



## D6.2 & D6.3: FINAL USABILITY REPORT AND SYSTEM EVALUATION

Grant Agreement number	7F14047
Project acronym	HaBiT
Project title	Harvesting big text data for under-resourced languages
Deliverable title	<b>D6.2 &amp; 6.3: Final Usability Report and System Evaluation</b>
Responsible partner	<b>Janne Bondi Johannessen, UiO and Björn Gambäck, NTNU</b>
Dissemination level	Public
Due delivery date	30 April 2017 (+ 30 days)
Actual delivery date	1 June 2017

Principal Investigator	<b>Karel Pala</b>
Project Promoter	<b>Masaryk University</b>
Tel	+420 549 49 5616
E-mail	<b>pala@fi.muni.cz</b>
Project website address	www.habit-project.eu



## Table of Contents

1. Introduction .....	5
2. Norwegian .....	6
200 most frequent words from the wordlist .....	6
Norwegian Bokmål .....	6
Norwegian Nynorsk .....	8
Language recognition .....	9
Random concordance pages .....	10
Norwegian Bokmål .....	10
Norwegian Nynorsk .....	10
Word sketches .....	11
Norwegian Bokmål .....	11
Thesaurus .....	12
Norwegian Bokmål .....	12
3. Oromo.....	13
Most frequent words.....	13
Concordances .....	14
Word sketches .....	15
4. Amharic .....	16
Word frequency list.....	16
Concordances .....	16
Word sketches .....	17
Thesaurus .....	20
5. Somali.....	22
Word frequency list.....	22
Concordances .....	22
Word sketches .....	23
6. Tigrinya .....	24
Word frequency list.....	24
Concordances .....	25
Word sketches .....	26
7. Conclusion.....	27
References .....	28



## 1. Introduction

There are several challenges with respect to the corpus resources developed in the HaBiT project – Harvesting big text data for under-resourced languages) financed by EEA/Norway Grants), in cooperation with the project Linguistic Capacity Building — Tools for the inclusive development of Ethiopia, financed by the Norwegian NORAD, Norwegian Agency for Development Cooperation.

This evaluation assesses usefulness and the usability of the resources and tools developed for Norwegian (both Bokmål and Nynorsk) and all the Ethiopian languages in the project: Amharic, Oromo, Tigrinya and Somali. The corpora are very different in size given the amount of material that is actually available on the web. The types of text also vary considerably. For some of the languages a problem in the corpus creation process was found already in the harvesting process, given that there are varieties that are very close to each other linguistically. For Norwegian this is true of the two varieties, Bokmål and Nynorsk, with the additional problem that Danish is also a variant that is easily confused with Bokmål. Further, Amharic can be confused with Tigrinya and with Ge'ez (the older clerical language still used by the Ethiopian church).

The evaluations reported here were performed on the versions of the resources and tools towards the end of the project. An initial evaluation of the usability of the tools and the reliability of the produced resources was also conducted during a workshop in Brno in February 2017, with participants from the teams at Addis Ababa University, Ethiopia, as well as the Norwegian groups from University of Oslo and NTNU.

This deliverable evaluates the resources and tools for each language in turn, and concludes with an overall summary in a table.

## 2. Norwegian

*Evaluated by: Janne Bondi Johannessen (UiO)*

For Norwegian, the most challenging task must be the *language recognition part*. In order to make internet corpora, the texts on the internet must be recognised as being from the intended language. Norwegian Bokmål and Nynorsk are very similar to each other, and Bokmål is also very similar to Danish. In addition, old Norwegian texts were written in Danish. A corpus that contains large chunks of texts from the other languages must be considered unsuccessful. This is therefore something that must be checked.

### 200 most frequent words from the wordlist

#### Norwegian Bokmål

[http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=notenten15\\_4\\_bokmal&refs=%3Dword&wlmaxitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include\\_nonwords=1&wlnums=frq&wltype=simple](http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=notenten15_4_bokmal&refs=%3Dword&wlmaxitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include_nonwords=1&wlnums=frq&wltype=simple)

We have evaluated the Habit Bokmål frequency list by comparing it with the NoWaC frequency list (available at the Text Laboratory at UiO, see Guevara 2010), see Table 1.

As Table 1 shows, the most frequent words are more or less the same in both lists, with pronouns and grammatical words topping the list: *og* 'and', *i* 'in', *er* 'is', *å* 'to' (infinitive marker), *på* 'on', *det* 'it', *som* 'which', *en* 'a', *til* 'to', *av* 'of' etc. These results are expected and shows that the method used has been successful. The fact that each word has approximately the double frequency of that of the existing list shows that the corpus is much bigger (double) than the existing one, which is a very welcome result. The Habit Bokmål corpus has 1,178,357,993 words.

Table 1: Comparison of the Bokmål frequency lists from HaBiT and NoWaC

Habit project Bokmål list		NoWaC frequency list
<u>word</u>	<u>Freq</u>	
.	<a href="#">66,302,248</a>	16145498 og
,	<a href="#">48,783,706</a>	15256386 i
og	<a href="#">38,459,956</a>	12206184 er
i	<a href="#">29,244,572</a>	11698509 det
er	<a href="#">26,587,096</a>	9571416 på
å	<a href="#">20,257,955</a>	8630777 som
på	<a href="#">20,092,511</a>	8293212 en
det	<a href="#">19,741,244</a>	8168237 til
som	<a href="#">19,704,300</a>	8131491 å
en	<a href="#">19,532,414</a>	7993433 av
til	<a href="#">17,771,732</a>	7491272 for
av	<a href="#">17,219,280</a>	6670288 med
for	<a href="#">15,515,275</a>	6639161 :
med	<a href="#">15,026,187</a>	6531429 -
har	<a href="#">12,734,327</a>	6301028 har
at	<a href="#">11,518,852</a>	6030608 at
ikke	<a href="#">10,277,938</a>	5235776 ikke
jeg	<a href="#">9,448,321</a>	5073332 jeg
de	<a href="#">8,835,895</a>	4612638 de
den	<a href="#">7,746,306</a>	4146748 )
kan	<a href="#">7,156,786</a>	4020263 den
du	<a href="#">7,094,265</a>	3980207 om
et	<a href="#">6,959,830</a>	3678914 (
om	<a href="#">6,950,325</a>	3239821 "
var	<a href="#">6,057,924</a>	3199327 et
så	<a href="#">5,799,457</a>	3102444 men
Det	<a href="#">5,634,376</a>	2935861 vi
		2879552 kan
		2874527 du
		2747899 fra
		2632350 så
		2471901 var
		2375553 ?
		2203208 ...
		2104178 dette
		2083048 han
		2051824 seg
		1858213 skal
		1848715 vil
		1818427 !
		1710749 eller
		1690743 også
		1547255 mer
		1508222 _
		1410569 man
		1376300 ut
		1349778 ble
		1331468 da
		1326745 være

## Norwegian Nynorsk

[http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=notenten15\\_4\\_nynorsk&refs=%3Dword&wlmaxitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include\\_nonwords=1&wlnums=frq&wltype=simple](http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=notenten15_4_nynorsk&refs=%3Dword&wlmaxitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include_nonwords=1&wlnums=frq&wltype=simple)

