystem & applications at LT



Disclaimer: I am not a seasoned Amharic NLP researcher. I have done my M.Sc. thesis on Amharic QA in 2008 but rarely participate on Amharic NLP then after. And by no means, this talk is not a detailed analysis of current Amharic NLP.

Outline

- NLP research and applications in Amharic
 - POS tagging, Morphological processing, Spell checking, Named entity recognition, and Questions Answering
- Our Contribution to Amharic NLP
 - POS, QA, Corpora development
- Challenges and bottlenecks in Amharic NLP
- Collaboration with the Ethiopistic department

Introduction

- Amharic writing system: version of the Ge'ez script known as ፊደል (Fidel)
- Ethiopic characters (fidels) have more than 380 Unicode representations (U+1200-U+137F)



	a/ä	u	i	а	е	ə	0		a/ä	u	i	а	е	(ə)	0
	[a/ε]	[u]	[i]	[a]	[e]	[¥]	[o/ɔ]		[a/ɛ]	[u]	[i]	[a]	[e/ɛ]	[i/u]	[o/ɔ]
h	υ	ሁ	ሂ	Y	ሂ	บ	ሆ	h/k	ň	ዅ	ኺ	ሻ	ኼ	ሽ	ኾ
[h]	ha	hu	hi	ha	he	h(i)	ho	[h]	hε	hu	hi	ha	he	h(i)	ho
1	٨	ሉ	ሊ	٩	∿	6	ک	w	D	ጨ	ዊ	ዋ	ዌ	ው	ዎ
[1]	le	lu	li	la	le	l(i)	lo	[w]	wε	wu	wi	wa	we	w(ʉ)	wo
h/ḥ	ሐ	ሑ	ሒ	ሓ	ሔ	ሕ	ሖ	•/¢	0	ዑ	ዒ	ዓ	જ	ò	8
[h]	ha	hu	hi	ha	he	h(i)	ho	[2]	?a	?u	?i	?a	?e	?i	?o
m	ത	ሙ	ሚ	ማ	ሚ	ም	ሞ	z	H	ŀł	H,	મ	ዜ	ห	Ч
[m]	mε	mu	mi	ma	me	m(i)	mo	[z]	z٤	zu	zi	za	ze	z(i)	zo
s/ś	ມ	ມະ	ખ્	ሣ	պ	ሥ	ሦ	zh/ž	H	ዡ	ዢ	મ	ዤ	ዠ	H
[s]	se	su	si	sa	se	s(i)	so	[3]	38	3u	3i	за	зe	3(i)	30
r	ሪ	ሩ	Տ	ራ	ሬ	ር	ሮ	У	የ	ķ	ዪ	ç	ዬ	e	ዮ
[r]	r٤	ru	ri	ra	re	r(i)	ro	[j]	jε	ju	ji	ja	je	j(i)	jo
s	٨	ሱ	ሲ	ሳ	ሌ	ስ	ሳ	d	ደ	ዱ	ዲ	ዳ	ዴ	ድ	ጾ
[s]	se	su	si	sa	se	s(i)	so	[d]	đ٤	du	di	da	đe	d(i)	do
sh/š	ሽ	ሹ	ሺ	ሻ	ሼ	ሸ	ሻ	j/ğ	Ĕ	ኟ	ኟ	፞፞፞፞፞	ጀ	ጅ	ጆ
ហ	ſε	∫u	∫i	∫a	∫e	∫(i)	∫o	[¢]	đzε	фu	фi	dза	фе	₫(i)	фo
k'/q	ቀ	ቂ	ቂ	ቃ	ቆ	ቅ	ዋ	g	1	Դ	l	3	เ	ๆ	1
[k']	k'e	k'u	k'i	k'a	k'e	k'(i)	k'o	[g]	g٤	gu	gi	ga	ge	g(i)	go

