

## Msimulizi as corpus for accurate search<sup>1</sup>

Arvi Hurskainen  
Department of World Cultures, Box 59  
FIN-00014 University of Helsinki, Finland  
[arvi.hurskainen@helsinki.fi](mailto:arvi.hurskainen@helsinki.fi)

### Abstract

In Technical Report 60<sup>2</sup>, I described the process of converting printed text into machine-readable form. This report is an extension to it, and here I will go into more detail in describing and demonstrating the capabilities of the search system based on analysed text. All the material on Msimulizi (years 1888-1896) that is available on SOAS web page was processed into machine-readable form, including manual editing of the whole text. The second round of editing was done on the basis of computational analysis, which points out the remaining scanning mistakes. The clean text was then converted into an analysed format, which is optimal for information retrieval. The report demonstrates especially such search tasks, which are hardly possible using conventional string search, due to the complex word structure of Swahili.

**Key Words:** *information retrieval, morphological analysis.*

### 1 Introduction

*Msimulizi* is the first printed periodical in Swahili language. It started to appear in 1888 with intervals of two months. Its aim was to provide a communication channel between various mission stations of the Universities, Mission to Central Africa. The stations are located in the area of the current Tanzania and Malawi. In addition, news in the area of current Uganda are also reported. The mission agency had also work there, which is why there was great interest to hear news from that area also.

The aim of the periodical was to have news reports from each station in each issue. Only seldom this fully realized, due to various reasons. The periodical includes also general news on current themes, such as conflicts between colonial powers (England, Germany, Portugal) and the local Arab administration.

The Universities' Mission to Central Africa was established as response to the alarming news on ruthless slave trade of Arabs in eastern and central Africa. Especially the reports of Livingstone gave impetus to action. *Msimulizi* has many reports on how slaves were still raided and sold to Arabs, and how mission agencies tried to get slaves to mission schools and give them Christian education. Former slaves played important roles in mission schools.

---

<sup>1</sup> The report is issued under licence CC BY-NC

<sup>2</sup> <http://www.njas.helsinki.fi/salama/printed-text-into-machine-readable-form.pdf>

According to *Msimulizi* texts, mission work was very much concentrated on education. Medical care was also part of the work, but it is referred to only when someone badly ill is moved to another station for treatment.

Training continued to up to six years of education, and examinations were arranged at certain intervals. *Msimulizi* reports on these examinations and mentions even the names of those who passed on each class.

*Msimulizi* reveals that Africa was dangerous to whites. Astonishingly many of the missionaries died during those years. There is no diagnosis of deaths, but fever is usually reported to be the reason of death. Obviously, malaria was the main cause of deaths.

*Msimulizi* follows contemporary historical events. In addition to Livingstone, who already had died that time, the periodical tells about Stanley, who travelled across the continent. Also, the whereabouts of Emin Pasha in Uganda were reported on.

## 2 Language

The Swahili language of the periodical is not even. The writing style and even the orthography varies depending on the writer. It is not known how much the reports were edited before printing. Yet different styles are apparent.

Considering the time of writing - about 140 years ago - the language has changed remarkably little since then. There are underdeveloped elements, such as the lack of month names (English forms are used) and the Arabic use of some numbers. Also, the naming of the week days is searching for its form.

Most of the *Msimulizi* text was written by advanced students, selected as reporters in each mission station. Reporting is detailed in relation to time. Usually the day of the month, and often also the week day, and sometimes also the time of the day, are mentioned in reporting on events. The events themselves were usually related to travelling and local festivities, which are usually Christian, but also sometimes traditional.

From the viewpoint of constructing an accurate information retrieval system for *Msimulizi*, it is important to be able to treat also the non-standard forms properly. The morphological analyzers normally recognize the standard forms only, and non-standard forms are treated with a guesser. This method is very unreliable, especially if the non-standard forms are verbs. Therefore, I have included the non-standard forms of *Msimulizi* into the morphological lexicon, so that they can be retrieved as accurately as the standard forms.

## 3 The structure of the enriched text

The accurate retrieval system described here uses the so-called enriched text as search target, in this case the *Msimulizi* text. The enriched text is produced so that the normal text is passed through the analysis and morphological disambiguation, and the result contains morphological interpretation for each morpheme of the word. In the case of verbs, the interpretation may be quite complex due to several consecutive morphemes.