Habit Nynorsk frequency list	
<u>word</u>	<u>Freq</u>
.	<a href="#">3,277,542</a>
,	<a href="#">2,329,138</a>
og	<a href="#">1,967,701</a>
i	<a href="#">1,709,665</a>
er	<a href="#">1,136,473</a>
på	<a href="#">1,007,837</a>
det	<a href="#">943,786</a>
til	<a href="#">916,306</a>
som	<a href="#">905,752</a>
å	<a href="#">747,507</a>
av	<a href="#">736,146</a>
med	<a href="#">730,444</a>
for	<a href="#">683,801</a>
ein	<a href="#">665,297</a>
har	<a href="#">587,469</a>
at	<a href="#">512,242</a>
dei	<a href="#">460,646</a>
var	<a href="#">424,225</a>
ikkje	<a href="#">412,301</a>
eg	<a href="#">392,888</a>
om	<a href="#">362,124</a>
den	<a href="#">321,121</a>
frå	<a href="#">320,828</a>
eit	<a href="#">317,410</a>
Det	<a href="#">289,144</a>
)	<a href="#">278,553</a>
:	<a href="#">271,313</a>
ei	<a href="#">246,774</a>
han	<a href="#">245,719</a>



kan	<a href="#">243,963</a>
(	<a href="#">233,281</a>
så	<a href="#">221,023</a>
seg	<a href="#">220,810</a>
men	<a href="#">214,5</a>

The first part of the Nynorsk list shows the same expected words as the Bokmål one, again indicating that the method is right. There is a lot less written in this variety of Norwegian than in the other, and therefore there is also a lot less in the web. Thus, while there are 38 million occurrences of the most frequent word, the conjunction *og* ('and' in Bokmål), there are only around 2 million occurrences of this word in the Nynorsk corpus. Some of the most frequent words have different form in Bokmål and Nynorsk, and we find them in the same position in the respective frequency lists, which is a good sign. For example, the article *ein* ('a') is number 12 in the Nynorsk list (with more than 700 000 occurrences) and the equivalent Bokmål *en* ('a') number 8 in the Bokmål list (with nearly 20 million occurrences). The Nynorsk version *ikkje* ('not') is number 17 in its list (with more than 400 000 occurrences), while the Bokmål negation *ikke* ('not') is number 15 in its list (with more than 10 million occurrences). The similarities between the position of the frequent words on the lists therefore indicate that the lists and therefore these two corpora are good reflections of the two languages.

### Language recognition

While the language recognition of Bokmål and Nynorsk look right when considering the top positions of the frequency lists, we need to check the corpora for the "wrong" words. Thus, we do not want to find a high number of Nynorsk words in the Bokmål corpus or vice versa. Neither do we want to find Danish words in the Bokmål corpus.

Looking for the Bokmål *ikke* ('not') in the Nynorsk corpus reveals that there are more than 4000 occurrences of this word. But given that there are more than 400 000 of the correct equivalent *ikkje* ('not'), we can accept that the Nynorsk corpus has 1 % Bokmål versions of this word. We also find nearly 20 000 occurrences of the Bokmål *en* ('a'), but this is still a small number compared with the Nynorsk *ein* ('a') of more than 700 000. Given that *en* ('a') in Bokmål to some extent replaces the feminine article *ei* 'a' as well, we must also look for this feminine article in the Nynorsk corpus, and find 263 000 occurrences. Thus the Bokmål article in the Nynorsk corpus makes up 20 000 compared with 963 000 equivalent Nynorsk articles, which is a low number, i.e. ca 2 %.

One might wonder why there are Nynorsk words in the Bokmål corpus at all. The reason is linked to the fact that the Nynorsk language norm is in minority in Norway, so that many of those that write Nynorsk also use Bokmål words in their Nynorsk texts as a result of influence by the majority. The texts show this clearly: Many are from informal blogs with a mix between the two norms.

Looking for Nynorsk words in the Bokmål corpus gives the following results. The nynorsk negation *ikkje* ('not') occurs nearly 38 000 times, but given that there are more than

10 million occurrences of the Bokmål *ikke* ('not'), only 0,4 % are Nynorsk. There are 31 000 occurrences of the Nynorsk *ein* ('a') in the Bokmål corpus, which makes up 0.1 % of the 31 million occurrences of the correct Bokmål occurrences of *en* ('a'). Again, we find that blogs are the cause of many of the occurrences from the other norm. The result is therefore very good.

We must also check the results compared with Danish. Searching for the Danish *næsten* ('almost') in the Bokmål corpus gives under 300 hits, as against the Bokmål *nesten* ('almost'), which occurs nearly 500 000 times. There are, then, 0.15 % of this Danish word in the corpus. There are less than 100 occurrences of the Danish *fuld* ('full') in the corpus, as against 420 000 occurrences of the Bokmål *full* ('full'), giving 0.0025 %. The amount of Danish is negligible. The few occurrences that do occur are found in texts like blogs and film titles.

It can be concluded that the language recognition module works well and has produced good quality web corpora.

### Random concordance pages

To check the concordance production, concordances for at least five words or phrases were generated for each language.

#### Norwegian Bokmål

[http://corpora.fi.muni.cz/habit/run.cgi/first\\_form?corpname=notenten15\\_4\\_bokmal](http://corpora.fi.muni.cz/habit/run.cgi/first_form?corpname=notenten15_4_bokmal)

#### **arbeide 'work', as a verb:**

More than 10 000 hits, that look right, though there are some instances of wrongly tagged words, so that there are some nouns in between the verbs.

#### **arbeid 'work', as a noun:**

More than 470 000 hits, that look right and rightly tagged.

#### **snill 'kind', as an adjective:**

More than 65 000 hits. These include many inflectional forms: *snill*, *snille*, *snilt*, *snilleste*, *snillere*.

#### **glad for 'happy about', i.e., a phrase:**

More than 54 000 hits that are good.

#### **tsjekkisk 'Czech', as a lemma:**

Nearly 4000 hits that look right. Both inflectional forms: *tsjekkisk*, *tsjekkiske*.

#### Norwegian Nynorsk

[http://corpora.fi.muni.cz/habit/run.cgi/first\\_form?corpname=notenten15\\_4\\_nynorsk](http://corpora.fi.muni.cz/habit/run.cgi/first_form?corpname=notenten15_4_nynorsk)

The Nynorsk corpus is not tagged.

#### **arbeide 'work':**

Over 5000 hits. Only this word form, since the corpus is not tagged.

#### **arbeid 'work':**

20 000 hits that look right. Only this word form.

#### **snill 'kind':**

1300 hits, only this form.

***glad for 'happy about':***

2700 hits. Only these word forms.

***tsjekkisk 'Czech'***

200 hits. Only this word form.

**Word sketches**

Word sketches are one page summaries of collocations and grammatical behaviour of words. To test these, word sketches for at least five frequent words (varying in the part of speech) for each language were generated (names have to be lowercased, e.g., “norge”, “oslo”)

**Norwegian Bokmål**

[http://corpora.fi.muni.cz/habit/run.cgi/wsketch\\_form?corpname=notenten15\\_4\\_bokmal](http://corpora.fi.muni.cz/habit/run.cgi/wsketch_form?corpname=notenten15_4_bokmal)

**tegne 'draw', word sketch:**

The **modifiers** (postpost) are very good: *dyster* 'gloomy', *deretter* 'then', *ferdig* 'finished', *først* 'first'

The **objects** are good: *forsikring* 'insurance', *membership* 'medlemskap', *portrett* 'portrait',: *tegne* -> psykolog,

The subjects are mostly, but not always good. Sometimes the head of a relative clause is misinterpreted as a subject:

en **tegneserie tegnet** for hånd = en cartoon drawn by ...

The **coordination** (tegne and/or) are very good: write, paint, explain, ...

The **pronominal** subjects are good.

**synges 'sing', verb, word sketch:**

The (postposed) **modifiers** are good: choir, solo, live, together...

The **objects** are good: song, duet, birthday song, melody, opera...

The **subjects** are good: bird, choir, artist, ...