			-									-			
qh	ቐ	ቒ	ቒ	ቓ	ቘ	ቐ	ኞ	ť'/ț	ጠ	ጡ	ጢ	ጣ	ጤ	ፕ	ጣ
[ĸ,]	в ,5	в'u	в,i	в,я	в,e	r,(i)	r,o	[ť]	ť٤	ťu	ťi	ťa	ťe	ť'(i)	ťo
b	ึก	ቡ	ቢ	ŋ	ቤ	ብ	ր	ch'/č	ጨ	ጩ	ጪ	ጫ	ጬ	ጭ	ጮ
[b]	bε	bu	bi	ba	be	b(i)	bo	[ť]	t∫'ε	t∫'u	t∫'i	t∫'a	t∫'e	t∫'(i)	ť∫'o
t	ተ	ቱ	ቲ	ታ	ቴ	ት	ቶ	p'/p	8	ጱ	ጲ	8	ጴ	ጽ	8
[t]	te	tu	ti	ta	te	t(i)	to	[p']	p'e	p'u	p'i	p'a	p'e	p'(i)	p'o
ch/č	ቸ	ቹ	ቺ	ቻ	ቼ	ቸ	ቾ	ts'/ș	ጸ	ጹ	ጺ	१	ጼ	ጽ	8
[ʧ]	t∫ε	t∫u	t∫i	t∫a	t∫e	t∫(i)	t∫o	[ts]	ts'e	ts'u	ts'i	ts'a	ts'e	ts'(i)	ts'o
h/ḫ	ኅ	ኁ	ኂ	ኃ	ኄ	ኅ	ኆ	ts'/ś	0	ፁ	ጚ	ን	ጜ	è	1
[h]	ha	hu	hi	ha	he	h(i)	ho	[ts]	ts'e	ts'u	ts'i	ts'a	ts'e	ts'(i)	ts'o
n		ኑ	ኒ	ና	ኔ	ን	ኖ	f	ፈ	4	ሌ	ፋ	ፌ	ፋ	ፎ
[n]	nε	nu	ni	na	ne	n(i)	no	[f]	fe	fu	fi	fa	fe	f(i)	fo
ny/ñ	ኝ	ች	፟፟፟ጟ	ኛ	ኜ	፝ጛ	ኛ	р	T	F	T	T	ጌ	T	7
[ŋ]	jnε	'nu	ni	'na	ne	ŋ(i)	лo	[p]	pε	pu	pi	pa	pe	p(i)	po
م ،	ħ	ኩ	ኪ	ካ	ኬ	ħ	た	v	ที	ሹ	ቪ	ሻ	ቬ	ቭ	ሸ
[?]	?a	?u	?i	?a	?e	?i	?o	[v]	v٤	vu	vi	va	ve	v(i)	vo
k	n	ኩ	h,	η	ռ	ħ	n								
[k]	kε	ku	ki	ka	ke	k(i)	ko								

ቊ	ቋ	ቌ	ቍ	ኍ	፟	ኌ	ኍ	ኵ	ኣ	ħ.	ኵ	ጒ
k'"i	k' ^w a	k' ^w e	k'"e	h ^w i	h ^w a	h ^w e	h ^w ε	k ^w i	k ^w a	k ^w e	k ^w ε	$\mathbf{g}^{\mathbf{w}}\mathbf{i}$
ኋ	l	ጕ	ሏ	Ŋ	મ્	ጧ	ሟ	ቷ	મ્	ጯ	ሯ	ቿ
g ^w a	g ^w e	g ^w e	l ^w a	b ^w a	z ^w a	t' ^w a	m ^w a	t ^w a	3 ^w a	t∫' ^w a	r ^w a	t∫ ^w a
ኟ	ጿ	ሏ	ኗ	ደ	ሏ	ሿ	ኟ	ሪ	ሻ	ጜ	ኸ	
ф ^w a	ts' ^w a	s ^w a	n ^w a	d ^w a	f ^w a	∫ ^w a	'n ^w a	r ^j a	m ^j a	f ^j a	?ε	
-		;;		Ī		÷			-			!
comma	full sto	ip / perio	d	colon		semi-co	olon	prefa	ce colo	n q (no	luestio Diong	on mark er used
õ	ĕ	<mark>ĵ</mark>	ĝ	Ž		į.	<u>7</u>	Ţ		ğ	Ĩ	[

