

An Amharic Speech Corpus for Large Vocabulary Continuous Speech Recognition

Solomon Teferra Abate, Wolfgang Menzel, Bairu Tafila

University of Hamburg, Faculty of Informatics
 {solomon,menzel}@nats.informatik.uni-hamburg.de

1. The Amharic Language

Amharic:

- is official language of Ethiopia.
- is a Semitic language family.
- has the largest number of speakers after Arabic - 22.5 million.
- has five dialectal variations: Addis Ababa, Gojjam, Gonder, Wollo, and Menz.
- has special features. For example: glottalized plosives ጸ, ጥ, ጸ, ጭ, and ቆ.
- has rich morphology -> many word forms.

Phonetics

Amharic has a set of 38 phones, seven vowels and thirty-one consonants.

Manner of Art/n	Voicing	Place of Articulation				
		Lab	Dent	Pal	Vel	Glo
Stops	Voiceless	ጥ[p]	ት[t]	ቸ[tʰ]	ከ[k]	አ[ʔ]
	Voiced	ብ[b]	ድ[d]	ጅ[dʒ]	ግ[g]	
	Glottalized	ጸ[pʰ]	ጥ[tʰ]	ጭ[tʰʰ]	ቆ[q]	
	Rounded				ከጎ [kʷ], ጎጎ [gʷ], ቆጎ [qʷ]	
Fricatives	Voiceless	ፍ[f]	ሰ[s]	ሸ[ʃ]		ሀ[h]
	Voiced		ዝ[z]	ሻ[ʒ]		
	Glottalized		ጸ[sʰ]			
	Rounded					ከጎ [hʷ]
Nasals	Voiced	ጦ[m]	ን[n]	ኝ[ɲ]		
Liquids	Voiced		ል[l], ር[r]			
Semi vowels	Voiced	ው[w]			ይ[j]	

Vowels

	front	center	back
high	አ[i]	እ[ɨ]	ኡ[u]
mid	ኣ[e]	አ[ə]	ኦ[o]
low		አ[a]	

The writing system

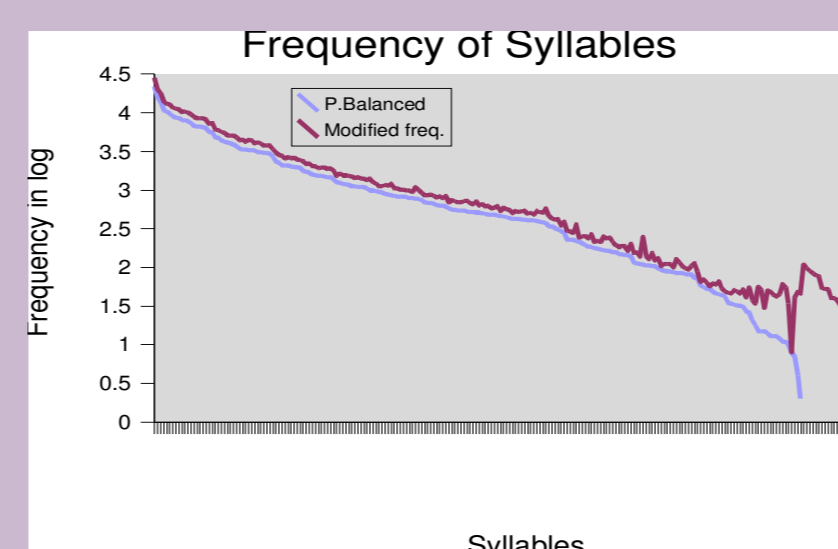
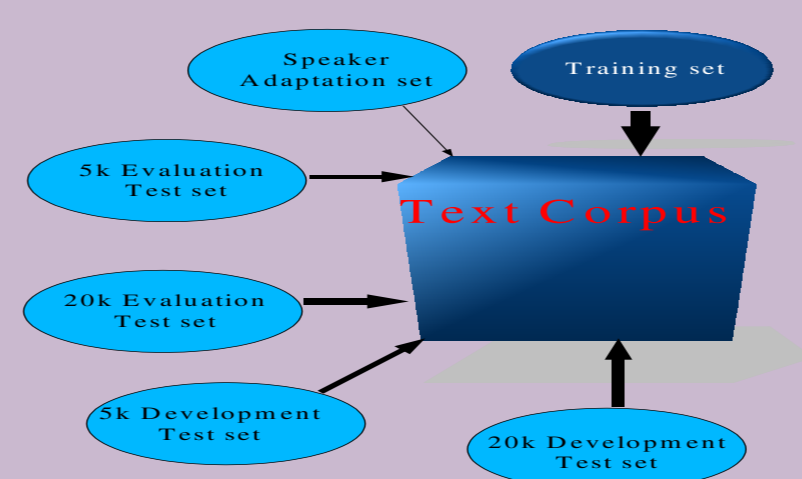
- Amharic writing system is phonetic.
- Amharic orthography is syllabic, because
 - there is, more or less, a one-to-one correspondence between the sounds of CV syllable and the graphemes.
- Amharic orthography consist of 276 symbols, including the redundant ones.
- Amharic has 234 distinct CV syllable sounds.