**Coordination** and/or is good: dance, play, speak...

**Pronominal** subjects are good

**elev 'pupil', noun, word sketch:**

**Modifiers** are good: single, minority-linguistic, good...

**Verbs** with *elev* as object are good: teach, help, let, motivate

**Verbs** with *elev* as subject are good: learn, participate, work, experience, make...

**Coordination** and/or is good: teacher, parent, student...

Elev **is**: eager, motivated, engaged, good...

**lat 'lazy', adjective, word sketch:**

**Modifiers** are good: too, intellectually, quite, very...

**Coordination** and/or is good: fat, useless, spoilt, stupid

**Subjects** of 'be lat' are good: youth, student, Norwegian...

**Verbs** modified by 'lat' are slightly strange: understand, call, film, feel...

**Nouns** modified by 'lat' are good: summer's day, holiday, morning...

**blå 'blue', adjective, word sketch:**

**Modifiers** are good: dark, glowing, light...

**Coordination** and/or is good: white, green, yellow, pink...

**Subjects** of 'be blå' are good: sky, colour, sea...

**Nouns** modified by 'blå' are good: sky, eye, sea...

**Verbs** modified by 'blå' are slightly strange: lighten, colour, unite...

## Thesaurus

### Norwegian Bokmål

[http://corpora.fi.muni.cz/habit/run.cgi/thes\\_form?corpname=notenten15\\_4\\_bokmal](http://corpora.fi.muni.cz/habit/run.cgi/thes_form?corpname=notenten15_4_bokmal)

<b>glad</b> ( <i>adjective</i> ) Norwegian Web 2015 (B)		
Lemma	Score	Freq
<a href="#">flink</a>	0.474	137,058
<a href="#">opptatt</a>	0.424	126,574
<a href="#">lykkelig</a>	0.422	47,064
<a href="#">stolt</a>	0.417	78,073
<a href="#">redd</a>	0.410	119,413
<a href="#">snill</a>	0.406	66,585
<a href="#">fornøyd</a>	0.396	71,989
<a href="#">hyggelig</a>	0.394	144,736
<a href="#">ung</a>	0.392	338,015
<a href="#">god</a>	0.390	4,549,101
<a href="#">veldig</a>	0.390	1,095,565
<a href="#">ivrig</a>	0.389	41,774
<a href="#">morsom</a>	0.376	167,941
<a href="#">rolig</a>	0.374	117,697
<a href="#">aktiv</a>	0.373	193,196
<a href="#">klar</a>	0.372	628,218
<a href="#">frisk</a>	0.363	148,294
<a href="#">fin</a>	0.358	800,680
<a href="#">heldig</a>	0.357	94,030
<a href="#">syk</a>	0.357	142,657
<a href="#">søt</a>	0.356	104,022
<a href="#">mye</a>	0.354	5,115,406
<a href="#">sikker</a>	0.349	489,864
<a href="#">positiv</a>	0.346	220,085



The thesaurus is good for finding alternative words, but the words are not necessarily synonyms. For example the adjective *glad* 'happy, glad' have as the most frequent suggestions: clever, busy, happy, proud, afraid, kind, contented. Another example is *skrive* 'write', with the "synonyms": read, tell, mean, see, say, speak.

Overall, the Norwegian resources have high quality.

### 3. Oromo

*Evaluated by: Fede Negesse (AAU)*

[http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=orwac16&refs=%3Dword&wlmitems=200&wlsort=f&wlatr=word&wlminfreq=5&wlmaxfreq=0&include\\_nonwords=1&wlnums=frq&wltype=simple](http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=orwac16&refs=%3Dword&wlmitems=200&wlsort=f&wlatr=word&wlminfreq=5&wlmaxfreq=0&include_nonwords=1&wlnums=frq&wltype=simple)

#### Most frequent words

After checking, the most frequent words in the web corpus of Oromo, were identified as being words that are mostly function words which exist in the language. The results make sense, but there are bound morphemes ( hin-, ni-, nu-, haa-, -e, -a, -u) that were separated from other morphemes, likely during tokenization.

Examples:

hin- hin-deemu ( I/he won't go)  
ni- ni-deemu (They will go)  
nu- nu-beeku( They know us)  
haa- haa-deemnu( Let's go)  
-e/-a/-u haa-taa'-u ( Let him sit )  
taa'-e ( He/ sat)  
taa'-a ( Sit or take a chair)

-ti ( copula) Kun mana kee-ti ( This is your house)  
in- ( allomorph of /ni/ ) An in-deema (I will go) for some dialects of Oromo.

The wrong tokenization has to do likely with the glottal stop /ʔ/, orthographically represented by an apostrophe '/'. The problem is significantly reduced now. It was rampant before the evaluation workshop in Brno.

## Concordances

An example of the concordance is given below.

Query **hin, deemu** 29 (5.70 per million)

Page 1 of 2 Go Next Last

<a href="#">bilisummaa.com</a>	bara keessa jirru kana waliinis gongumaa <b>hin deemu</b> . </p><p> Odoo seerri humnaan amantii ofii
<a href="#">omicwomo.info</a>	ittiin yaalamus argachaas hin jiru argatuufis <b>hin deemu</b> ka nama sodaachisu mootummaan warra du'a
<a href="#">finfinnetribune.com</a>	hanqinni kun kan Hogganaatii? kana irra <b>hin deemu</b> . </p><p> Yaadoota kanan kaasuuf, Ani Hoggana
<a href="#">ayyaantuu.net</a>	kee lameenuu, mata duree kenniteef waliin <b>hin deemu</b> . Tasuma barreeffama kee ijaa fi ijoo hin
<a href="#">oromedia.net</a>	otoo guyya sadan tokko waliin hidabaarsin <b>hin deemu</b> .Baqalaa Nadhii anaf obboleecha jechuu dha
<a href="#">qeerroo.org</a>	kenyaa hamma Wayyaaneen jiru hiikkachuuf <b>hin deemu</b> ,kanaaf manni Hidhaa keessii fi alli isaa
<a href="#">qeerroo.org</a>	qabeenyaa habashaa kun naannoo kana gadhiisee <b>hin deemu</b> taanaan fuuldura lubbuu isaa sodaachisaa
<a href="#">gubirmans.com</a>	kaanne qaba. Bilisummaa sabicha irra fagaatee <b>hin deemu</b> . Kan hojiin farraa ta'e illee afaaniin
<a href="#">blogspot.cz</a>	akka danda'amu mul'isaa jirti. Kana irra <b>hin deemu</b> . Shira diinaa kana fashalsuuf waan hojjachuu
<a href="#">gadaa.net</a>	qawweedhaan dhufe, humnaan malee namarraa <b>hin deemu</b> . Kana isaanuu ifaa ifatti dubbachaa turan
<a href="#">voaafaanoromoo.com</a>	dhihaatanii, akkasitti immoo haasaan nagaa <b>hin deemu</b> . . . " jedhan. </p><p> sagataan keessan
<a href="#">jw.org</a>	ifni dungichaa dhaamu ifni sun eessayyuu <b>hin deemu</b> . Salphaadhumatti dungoon sun ifa kennuusaa
<a href="#">ulfo.org</a>	ergarama mana hidhaatti irra geyerraa haasawuuf <b>hin deemu</b> . Mana hidhaa keeysatti addaachuu, daaruu
<a href="#">oromoliberationfront.org</a>	kophee malee deema; garuu mataa gadi qabatee <b>hin deemu</b> , jedhe. </p><p> Yeroo sanatti qabsoo bilisummaa
<a href="#">voaafaanoromoo.com</a>	dhabsiisan. Walitti bu'insi sunis wal waraansatti <b>hin deemu</b> jedhan. </p><p> Yunivarsitii Jimmaa fi Mormii
<a href="#">gadaa.com</a>	hari'uun Ulee lafa hin godhu. Harka duwwaa <b>hin deemu</b> . Saroonni yoo itti dhufte qarrifaa ishii
<a href="#">oromiatimes.org</a>	waraana naquun maal nu akeeka ??? kana irra <b>hin deemu</b> . garuu waan hamileen yaadinu fi saganteeffannee
<a href="#">gadaa.net</a>	akkuma yaadannu warreen amantaa keenya waliin <b>hin deemu</b> jedhaa turan ni jirani. Isaan kunneen maaliif
<a href="#">readbag.com</a>	nam tokkeenis wal abaaruun furmaata ta'uuf <b>hin deemu</b> . Hegereen Oromoo bareeduuf, ar'a hujii
<a href="#">gadaa.com</a>	garboomsaa kamiyyuu waliin araara uumuuf <b>hin deemu</b> ; araarri kan dhufu, mirgii fi dantaan ummata