ŏ	ğ	1	У.	Ģ	8	4	\$	Ņ	-
አንዶ	ሁለት	ሶስት	አራት	<mark>አ</mark> ምስት	ስዶዶስት	ሰባት	ስምንት	ዘጠኝ	አስር
and	hulätt	sost	aratt	amməst	səddəst	säbatt	səmmənt	zäţäňň	asser
1	2	3	4	5	6	7	8	9	10
Χ̈́	ល្អ	ମୁ	Ŷ	Ţ,	Ĝ	Ť	ĩ	P	፼
08	ሰሳሳ	አርባ	ሀምሳ	ስልሳ	ሰባ	ሰማን <i>ያ</i>	ዘጠና	መቶ	T.
haya	sälasa	arba	hamsa	səlsa	säba	sämaňa	zäțena	mäto	ši
20	30	40	50	60	70	80	90	100	1000

Amharic morphology [1/2]

 Amharic is morphologically complex Example Noun: መሬት (meret) ~ land መሬቶች (meretoch) ~ lands መሬቶቻችን (meretochachin)~ our lands ስለመሬቶቻችን (silemeretochachin)~ about our lands ስለመሬቶቻችንስ? (silemeretochachins)~ what about our lands?

Amharic morphology [2/2]

Example Verb

መስበር (mesber) ~ break

ሰበራቸው (seberachihu)~ you break (s/t)

ተሰበራችው (teseberachihu) ~ you are broken

ተሰባበራችው (tesebaberachihu) ~ you are broken (to portions)

<mark>ስላል</mark>ተሰባበራችው (silaltesebaberachihu) ~ as you are not broken

ስላልተሰባበራችውም (silaltesebaberachihum) ~ and as you are not broken

A single verb may consist time (past, present, and future), gender (male and female), action (command, statement, invitation) and negation (not).

HORNMORPHO 2.5^[1]

 morphological analysis and generation of Amharic and Oromo verbs and nouns and Tigrinya verbs

```
>>> 13.anal('am', 'ስላልተሰባበራቸሁም')
Loading morphological data for Amharic ...
word: ስላልተሰባበራቸሁም
POS: verb, root: <sbr>, citation: ተሰባበረ
subject: 3, sing, masc
object: 2, plur
grammar: perfective, iterative, passive, relative, definite, negative
preposition: sIle, conjunctive suffix: m
POS: verb, root: <sbr>, citation: ተሰባበረ
subject: 2, plur
grammar: perfective, iterative, passive, relative, negative
preposition: sIle, conjunctive suffix: m
```

- No notion of "spelling" in Amharic^[1]
 - "if a word sounds right when read aloud then it was rightly written"
 - "ውኃ" vs "ዉሀ", "ታህሳስ" vs "ታኅሣሥ" -> acceptable
 - "<mark>ዓ</mark>ዲሥ ዐበባ <u>ዒ</u>ትዮጵያ" vs "አዲስ አበባ ኢትዮጵያ not acceptable

- Different levels (elementary, intermediate, advanced)
- "**ባ**ል", "መልክት", "<mark>እንቁ</mark>ጣጣሽ", "<mark>አ</mark>ይን", "<mark>አ</mark>ሳ", "**ያ**ገር"
- "<mark>በዓ</mark>ል", "መልዕክት", "<mark>እንቍ</mark>ጣጣሽ", "<mark>ዓ</mark>ይን", *አሣ*"፤"የአገር"
- "በዐል", "መል<mark>እ</mark>ክት", "<mark>ዕንቍ</mark>ጣጣሽ", "ዐይን", "ዐሣ", "የሀገር"