	Used graphemes													
	a	u	i	Ä	E	e	o	ua	ui	uA	uE	ue		
h	ሀ	ሁ	ሂ	ሃ	ሄ	ህ	ሆ	ሐ	ሑ	ሒ	ሓ	ሔ	ሕ	ሖ
l	ለ	ሉ	ሊ	ላ	ሌ	ሎ								
m	መ	ሙ	ሚ	ሚ	ሚ	ሚ	ሚ	ሚ	ሚ	ሚ	ሚ	ሚ	ሚ	ሚ
r	ረ	ሩ	ሪ	ራ	ሮ	ሮ								
s	ሰ	ሱ	ሲ	ሳ	ሴ	ስ	ሶ	ሷ						
S	ሰ	ሱ	ሲ	ሳ	ሴ	ስ	ሶ	ሷ						
^	ሰ	ሱ	ሲ	ሳ	ሴ	ስ	ሶ	ሷ						
q	ቀ	ቁ	ቂ	ቃ	ቄ	ቅ	ቆ	ቇ	ቈ	቉	ቊ	ቋ	ቌ	ቍ
b	ቦ	ቧ	ቨ	ቩ	ቪ	ቫ	ቬ	ቭ	ቮ	ቯ	ተ	ቱ	ቲ	ታ
v	ቨ	ቩ	ቪ	ቫ	ቬ	ቭ	ቮ	ቯ	ተ	ቱ	ቲ	ታ	ቴ	ት
t	ተ	ቱ	ቲ	ታ	ቴ	ት	ቶ	ቷ	ቸ	ቹ	ቺ	ቻ	ቼ	ች
c	ቸ	ቹ	ቺ	ቻ	ቼ	ች	ቾ	ቿ	ቻ	ቼ	ች	ቾ	ቿ	ቻ
n	ነ	ኑ	ኒ	ና	ኔ	ኖ	ኘ	ኙ	ኚ	ኛ	ኜ	ኝ	ኞ	ኟ
N	ነ	ኑ	ኒ	ና	ኔ	ኖ	ኘ	ኙ	ኚ	ኛ	ኜ	ኝ	ኞ	ኟ
H	አ	ኡ	ኢ	ኣ	ኤ	ኦ	ኧ	ከ	ኩ	ኪ	ካ	ኼ	ኽ	ኾ
k	ከ	ኩ	ኪ	ካ	ኼ	ኽ	ኾ	኿	ኻ	ኼ	ኽ	ኾ	኿	ኻ
K	ከ	ኩ	ኪ	ካ	ኼ	ኽ	ኾ	኿	ኻ	ኼ	ኽ	ኾ	኿	ኻ
w	ወ	ዐ	ዑ	ዒ	ዓ	ዔ	ዕ	ዖ						
z	ዘ	ዙ	ዚ	ዛ	ዞ	ዟ	ዠ	ዡ						
Z	ዘ	ዙ	ዚ	ዛ	ዞ	ዟ	ዠ	ዡ						
y	የ	ዮ	ዿ	የ	የ	የ	የ	የ						
d	ደ	ዱ	ዲ	ዳ	ዴ	ድ	ዶ	ዷ						
D	ደ	ዱ	ዲ	ዳ	ዴ	ድ	ዶ	ዷ						
g	ገ	጑	ጒ	ጓ	ጔ	ጕ	጖	጗	ጘ	ጙ	ጚ	ጛ	ጜ	ጝ
T	ተ	ቱ	ቲ	ታ	ቴ	ት	ቶ	ቷ	ቸ	ቹ	ቺ	ቻ	ቼ	ች
C	ቸ	ቹ	ቺ	ቻ	ቼ	ች	ቾ	ቿ	ቻ	ቼ	ች	ቾ	ቿ	ቻ
P	ጸ	ጹ	ጺ	ጻ	ጼ	ጽ	ጾ	ጿ						
x	ረ	ሩ	ሪ	ራ	ሮ	ሮ								
f	ፍ	ፋ	ፅ	ፆ	ፇ	ፈ	ፉ	ፊ	ፋ	ፅ	ፆ	ፇ	ፈ	ፉ
p	ጥ	ፑ	ፒ	ፓ	ፔ	ፕ	ፖ	ፑ						

2. Preparation of Amharic Speech Corpus

Text Corpus Preparation

- It is a read speech corpus.
- Text was aquired from EthioZena website.
- Preprocessing:
 - Spelling and grammar errors have been corrected;
 - Abbreviations have been expanded;
 - Foreign words have been eliminated;
 - Numbers have been textually transcribed; and
 - Concatenated words have been separated.