Page 1 of 2 Go Next Last

## Word sketches

[http://corpora.fi.muni.cz/habit/run.cgi/first\\_form?corpname=orwac16](http://corpora.fi.muni.cz/habit/run.cgi/first_form?corpname=orwac16)

The search for words or phrases generally produces good results.

The hits from a search for word sketches is illustrated below.

**deemu** Oromo WaC [2016] freq = **696** (136.69 per million)

modifiers of "deemu"	356	1.00	verbs with "deemu" as object	35	0.90	nouns modified by "deemu"	429	1.30
khana	<u>3</u>	8.08	jechuu	<u>2</u>	4.19	ta'uuf	<u>8</u>	8.44
taasisa	<u>3</u>	7.57	qabnu	<u>2</u>	4.19	guddachaa	<u>6</u>	8.40
sakatta'aa	<u>2</u>	7.48	qaba	<u>5</u>	4.03	dabalaa	<u>6</u>	8.30
alaalatti	<u>2</u>	7.48	yommuu	<u>4</u>	3.97	tahuuf	<u>4</u>	8.04
mara	<u>5</u>	7.45	kanaan	<u>4</u>	3.86	dhihaachaa	<u>4</u>	8.00
ammatti	<u>2</u>	7.31	gochuu	<u>2</u>	3.79	jabaachaa	<u>4</u>	7.85
rasaasaan	<u>2</u>	7.14	jira	<u>5</u>	3.00	baadhatee	<u>3</u>	7.82
hundatti	<u>3</u>	7.08	jechuun	<u>2</u>	2.79	dhalachuuf	<u>3</u>	7.78
gochuuf	<u>5</u>	6.98	qabu	<u>2</u>	2.13	dhufuuf	<u>4</u>	7.77
ennaa	<u>3</u>	6.95				garamitti	<u>3</u>	7.77
eessa	<u>2</u>	6.85	verbs with "deemu" as subject	<u>19</u>	<b>1.00</b>	ennaa	<u>6</u>	7.75
godhanii	<u>2</u>	6.68	qaba	<u>4</u>	4.09	karaarra	<u>3</u>	7.70
godha	<u>3</u>	6.52	jira	<u>8</u>	3.95	fudhatanii	<u>4</u>	7.56
qeerroon	<u>4</u>	6.46	wal	<u>2</u>	3.17	yommuun	<u>3</u>	7.53
qabanii	<u>2</u>	6.44				galuuf	<u>3</u>	7.51
daandii	<u>2</u>	6.24	"deemu" is ...	<u>3</u>	<b>3.90</b>	fuudhee	<u>3</u>	7.40
godhe	<u>2</u>	6.06	kanatti	<u>3</u>	4.96	dabree	<u>3</u>	7.34
sababa	<u>2</u>	5.86	"deemu" is a ...	<u>42</u>	<b>1.40</b>	baatee	<u>3</u>	7.29
tahuu	<u>3</u>	5.81	alaalatti	<u>2</u>	10.02	kheessa	<u>2</u>	7.24
turan	<u>8</u>	5.05	dhugumatti	<u>2</u>	9.19	qottoon	<u>2</u>	7.24
waliin	<u>4</u>	4.84	ammatti	<u>2</u>	8.71	babaldhachaa	<u>2</u>	7.23
kanatti	<u>3</u>	4.80	hundatti	<u>3</u>	7.11	diigamuuf	<u>2</u>	7.23
danda'	<u>2</u>	4.65	godha	<u>3</u>	6.57	oolfamuuf	<u>2</u>	7.23
humna	<u>2</u>	4.47				babal'ataa	<u>2</u>	7.21
nama	<u>7</u>	4.25				fudhatamuuf	<u>2</u>	7.20

## 4. Amharic

*Evaluated by Derib Abo (AAU)*

### Word frequency list

[http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=amwac16&refs=%3Dword&wlmaxitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include\\_nonwords=1&wlnums=frq&wltype=simple](http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=amwac16&refs=%3Dword&wlmaxitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include_nonwords=1&wlnums=frq&wltype=simple)

There are several punctuation marks that are generated as words. Three different letters that should have been written together with a word are presented separately and a symbol for addition is also generated as the most frequent word.

The words that are generated as most frequent are likely to be most frequent specially based on the type of data collected. On the top are the stop words and the verb ነፃ, which is likely to be the case for any corpus on Amharic. Since online sources in Amharic have political and religious content, the content words are from these two domains and related subjects.

### Concordances

[http://corpora.fi.muni.cz/habit/run.cgi/first\\_form?corpname=amwac16](http://corpora.fi.muni.cz/habit/run.cgi/first_form?corpname=amwac16)

The concordance is very good. I have tried to display concordance for the following words: መንግስት 'government', ቤት 'house, home', አንዱ 'the one', ሳይንስ 'science', የተለያዩ 'different'.

All sentences on the first page of the concordance are correct.



## Word sketches

[http://corpora.fi.muni.cz/habit/run.cgi/wsketch\\_form?corpname=amwac16](http://corpora.fi.muni.cz/habit/run.cgi/wsketch_form?corpname=amwac16)

Generally, the modifiers (modifiers and modified by) and the and/or sketches are good. The *subject* sketches and the *verbs with* sketches seem to be a bit problematic because Amharic sentences are complex sentences with several verbs in one sentence. There are spelling errors in the data by the writers and that has contributed its share, specially for the word ሰራ '[he] worked'. Compound nouns could be treated as modifiers sometimes when they do not refer to the specific proper or common noun.

Word sketch of a noun:

ሰው 'man, human being'