In the final enriched text, most morphological information is removed. Only the word lemma and its POS class code are retained, and all the rest is removed. As a result, the enriched text is clearly structured, because after each word form comes the lemma of that

word and its POS tag. Such text form is also easy to convert into more readable form in information retrieving phase.

The original text in *Msimulizi* is located under various headings. Most headings include the name of the mission station, such as HABARI ZA MISOZWE (News from Misozwe), but also topic names, such as HABARI ZA VITA (War news) are included. In the edited and modified version, each sentence is on its own line. In the beginning of the sentence, there is an identification code area of eight spaces. The identification code has the form 'NYA-88-', where the first part means NJASSA and the last part the year 1888.

The identification code has two functions. First, it helps in retrieving material on a certain subject (in this case news from Njassa area), or from a certain year. Second, the code can be produced in each search result, if needed.

The list of topic names and their abbreviations are in (1).

(1)

BOO - BOONDE	mku - MKUZI
chi - CHISUMULU	MLI - MLIMA WA MLINGA
CHI - CHITANGALI	CHI - MLOLELA
CHU - CHUAKA	MSA - MSARAKA
HAD - HADITHI	MSU - MSUMBA
KIC - KICHELWE	MWI - MWITI
KIL - KILIMANI	NEW - NEWALA
KIU - KIUNGANI	NAN - NANYANGA
KOR - KOROGWE	NYA - NYASSA
KOT - KOTAKOTA	PAC - PACHIA
KUF - KUFARIKI	PWA - PWANI
LIK - LIKOMA	SAF - SAFARI
MAC - MACHEMBA	SHA - SHAMBA LA MBWENI
MAG - MAGILA	UGA - UGANDA
MAS - MASASI	UNA - UNANGU
MBW - MBWENI	VIT - VITA
MIS - MISOZWE	vit - VITENDAWILI
MIW - MIWA	WAR - WARENO
MKU - MKUNAZINI	YER - YERUSALEMI

In case the section of *Msimulizi* does not fall in any of those categories, the category name MSI (MSIMULIZI) is used. Therefore, each sentence has the two-part identification code.

An extract from the enriched *Msimulizi* text is in (2).

(2)

KIU-88- Lakini {lakini\_CONJ} zamani {zamani\_ADV} hakutaka {taka\_V} hatta {hatta\_ADV} kidogo {kidogo\_AD-ADJ} , {,\_COMMA} sasa {sasa\_ADV} anataka {taka\_V} sana {sana\_AD-ADJ} .  
KIU-88- Katika {katika\_PREP} Oct. {Oct.\_PROPN} 4 {4\_NUM} tulimsafirisha {safirisha\_V} Christopher {Christopher\_PROPN} Kiwapo {wapo\_V} mwalimu {mwalimu\_N} wa {wa\_GEN-CON} Kiungani {kiunga\_N} kwenda {kwenda\_V} Newala {Newala\_N} kufundisha {fundisha\_N} kule {kule\_PRON} .

KIU-88- Amefuatana {fuatana\_V} na {na\_PREP} Charles {Charles\_PROPN} W. {W\_N} Rehani {rehani\_N} , {,\_COMMA} mpishi {mpishi\_N} , {,\_COMMA} aliyemwoa {oa\_V} siku\_hizi {siku\_hizi\_ADV} Hope {Hope\_PROPN} Chela {Chela\_PROPN} , {,\_COMMA} Mbweni {mbwe\_N} .  
KIU-88- Twasikia {sikia\_V} wameshuka {shuka\_V} salama {salama\_ADV} Lindi {Lindi\_N} .  
KIU-88- Tunatumaini {tumaini\_V} tutaingia {ingia\_V} Kanisa {kanisa\_N} letu {etu\_PROPN} jipya {pya\_ADJ} mwezi {mwezi\_N} ukiandama {andama\_V} .  
KIU-88- Bwana {Bwana\_N} Askofu {askofu\_N} ameleta {leta\_V} simu {simu\_N} kusema {sema\_V} atafika {fika\_V} Unguja {Unguja\_N} Oct. {Oct.\_PROPN} 31 {31\_NUM} .