- Metaphone algorithm and an edit distance algorithm^[1]
 - 1. Get WORD in the dictionary
 - 2. Compute Amharic Metaphone code for *WORD*, *W_CODE*
 - IF WORD has no similar W_CODE with a previously hashed word:

Add *W_CODE* as a key and *WORD* as value in a hash table, *H_DICT*

ELSE:

Add WORD as the next value in the respective key

4. IF the dictionary has another WORD:

GOTO 1

[1] http://thirdworld.nl/development-of-an-amharic-spelling-corrector-for-tolerant-retrieval

Results

Rank of Correct Suggestions	Count	Percentage
Top One	55	67.1%
Top Three	61	74.4%
Top Five	67	81.7%
N/A	15	18.3%
Total	82	100%

POS tagging [1/3]

- Getachew (2001) : Hidden markov models
 definition of a tagset of 25
- Adafre (2005): stochastic model based on conditional random fields
 - revised Getachew's tagset and reduced it to 10
 - average accuracy of 74%
 - used dictionaries of affixes
 - some 15,000 entries with their POS tags (Noun, Verb, Adjectives, Adverb, and Adposition)
- (Demeke and Getachew, 2006):
 - corpus of 1065 news articles
 - 31 parts of speech
 - Multi-word expressions use separate tag

POS Tagging [2/3]

- Gamback et al. (2009): using TnT, SVMTool, and Mallet
 - Accuracy of 85.56% for TnT, 88.30% for SVM and 87.87% for MaxEnt
- Tachbelie and Menzel (2009): using TnT and SVMTool models
 - accuracies of 82.99% for TnT and 84.44% for SVM

POS Tagging

- Binyam (2011)
 - clean the WIC data
 - Use vowel patterns and the radicals as POS features
 - State-of-the-art tagging models (CRF++, LIBSVM, and TnT)

POS tagging problems

- The corpus used is usually small in size
- The quality of the corpus is poor
 - A lot of inconsistency
- Due to the agglutinative nature of the language
 - Fine-grained tagsets difficult to annotate large data set
 - Coarse-grained tagsets more ambiguous and might be wrong
- Lack of annotation guidelines
 - Computer programmers do not know much of the linguistic theory
 - Linguists do not know much of the programming

Definitive QA (Teshome 2013)^[1]

- Based on surface text pattern method
 - Design patterns to discover a set of definitionrelated text patterns from the Amharic legal corpus.
 - Extract a collection of concept-description pairs from a target document
 - Apply the definition extraction to return answer to a given question.
- Achieved F-measure of 78.8%

[1] http://etd.aau.edu.et/bitstream/123456789/8300/1/Wondwossen%20Teshome%20201 3.pdf

QA for list questions (Bete 2013)^[1]

- Closed domain (Ethiopian tourism)
- Hypothesis:
 - answers to a list questions have same semantic entity class
 - Answers that co-occur within the sentences of the documents are related to the target
 - The question and sentences containing the answers share similar context.

Amharic Named Entity recognition ^[1]

- Machine learning approaches
 - Conditional random fields
 - 80.66% of F-measure achieved
- Amahric NE
 - no case information like English
 - Clue words for each NE classes are used
 - ፕሬዚዳንት President, አቶ Mr. , ወ/ሮ Mrs
 - ካፒታል ሆቴል Capital Hotel, ሐረማያ ዩኒቨርሲቲ Haramaya University ፕሬዚዳንት ቢሮ, Office of president, የተባበሩት መንግስታት ድርጅት, United Nations Organization

[1] <u>http://etd.aau.edu.et/bitstream/123456789/8579/1/Besufikad%20Alemu%20final.pdf</u> ₂₃

NER challenges for Amharic

- Ambiguity:
 - Odle (Sun) person name as well the object sun
 - Some names are used as both person and location names
- Spelling variations
 - "ፀሐይ", "ጻሐይ", "ፀህይ", "ጻህይ" ~ sun
- Lack of capitalization

Our contribution

Factoid QA system for Amharic

- The first QA system for Amharic
- Components
 - Document pre-processing
 - Question processing
 - Document retrieval
 - Sentence/paragraph re-ranking
 - Answer selection modules.