Amharic WaC [2013 + 2015 + 2016] freq = **41,861** (2,063.41 per million)

modifiers of "ሰው"	21,181	1.20	verbs with "ሰው" as	39,971	2.10	"ሰው" and/or ...	6630.80	nouns modified by "ሰው"	19,963	1.20
አንድ 'one'	4,345	11.15	subject			ድርጅት 'organization'	178.89	ሰራ ሽ A		
ማንም 'anyone'	490	9.41	ይችላል 'can be'	441	7.70	አንድ 'one'	208.68	compound		
አንድን 'one, object'	251	8.49	ሆኖ 'being'	469	7.68	እንስሳ 'animal'	88.55	noun with ሰው to mean man made	191	8.27
የአንድ 'that of one'	235	8.31	አንደሆነ 'that being'	372	7.12	ተቋም	88.40	ሆኖ 'being'	263	8.19
በአንድ 'by one'	253	8.06	የለም 'not'	227	7.10	'institution'	118.35	" punctuation	353	7.44
ዓይነት 'type, kind'	232	7.91	ብሎ 'he having'	319	7.01	ነገር 'thing'	98.34	" punctuation	124	7.12
ሌላ 'for other'	126	7.56	said			አካል 'body'	98.24	ነኝ 'I am'	98	7.12
ከአንድ 'from one'	140	7.53	አልነበረም 'it was not'	188	7.01	"	75.67	* symbol	100	7.09
ብዙ ,many'	201	7.46	አይደለም 'it is not'	278	6.96			ሆኖ 'that he is'	133	7.03
ታላቅ 'great'	134	7.45	ሌላ 'he said', 'he is present'	236	6.92			በሞት 'by death'	83	7.01
የአንድን 'that of one'	113	7.41	አይችልም 'it cannot'					የለም 'there is not'	73	6.74
ለአንድ 'for one'	117	7.35	ባይሆን 'it cannot'	171	6.86			" punctuation	221	6.68
ፍጹም 'perfect'	106	7.28	ባይሆን 'it cannot'	200	6.83			በሆኖ 'by being'	71	6.59
አኗሪ	102	7.27	ባይሆን	265	6.79			አንድ 'one'	117	6.57
አንድም	101	7.22	ከሆነ	185	6.65			ሕይወት 'life'	73	6.56
የመጀመሪያው	104	7.16	ይሆናል	202	6.62			ማን 'who'	65	6.52
ከፋ	103	7.16	ሊሆን	170	6.60			ልብ 'heart'	63	6.44
የሚባል	102	7.15	ነበር	143	6.44			ከእግዚአብሔር	56	6.44
ሌላ	100	7.13	ሆነ	143	6.37			'from God'		
አምላክ	119	6.96								
"	251	6.88								

## Word sketch of a verb:

ሰራ 'worked'

The ones highlighted in yellow are the result of the spelling error that the writers made. They wrote ሰራ 'worked, verb' in stead of ሰራ 'work, noun'

**Amharic WaC [2013 + 2015 + 2016] freq = 194** (9.56 per million)

<u>modifiers of "ሰራ"</u>	<u>108</u>	<u>1.20</u>		
በእርሻ 'in farming'	<u>2</u>	8.16	<u>nouns modified by "ሰራ"</u> <u>65</u>	<u>0.80</u>
ሁብረት 'in unions'	<u>2</u>	5.70	ፈጣሪ 'creator, God'	<u>2</u> 6.00
ታሪክ 'history'	<u>6</u>	5.32	አስኪያጅ 'excutive'	<u>2</u> 5.03
የልማት 'developmental'	<u>4</u>	4.66	.. punctuation	<u>2</u> 4.39
ሰራ 'work, Noun'	<u>3</u>	4.14	? punctuation	<u>12</u> 3.81
አይነት 'kind, type'	<u>2</u>	3.47		
ጥሩ 'good'	<u>2</u>	3.38		
ነገር 'thing'	<u>5</u>	3.30		
ሥራ , 'work, Noun'	<u>3</u>	2.87		
ከፍተኛ 'higher'	<u>2</u>	2.14		

### subjects of "ሰራ" 66 0.70

በሌሎች 'by'		
This is instrumental, not subject others'	<u>2</u>	5.42
ወገኖች 'sides'		
This is instrumental, not subject	<u>3</u>	4.75
ሰራ ;work, Noun'	<u>4</u>	3.56
		0.16

### verbs with "ሰራ" as subject 161 1.60

ለተሰማሩ ' for those working in..'	<u>2</u>	7.37
ቢባል ' if [it is ] said so'	<u>3</u>	6.15
የተሰማሩ 'those working in ..'	<u>2</u>	5.31
የሚሰሩ 'those who work'	<u>2</u>	5.16
ቢያንስ 'at least'	<u>2</u>	4.45
አይደሉም 'they are not'	<u>2</u>	3.72
ማለት 'means'	<u>3</u>	3.56
የሚለው 'means'	<u>2</u>	1.93
ተብሎ 'alled, referred to'	<u>2</u>	1.73

ጊዜ 'time' This is  
an adverbial, not  
a subject

Word sketch of an adjective:

አዲስ 'new'

**Amharic WaC [2013 + 2015 + 2016] freq = 20,015 (986.58 per million)**

**modifiers of "አዲስ" 393 0.10**

በሽራተን compound noun with አዲስ meaning Sheraton Addis	<u>11</u>	9.69
ትናንት 'yesterday', adverbial	<u>21</u>	7.94
ዘንድሮ 'this year', adverbial	<u>7</u>	7.30
ሰዎች 'these days', adverbial	<u>19</u>	7.28
በተለይ 'specially'	<u>24</u>	6.68
በተለይ 'specially'	<u>31</u>	6.57
ዛሬ 'today', adverbial	<u>52</u>	6.53
ገና 'not yet'	<u>18</u>	6.51
የ I should be attached to a word, cannot stand alone.	<u>17</u>	6.22
ነገ 'tomorrow'	<u>9</u>	6.22
ሁልጊዜ 'every time', 'always'	<u>5</u>	5.84
ዛሬም 'even today', adverbial	<u>5</u>	5.59
አሁን 'now', adverbial	<u>29</u>	5.49
እንደገና	<u>6</u>	5.29
ለምን	<u>5</u>	3.84
አቶ 'Mr.' title of a person.		
አዲስ can be the name of a person.	<u>8</u>	3.67

**nouns modified by "አዲስ" 12,916 4.70**

<u>አበበ</u> 'compound noun with አዲስ meaning Addis Ababa	<u>4,285</u>	12.46
<u>ዘመን</u> 'time'	<u>648</u>	9.68
<u>አበበን</u> "compound noun with አዲስ meaning Addis Ababa	<u>264</u>	9.34
<u>ነገር</u> 'thing'	<u>479</u>	8.83
<u>ዓለም</u> 'compound noun with አዲስ meaning Addis Alem	<u>288</u>	8.78
<u>ምዕራፍ</u> 'chapter'	<u>122</u>	8.08
<u>ጊዜጣ</u> 'newspaper'	<u>100</u>	7.67
<u>ትርጉም</u> 'interpretation'	<u>101</u>	7.59
<u>ከተማ</u> 'town', also 'compound noun with አዲስ meaning Addis keteme, one of the districts of Addis Ababa	<u>133</u>	7.54
ኪዳን 'testament', 'compound		

## Thesaurus

[http://corpora.fi.muni.cz/habit/run.cgi/thes\\_form?corpname=amwac16](http://corpora.fi.muni.cz/habit/run.cgi/thes_form?corpname=amwac16)

Though there are good similar words in the thesaurus, the majority of them are just modifiers or collocations with the word presented.

Example መንግስት 'government' has the following similar words:

---

[መንግሥት](#) 'government';

[ህዝብ](#) 'people'

[ፓርቲ](#) 'party'

[ኢህአዴግ](#) 'EPRDF'

But the following are collocations

[አንድ](#) 'one', [የኢትዮጵያ](#) 'that of Ethiopia', [ጥያቄ](#) 'question'.

There are punctuation marks and symbols that come as words in the thesaurus as well. They need to be taken care of.