The conversion of the *Msimulizi* text into sentence-per-line format was challenging, because the punctuation was non-systematic, and often the normal full stop at the end of the sentence was replaced with a comma or semicolon, or even with an empty space. Full stops after month days also disturbed the process, and a lot of post editing was needed for correcting wrong sentence splitting. But this was due to the non-conventional use of punctuation marks, and this should not be found in most texts.

#### 4 Search methods

The enriched form of text makes it possible to direct search on the surface form or the analyzed form of the word. Both search methods are needed in practice. However, considering the complex word structures in Swahili, the possibility of using the analyzed form as a search key is useful.

Consider the search task, where you want to find information on deaths in *Msimulizi*. Obvious search keys are words, which mean death or dying. In Swahili, such words are the verbs *kufariki* and *kufa*, as well as the noun *kifo* (pl. *vifo*). It appears that the noun *kifo* does not appear at all. The verb *kufa* appears 164 times, and the verb *kufariki* 90 times.

Especially the verb *kufa* is problematic, because its only constant element is the stem *f* and all the rest is inflection. We can see all the inflected forms of this verb in (3).

(3)

2 Alikufa [fa_V]	1 hufa [fa_V]
1 Alipokufa [fa_V]	1 iliokufa [fa_V]
1 Kikafa [fa_V]	1 imekufa [fa_V]
3 afe [fa_V]	6 kafa [fa_V]
25 akafa [fa_V]	22 kufa [fa_V]
2 akifa [fa_V]	5 nife [fa_V]
11 alikufa [fa_V]	1 nikifa [fa_V]
2 alipokufa [fa_V]	1 ningekufa [fa_V]
6 aliyekufa [fa_V]	1 tukifa [fa_V]
34 amekufa [fa_V]	1 tulikufa [fa_V]
1 angalikufa [fa_V]	2 tunakufa [fa_V]
1 angekufa [fa_V]	1 ufe [fa_V]
5 atakufa [fa_V]	8 vifao [fa_V]
2 hafi [fa_V]	2 wafe [fa_V]
1 hajafa [fa_V]	9 wakafa [fa_V]
3 hakufa [fa_V]	3 waliokufa [fa_V]
1 hawakufa [fa_V]	10 wamekufa [fa_V]

```
3 wanakufa [fa_V]
1 wangekufa [fa_V]
1 wasife [fa_V]
```

The list shows that without analysed text form the search task would be impossible. We can produce a similar list of the occurrences of the verb *kufariki* (4).

```
(4)
2 Alifariki [fariki_V]
1 Alipofariki [fariki_V]
1 Amefariki [fariki_V]
23 akafariki [fariki_V]
8 alifariki [fariki_V]
1 aliofariki [fariki_V]
1 alipofariki [fariki_V]
1 alivyofariki [fariki_V]
5 aliyefariki [fariki_V]
35 amefariki [fariki_V]
5 kafariki [fariki_V]
2 kufariki [fariki_V]
6 waliofariki [fariki_V]
3 wamefariki [fariki_V]
```

We see that the list for *kufariki* is much shorter, which means that less forms are constructed with this stem. This may be partly due to the fact that the basic verb is *kufa*, while the verb *kufariki* is a loan word from Arabic and its connotation is more solemn. Yet one could consider using the normal string search for finding the occurrences, because the stem is *fariki*, including also the final vowel, as loan word verbs do in Swahili.

We can retrieve also other monosyllabic verbs from *Msimulizi*. Such verbs include: *kula* (to eat), *kuwa* (to be), *kunya* (to rain), *kucha* (to rise), *kuja* (to come), and *kupa* (to give). These verbs have also extended forms, but we will concentrate here on base forms (5).

(5) *kula* - to eat

```
2 Chala [la_V]
2 Kapala [la_V]
1 Kaule [la_V]
2 Mpale [la_V]
1 Mtakula [la_V]
17 Mtaula [la_V]
1 Pala [la_V]
1 Wakala [la_V]
1 Zaila [la_V]
5 akala [la_V]
2 akila [la_V]
4 alao [la_V]
1 alikula [la_V]
1 amekula [la_V]
2 anakula [la_V]
1 anawala [la_V]
1 asile [la_V]
4 hali [la_V]
1 hamli [la_V]
1 hatukula [la_V]
1 hawakula [la_V]
2 hawali [la_V]
1 hazikula [la_V]
6 hula [la_V]
1 kuila [la_V]
35 kula [la_V]
1 kuleni [la_V]
2 kumla [la_V]
```