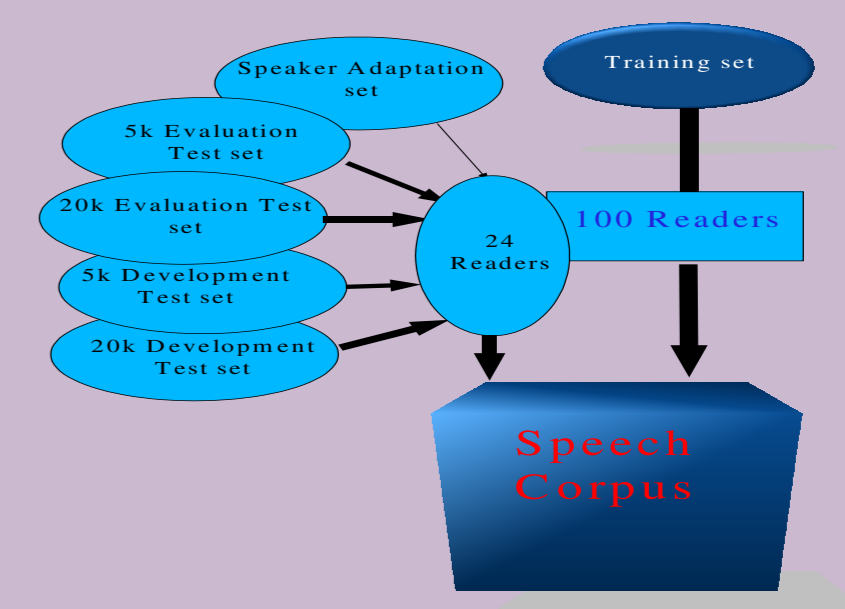


Recording of the speech

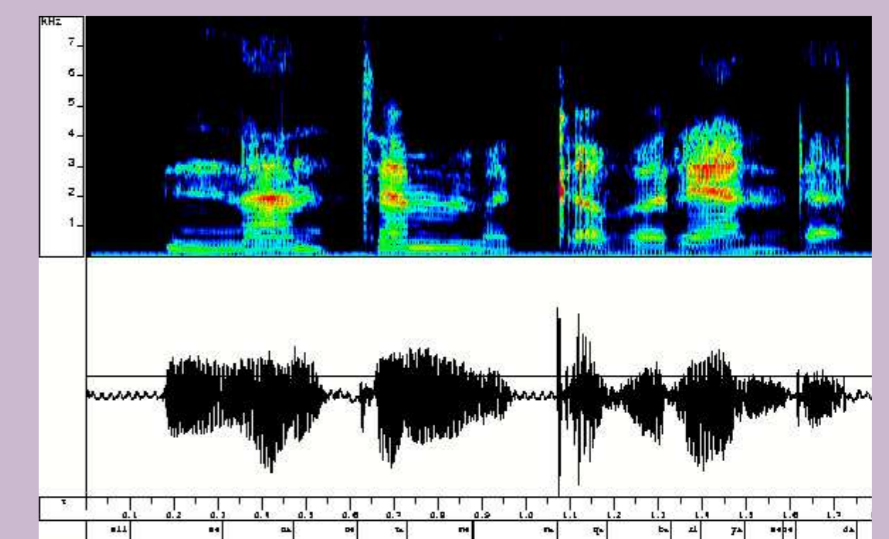
We have done the recording using a laptop and a noise canceling

microphone.

Age Range	Training set		Test sets	
	Male	Female	Male	Female
Speakers of the Addis Ababa dialect				
18-23	18	18	3	3
24-28	12	12	3	3
29-40	5	5	3	3
Older than 40	5	5	1	1
Total	40	40	10	10
Speakers of the other four dialects				
18-23	10	3		4
24-28	6	1		
Total	16	4		4
Grand Total	56	44	14	10



- Size of the corpus:
 - 20 hours of training speech
 - 100 speakers
 - 10850 sentences
 - Vocabulary of training speech is 28666
 - The training speech covers 233 of the Amharic CV syllables
- ### Segmentation of the speech
- Segmentation is done semi-automatically at word- and syllable-level.



3. Availability

- Used for our research of developing automatic speech recognizer for Amharic.
- Part of it is used by students of Addis Ababa University.
- It will be made available through a third party.

4. Acknowledgments

Daniel Yacob who made the archive of EthioZena available to us. He also provided his SERA to Ethiop (namely g2) conversion tool. DAAD (Deutscher Akademischer Austauschdienst) financed the project.

God bless you