Example of the thesaurus results for *government*:

*መንግስት* 'government'

**Amharic WaC [2013 + 2015 + 2016] freq = [17,661](#)** (870.54 per million)

Lemma	Score	Freq
<a href="#">መንግሥት</a>	0.490	16,300
<a href="#">'government;</a>		
<a href="#">“ punctuation</a>	0.333	102,859
<a href="#">ህዝብ 'people'</a>	0.331	15,175
<a href="#">ፓርቲ 'party'</a>	0.327	12,220
<a href="#">ኢህአዴግ</a>	0.319	6,503
<a href="#">'EPRDF'</a>		
<a href="#">አንድ 'one'</a>	0.315	50,442
<a href="#">ጊዜ 'time'</a>	0.312	64,423
<a href="#">ቤት 'house'</a>	0.310	37,620
<a href="#">ሁኔታ</a>	0.303	26,643
<a href="#">'situation'</a>		
<a href="#">የኢትዮጵያ</a>		
<a href="#">'that of</a>	0.301	24,067
<a href="#">Ethiopia'</a>		
<a href="#">ሕዝብ</a>	0.299	11,645
<a href="#">'people'</a>		
<a href="#">ጉዳይ 'affair',</a>	0.299	23,272
<a href="#">'issue'</a>		
<a href="#">” punctuation</a>	0.298	108,069
<a href="#">መንገድ 'road'</a>	0.295	24,393
<a href="#">ኢትዮጵያ</a>	0.294	20,687
<a href="#">'Ethiopia'</a>		
<a href="#">በኢትዮጵያ</a>	0.294	12,532
<a href="#">'That of</a>		
<a href="#">Ethiopia'</a>		
<a href="#">ድርጅት</a>	0.294	9,368
<a href="#">'organization'</a>		
<a href="#">ወቅት</a>	0.293	19,982
<a href="#">'season'</a>		
<a href="#">ሰው 'man'</a>	0.293	41,861
<a href="#">ወያኔ</a>		
<a href="#">'derogatory</a>	0.292	7,085
<a href="#">name of</a>		
<a href="#">TPLF'</a>		
<a href="#">አገር 'country'</a>	0.290	15,184

## 5. Somali

Evaluated by Ahmed Yusuf

### Word frequency list

[http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=sowac16&refs=%3Dword&wlmitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include\\_nonwords=1&wlnums=frq&wltype=simple](http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=sowac16&refs=%3Dword&wlmitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include_nonwords=1&wlnums=frq&wltype=simple)

Almost all words on the frequency list are correct Somali words. There is the word ‘home’ in the top list, which might be the ‘home’ button from the home page source. When in addition inspecting the Somali WaC [2016], a few errors were found, for example nouns modified by ‘dariiq’; however, except for *nin*, *waana*, *meel*, and *mid*, all the rest would definitely be verbs together with ‘dariiq’. The word *dariiq* (‘way, road, path’) is itself a noun, from Arabic طريق. In the list there is only one noun: *nin* (‘man’). This needs to be further checked.

### Concordances

[http://corpora.fi.muni.cz/habit/run.cgi/first\\_form?corpname=sowac16](http://corpora.fi.muni.cz/habit/run.cgi/first_form?corpname=sowac16)

An example of a Somali concordance is given below.

Query **dariiq** 709 (8.89 per million)

Page 1 of 36 Go [Next](#) | [Last](#)

<a href="#">infopankki.fi</a>	labaatankii saac. </p><p> Hosteelku caadi ahaan waa <b>dariiq</b> ka jaban hoteelka, lakiin heerka tayada
<a href="#">haatuf.net</a>	siyaasada taasina ay shirka ka dhigtay mid <b>dariiq</b> khaldan ku socda sidaa darteed ayay yidhaahdeen
<a href="#">shabellenews.com</a>	ay dib ugu soo noqon karto ay doonayaan, <b>dariiq</b> kasta haloo mare,taasi oo keeni karta in
<a href="#">hablaha.com</a>	gurigiisa, ooriidisa dhabta ah ee haddii aad <b>dariiq</b> wada martaan aanad ka fikiraynin in lagu
<a href="#">haatuf.net</a>	umada reer Somaliland, waxay usoo mareen <b>dariiq</b> dheer, waxaynu umad ahaan ku difaacnaa
<a href="#">hoygasuugaanta.com</a>	gabyeen rag kaloo gabayaa waaweyn ahaa oo uu <b>dariiq</b> cad marsiiyey. Gabayaagii Cabdi Gahayr
<a href="#">allgedo.com</a>	lasiyo wey wada gadgadn waana qaadicul <b>dariiq</b> aan wadini aheyn majeerteen muxuu kadali
<a href="#">puntlandtop.net</a>	baxaa, dhawr bilood. </p><p> Waxaa jira dhawr <b>dariiq</b> oo dhakhtarku u baari karo hadaad sonkor
<a href="#">halbeegnews.com</a>	waxa Ummada reer Somaliland ay usoo martay <b>Dariiq</b> aad u dheer oo dhib badnaa iyo dedaal xooggan
<a href="#">gabiley.net</a>	siyaasiyan-na u fiicneyn in uu waddan qudhi <b>dariiq</b> wax usoo maraan u noqdo. Wuxu lagama-maarmaan
<a href="#">gabiley.net</a>	gaajoonaya loo waddo hellaan marin amma <b>dariiq</b> ay maraan oo si dhakhso ah ku gaadhi karaan
<a href="#">gabiley.net</a>	Meles oo waxa Berbera Corridor lagu daray <b>dariiq</b> -yada muhiimada dhaqaale u leh Afrika-da
<a href="#">gabiley.net</a>	1961 inqilaabkii uu sababay midowgii aan <b>dariiq</b> caddaaladeed iyo wax wada lahaanshaha loo
<a href="#">allsomali24.org</a>	</p><p> 3- Beeshu waxay is ticmaali doontaa <b>dariiq</b> kasta oo ay xaqeeda ku heleyso. </p><p> 1-B
<a href="#">tooshnews.net</a>	aaminsan in khilaafaadka Dawliga ah lagu xalilo <b>dariiq</b> Nabadeed iyo Wadahadal, waxaanay aaminsantahay
<a href="#">somaliland.org</a>	Mustafe Cabdi Ciise (Mustafe Shiine) “Nin <b>dariiq</b> xariiqdaa garan-kara dayoobay, dhidarkii
<a href="#">somaliland.org</a>	waxa ay cadeynaysaa in qofka leh ama jeexda <b>dariiq</b> uu doonayo in uu meel [...] </p><p> Maayarka
<a href="#">somaliland.org</a>	dibna reer Somaliland taas mid la mid ah <b>dariiq</b> loogu jeexayaa, in lagu sameeyo tartan
<a href="#">somalilandpost.net</a>	Ismaaciil (Saylici) waxa uu sheegay in saddex <b>dariiq</b> ay sahlayaan in ardaydu ay waxbarasho helaan
<a href="#">tooshnews.net</a>	weeyi inaynu hal u jeedo yeelanyo inaynu hal <b>dariiq</b> u wada marno aayaha mustaqbalka caruurteena

Page 1 of 36 Go [Next](#) | [Last](#)

## Word sketches

An example of a Somali word sketch is given below.