Figure 1: Architecture of the system



Results

- 12000 question sets have been collected from the Web, Ethiopian television games and from questionnaire respondents.
- Question classification 89%
- Retrieval
 - Document based 97% recall
 - Sentence based 93 % recall
- Answer selection 72% recall

POS Tagging

- POS Automation ^[1]
 - Using WebAnno automation and correction module
 - manually annotate Amharic documents
 - Designed 11 POS tags equivalent to the universal POS tags

[1] <u>http://www.aclweb.org/anthology/P/P14/P14-5016.pdf</u>

POS tagging automation

- Initially manually annotated 21 sentences
- Iterative automatic annotation suggestion until 300 sentences are annotated.
- We obtained an F-score of 0.89 with the final model.

Annotation
PUNC PUNC NOUN NOUN NOUN NOUN ADV VERB PUNC 6 <
PRON VERB NOUN PRON NOUN NOUN NOUN NOUN VERB NOUN VERB 7 ይህንን ለማሳካትም ኢትዮጵያ ምንም ዓይነት የውጭ ዕርዳታ ለግድቡ እንዳታገኝ ግሬት እናደርጋለን
Suggestion
PUNC PUNC NOUN NOUN NOUN NOUN NOUN VERB PUNC 6 ሩ ሩ ግብፅ የህዳሴው ግድብ ግንባታን ፈፅሞ አትፈቅድም ።
PRON NOUN PRON NOUN PRON NOUN PRON NOUN NOUN NOUN NOUN NOUN NOUN VERB NOUN NOUN <t< td=""></t<>

Unsupervised POS tagging (Yemisrach - 2015)

- Co-advising a student master thesis work at Gondar University – Ethiopia
 - Develop fine-grained Amharic tagsets
 - Incorporating the unsupervised POS tags as features

Amharic text corpora [1/2]

- From tweeter using tweeter API
 - Download tweets everyday
 - 713683 tweets in three years

Year	Users per week	Tweet per week
2014	896	9438
2015	1234	6330
2016	1675	9438

Amharic tweet trends

Amharic text corpora [2/2]

- Collecting Amharic texts from the web using focused crawler
 - Around 1 million clean sentences are collected in one week time
 - Used for:
 - Unsupervised POS tagging feature
 - Spell checking word form preparation

- At CIDLeS Summer School 2014: Coding for Language Communities
 - Corpus based
 - Dictionary preparation— corpus frequency list
 - About 500,000 word forms
 - Moderate affix file preparation
 - Around 88% coverage on random web text

The main bottlenecks of Amharic NLP

- Lack of dedicated NLP research group
 - No cooperation or communication between researchers across different institutes
- Lack of Amharic keyboard (or it is complex)
 - Most social media users use the transliteration in stead of the Amharic script (Fidel)
- Most research are done as a master or doctoral thesis and
 - No continuity on the same topic
 - Resources are not available open source

Tasks and collaborations [1/2]

- Amharic tokenizer / segmenter (compound splitting, Sentence demarcation...)
- Pre-processing and normalization approaches
- Development of different corpora(annotation, crawling, cleaning...)
- Standardizing existing NLP tools
 - Guidelines and documentations
 - Open source and free
- Collaborate with different researchers and educators working on Amharic NLP
- Something else ???

Tasks and collaborations [2/2]

- LT group can
 - Help adapting existing NLP tools and approaches to Amharic
 - Give technical help for Master and PHD students working on Amharic NLP

Thank you