### dariiq Somali WaC [2016] freq = 709 (8.89 per million)

modifiers of "darii"	416	1.00	verbs with "darii" as object	90	1.60	nouns modified by "darii"	421	1.00
kaldan	<u>13</u>	8.06	garan	<u>3</u>	3.37	qaadicul	<u>5</u>	8.58
xariiqdaa	<u>3</u>	7.86	hal	<u>9</u>	2.66	jeexda	<u>3</u>	7.85
sax	<u>13</u>	7.48	ahayd	<u>5</u>	1.29	jeexay	<u>3</u>	7.23
saxa	<u>4</u>	7.45	waa	<u>35</u>	1.22	maro	<u>11</u>	7.20
toosan	<u>6</u>	7.09	yahay	<u>7</u>	0.13	maraya	<u>4</u>	5.99
qaldan	<u>7</u>	6.89				labadaas	<u>4</u>	5.75
kastana	<u>4</u>	6.63	verbs with "darii" as subject	<u>11</u>	0.70	marin	<u>5</u>	5.37
dastuuri	<u>3</u>	6.54	garan	<u>3</u>	5.45	helaan	<u>6</u>	4.96
hebel	<u>5</u>	6.47				martay	<u>3</u>	4.81
sharci	<u>30</u>	6.46	"darii" and/or ...	<u>183</u>	1.20	hesho	<u>3</u>	4.77
yada	<u>3</u>	5.53	wado	<u>3</u>	6.18	maray	<u>7</u>	4.29
qudha	<u>3</u>	5.43	jirin	<u>4</u>	6.10	helo	<u>5</u>	3.95
wanaagsan	<u>19</u>	5.05	heer	<u>3</u>	4.28	taagan	<u>4</u>	3.58
nabadeed	<u>3</u>	5.03				kara	<u>3</u>	3.25
dheer	<u>25</u>	4.93				helay	<u>4</u>	3.24
cad	<u>9</u>	4.85				sameeyo	<u>3</u>	3.22
fudud	<u>3</u>	4.83				ahayn	<u>4</u>	3.07
furax	<u>3</u>	4.52				jiro	<u>4</u>	2.86
xun	<u>6</u>	4.49				leedahay	<u>3</u>	2.85
walba	<u>15</u>	4.25				jirin	<u>4</u>	2.83
qofka	<u>3</u>	3.40				nin	<u>7</u>	2.70
fiican	<u>3</u>	3.13				waana	<u>3</u>	2.10
kale	<u>35</u>	2.58				meel	<u>4</u>	2.03
cusub	<u>18</u>	2.39				mid	<u>8</u>	1.57
iska	<u>5</u>	1.81						

## 6. Tigrinya

*Evaluated by Shimelis Mazengia and Tsegay (Tigrinya Team, AAU)*

### Word frequency list

[http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=tiwac16&refs=%3Dword&wlmaxitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include\\_nonwords=1&wlnums=frq&wltype=simple](http://corpora.fi.muni.cz/habit/run.cgi/wordlist?corpname=tiwac16&refs=%3Dword&wlmaxitems=200&wlsort=f&wlattr=word&wlminfreq=5&wlmaxfreq=0&include_nonwords=1&wlnums=frq&wltype=simple)

1. Stop words and punctuations are presented as the most frequent words. This shows that the text needs cleaning and normalization. For instance, the full stop in Tigrinya has not been normalized. At times, it appears as a single punctuation mark (: 116,579 times). At other times, it appears as a combination of two two-dots (:: 545 times).

2. The most frequent nouns in the first list of 200 words are:

ኣምላኽ	<u>9,089</u>
ሓደ	<u>8,791</u>
ቅዱስ	<u>8,330</u>
ሰብ	<u>7,862</u>
ህዝቢ	<u>4,896</u>

3. The most frequent verbs in the first list of 200 words are auxiliary verb as seen below:

ዘሎ	<u>7212</u>
ኣሎ	4303
ኢሎ	4228
ማለት	4145
ዝነበረ	4098

4. The first main verb was encountered in the third list of 200 words which is ጽግሩ (632 times).

5. In the word list, there are numerals included in the category of the most frequent words.

6. There are some dependent morphemes and also letters which appear as frequent words (ን ይ ዶ ሞ).

7. The first element of some compound nouns (which is modified for compounding) is presented as independent frequent word (ስነ ደቀ). These were found in the third list of 200 words.



## Concordances

[http://corpora.fi.muni.cz/habit/run.cgi/first\\_form?corpname=tiwac16](http://corpora.fi.muni.cz/habit/run.cgi/first_form?corpname=tiwac16)

The five nouns in 2 and five verbs in 3 above were checked in the concordance. They were found to be in the correct Tigrinya contexts. They all make sense in their collocations.

An example of a concordance in Tigrinya is given below:

Query **ዘሎ** 7,212 (2,848.97 per million)

Page 1 of 361 Go [Next](#) | [Last](#)

zaratbebat.com	ብዛዕባ መፈጠር ወዲ ሰብ ናይ'ዚ ግዜ'ዚ ዓውደ ፍልጠት ኣቐራቢዎ ከም <b>ዘሎ</b> መረዳኣታ'ኳ ... ምሉእ ትሕዝቶ <b>ዘሎ</b> 12/16/2015   Comments
nselam.com	ተቓውሞ ከስፖም ኣይኮነን ምንቕ ንኹብል ኣውን ከልኪልዎን ዓፈትዎን <b>ዘሎ</b> ሓላቲ ፍጹም ዘዓቀራርብን ዘዓረዳድኡን በታኒ ሓሳብ ገለ ውድባት
nselam.com	ዝመርሖን ዝዘውሮን ፖጅላ ህግጻፍ ኣብ ናይ ቀቢጸ ተስፋ ኩነታት ኢዩ <b>ዘሎ</b> :: ልዕሊ እዚ በሲኡ ፈጋዕጋዕ ዝበለ ኣዋን ንጺምበ ተቓውሞ ዝበለጸን
nselam.com	ብመትከልን ዕላማን ተኣማሚንካ ኣብ ምስራሕ ስግኣትን ትዕዝብትን ከም <b>ዘሎ</b> ዝከሓድ ኣይኮነን:: ስለዚ ከምቲ ናይ ዝሓለፈ ንስራሕ ዘይኮነ ንስለ
nselam.com	ሕቶ ንስለ ሕቶ ኣይኮነን ኣልዒለዮ ዘለኹ:: ኣብ ጸምበ ተቓውሞ <b>ዘሎ</b> ከፍወስ ዘይተኸለለ ከሮኒክ ሕማም ስለ ዝኾነ ኢዩ:: ፖጅላ ህግጻፍ
nselam.com	ስለ ዝኾነ ኢዩ:: ፖጅላ ህግጻፍ ነቲ ፍትሒ ይንገስ ኢኡ ዝቃወሞ <b>ዘሎ</b> ኤርትራዊ ንገለ ምስ ሃይማኖት ንገለ ምስ ወገን ኣና ጠቀነ ዕድመ
nselam.com	ማእለርትን ግልጽነትን መሪርዎ እግሩ ናብ ዝመረሐ ፋሕ ብትን ኣብ ዝብለሉ <b>ዘሎ</b> እዋን: ንፍትሕን ዲሞክራሲን ንቃለስ ኣለና ካብ ዝብሉ መራሕቲ
gov.et	የብሉን:: <b>ዘሎ</b> መ.ግዝያዊ ኮሚቴ ብሓባር ብምጥጋን ኣብ ዕዳጋ ካብ <b>ዘሎ</b> ቀረብ ዘሓሸ ፅሬትን ዋጋን ዘለዎ ወይ ዘለዎም መቐረብቲ ብምምራፅ
nselam.com	2003ዓ.ም ጀሚሩ ዝፀገዐ ይኸውን:: <b>ዘሎ</b> ነዚ ንምግባር ዚጸፍኣና <b>ዘሎ</b> ምድራዊ ወይ ሓላፊ ነገር ንምዕታር ወይ ክብርን ዝናን ተመነና
nselam.com	ዩ:: ስለዚ ሃገር ናይ ሰባት ድኻ ትኸውን ከምዘላ ብግልጺ ዚረኣ <b>ዘሎ</b> ኩሉ: ኣብ ዝኾነማሳበራዊ ፖዳይ ከም እንዕዘቦ: እቲ መንእሰይ
nselam.com	ኣይከውንን: ሃገር ከኣ ሃገር ኣይትኸውንን:: <b>ዘሎ</b> እዚ ዘብለና <b>ዘሎ</b> ድማ: ኣብ ቅድሚ ዓይንና ይፍጸም ዘሎ: ናይ ሰባት ቀጻሊ ዋሕድን
nselam.com	ኣይትኸውንን:: <b>ዘሎ</b> እዚ ዘብለና ዘሎ ድማ: ኣብ ቅድሚ ዓይንና ይፍጸም <b>ዘሎ</b> : ናይ ሰባት ቀጻሊ ዋሕድን ዋሕዝን ብሕቲ ዘይኮነስ ንጽባሕ ከጋጥም
nselam.com	ዘይኮነስ ንጽባሕ ከጋጥም ዘለዎ ጸንታ እዩ እቲ ዚረኣየናን ዘፍረሓናን <b>ዘሎ</b> :: ሓጺ ስጺተይና ኣብ ዝኾኑ ኸይዱ:ናይ ሃገሩ ናፍቐትን ዝኾርን
nselam.com	ሕጂ ንኡስን ማእከላይን ዝበልፍኡ ወለዶ ጥራሕ ኣይኮነን ዚጠፍእ <b>ዘሎ</b> : እቲ ኪትከኦ ዝነበሮ: ማለት እቶም ብሕፃናቶም ዝኾኑ: ኣብ
nselam.com	ትርጉም ምርኩብ ይከኣል..."17:: <b>ዘሎ</b> ነዚ ንምግባር "ዚጸፍኣና <b>ዘሎ</b> ምድራዊ ወይ ሓላፊ ነገር ንምዕታር ወይ ክብርን ዝናን ተመነና
nselam.com	? ዕላማኡኹ? ናብ'ቲ ሓቀይና ታሕጓስ ዚወስድ መንገዱ ኣበይ እዩ <b>ዘሎ</b> ? ብዛዕባ ሞት/ፍርዲ:ከንዩው ሞት ብዛዕባ ዘሎ ዓስቢ ተግባራትካ
nselam.com	መንገዱ ኣበይ እዩ ዘሎ? ብዛዕባ ሞት/ፍርዲ:ከንዩው ሞት ብዛዕባ <b>ዘሎ</b> ዓስቢ ተግባራትካ ብዚርኢኹ እቲ ሓቁ ኣየናይ እዩ? ንህላዌና ከቢብዎ
nselam.com	ምጁኦ በሓቁ "ዜሕንኹ ተቐባልነት ዘይብሉ: ምስ'ቲ ወድሰብ በጺሕዎ <b>ዘሎ</b> ናይ ምዕባሌን ሥልጣኔን ጸረጃ ዘይኸይድ" ኪብሉ ርክስ ሊቃነጳጳሳት
nselam.com	ከኣ: ህልው መነባብሮ ዘይቐወር ኮይንዎ: በቲ ኣብ ስግር-ማዕዶ <b>ዘሎ</b> ንብረትን: ዘይረኣዮ ሕይወትን ተሳሒቡ: እቲ ጸላይ ለውጢ መንእሰይ
nselam.com	ዘይረኣዮ ሕይወትን ተሳሒቡ: እቲ ጸላይ ለውጢ መንእሰይ ዚገብሮ <b>ዘሎ</b> ፖዕዞ: ኪጅምር ከሎ ብሕጋዊ መንገድን ኣገባብን ዝነበረ: ሕጂ

Page 1 of 361 Go [Next](#) | [Last](#)

## Word sketches

A word sketch is given below:

# ዘሎ

Tigrinya WaC [2016] freq = [7,212](#) (2,848.96 per million)

modifiers of "ዘሎ"	336	0.40
ሕጂ	<a href="#">123</a>	9.95
ሎሚ	<a href="#">82</a>	9.49
ቅድሚ	<a href="#">37</a>	8.35
ገና	<a href="#">10</a>	8.19
ድሕሪ	<a href="#">25</a>	8.02
ደው	<a href="#">13</a>	7.96
ልዕሊ	<a href="#">22</a>	7.75
ምሳይ	<a href="#">4</a>	7.73
ብሓባር	<a href="#">4</a>	7.18
ሰለስተ	<a href="#">4</a>	6.88

objects of "ዘሎ"	5,342	2.90
ስርዓት	<a href="#">101</a>	8.82
ህዝቢ	<a href="#">78</a>	8.21
ኩነታት	<a href="#">52</a>	8.00
ፍልልይ	<a href="#">42</a>	7.94
ሰማይ	<a href="#">45</a>	7.87
አዋን	<a href="#">57</a>	7.81
ፖለቲካዊ	<a href="#">41</a>	7.72
ሰማያት	<a href="#">37</a>	7.70
፡	<a href="#">61</a>	7.64
ህዝብና	<a href="#">35</a>	7.62
አበይ	<a href="#">33</a>	7.57
ቃልሲ	<a href="#">36</a>	7.53
ሃገርና	<a href="#">33</a>	7.51
ዓለም	<a href="#">44</a>	7.45
ግዜ	<a href="#">49</a>	7.42
ጉዳይ	<a href="#">33</a>	7.37
ሰብ	<a href="#">62</a>	7.35
ጸገማት	<a href="#">27</a>	7.26
ሓቂ	<a href="#">31</a>	7.17
ጉጅለ	<a href="#">25</a>	7.09
ጽሑፍ	<a href="#">25</a>	7.08
ዝምድና	<a href="#">23</a>	7.06
ሕገ	<a href="#">27</a>	6.98
ጸገም	<a href="#">23</a>	6.98
ሽግር	<a href="#">22</a>	6.98

subjects of "ዘሎ"	553	1.30
መንጎ	<a href="#">52</a>	10.95
ኤርትራን	<a href="#">12</a>	8.98
ዝባን	<a href="#">7</a>	8.62
ሓጢአት	<a href="#">7</a>	8.51
ጽፍሪ	<a href="#">5</a>	8.19
መንግሥ	<a href="#">5</a>	8.15
ክርስቶስ	<a href="#">10</a>	8.08
ዓይኒ	<a href="#">5</a>	8.03
ደንበ	<a href="#">5</a>	7.97
ደምበ	<a href="#">5</a>	7.92
ምእንታኹም	<a href="#">4</a>	7.85
ተራ	<a href="#">4</a>	7.66
ሰማይን	<a href="#">4</a>	7.64
ስድራ	<a href="#">4</a>	7.43
ዓዲ	<a href="#">4</a>	7.42
ጉጅለ	<a href="#">4</a>	7.37
ስርዓት	<a href="#">6</a>	7.34

pronominal subjects of "ዘሎ"	25	0.80
ንሱ	<a href="#">14</a>	7.23
ኣነ	<a href="#">5</a>	6.22

## 7. Conclusion

There is no doubt that the four Ethiopian corpora with their web based search and results interfaces are a major contribution to the possibilities for linguistic research into these languages in the future. The same is true for the Norwegian resources. The table below summarizes the results.

	Norwegian Bokmål	Norwegian Nynorsk	Amharic	Oromo	Somali	Tigrinya
<b>No. of Words in Corpus</b>	1,178,357,993	54,511,854	17,320,000	4,249,953	71,871,585	2,087,613
<b>No. of Documents</b>	3,443,807	214,379	33,542	8,851	385,338	1,907
<b>Grammatically Tagged</b>	tagged	no	yes	no	no	yes
<b>Concordance</b>	yes	yes	yes	yes	yes	yes
<b>Word List</b>	yes	yes	yes	yes	yes	yes
<b>Word Sketches</b>	yes	no	yes	no	yes	yes
<b>Thesaurus and Word Cloud</b>	yes	no	yes	no	no	yes

## References

Guevara, Emiliano Raul (2010). [NoWaC: a large web-based corpus for Norwegian.](#)  
In *Proceedings of the NAACL HLT 2010 Sixth Web as Corpus Workshop*, Association for Computational Linguistics page 1 - 7.