2	kutawala [la_V]	1	vyakula [la_V]
1	lala [la_V]	1	waile [la_V]
1	laleni [la_V]	1	wakaila [la_V]
1	nimekula [la_V]	15	wakala [la_V]
1	nitazila [la_V]	1	wakila [la_V]
47	tukala [la_V]	1	walapo [la_V]
1	tukamla [la_V]	3	wale [la_V]
1	tule [la_V]	2	walikula [la_V]
1	tulipokula [la_V]	1	waliokula [la_V]
1	tumekula [la_V]	2	wamekula [la_V]
1	tunakula [la_V]	1	wanile [la_V]
2	tutakula [la_V]	1	watakula [la_V]
1	twalikula [la_V]	6	yale [la_V]
2	ulalo [la_V]	1	zikala [la_V]
2	unakula [la_V]		

All the forms are theoretically legal forms of the verb *kula*, but there are some words, which are in fact proper names in text. Disambiguation had failed in those cases. Also such words as *wale*, *yale* and *lala* are obviously not forms of the verb *kula*. However, the vast majority of hits are correct.

(6)

kuwa - to be

2	Akawa [wa_V]	3	Walikuwa [wa_V]
7	Alikuwa [wa_V]	3	Walipokuwa [wa_V]
1	Alipokuwa [wa_V]	1	Yakuwa [wa_V]
1	Asiwe [wa_V]	1	Yalikuwa [wa_V]
1	Hakuwa [wa_V]	1	ahilikuwa [wa_V]
1	Hapakuwa [wa_V]	1	ahtulikuwa [wa_V]
20	Ikawa [wa_V]	48	akawa [wa_V]
60	Ilikuwa [wa_V]	24	akiwa [wa_V]
3	Imekuwa [wa_V]	2	alikuwa [wa_V]
2	Itakuwa [wa_V]	210	alimokuwa [wa_V]
1	Itakuwaje [wa_V]	2	aliokuwa [wa_V]
5	Kiwayo [wa_V]	5	alipokuwa [wa_V]
1	Lilikuwa [wa_V]	17	alivyokuwa [wa_V]
1	Mwe [wa_V]	1	alivyokuwa [wa_V]
2	Nalikuwa [wa_V]	39	alivyokuwa [wa_V]
1	Niliokuwa [wa_V]	15	amekuwa [wa_V]
1	Nilipokuwa [wa_V]	2	anakuwa [wa_V]
6	Palikuwa [wa_V]	2	angalikuwa [wa_V]
4	Tukawa [wa_V]	2	angekuwa [wa_V]
2	Tukiwa [wa_V]	2	asiwe [wa_V]
2	Tulikuwa [wa_V]	1	atakayekuwa [wa_V]
1	Tulipokuwa [wa_V]	9	atakuwa [wa_V]
2	Tuwe [wa_V]	9	awe [wa_V]
2	Ulikuwa [wa_V]	3	haijawa [wa_V]
1	Ulipokuwa [wa_V]	3	haikuwa [wa_V]
2	Uwe [wa_V]	1	hajawa [wa_V]

3 hakuwa [wa\_V]  
2 halikuwa [wa\_V]  
6 hapakuwa [wa\_V]  
2 hatukuwa [wa\_V]  
1 hatuwi [wa\_V]  
1 hawakuwa [wa\_V]  
1 hawatakuwa [wa\_V]  
2 hayakuwa [wa\_V]  
2 huwa [wa\_V]  
2 ijapokuwa [wa\_V]  
1 ikauwa [wa\_V]  
161 ikawa [wa\_V]  
1 ikawaje [wa\_V]  
1 ilikokuwa [wa\_V]  
292 ilikuwa [wa\_V]  
1 iliokuwa [wa\_V]  
8 ilipokuwa [wa\_V]  
4 ilivyokuwa [wa\_V]  
5 iliyokuwa [wa\_V]  
20 imekuwa [wa\_V]  
3 inakuwa [wa\_V]  
2 ingalikuwa [wa\_V]  
2 isiwe [wa\_V]  
1 itakavyokuwa [wa\_V]  
20 itakuwa [wa\_V]  
2 itakuwaje [wa\_V]  
1 iwayo [wa\_V]  
12 iwe [wa\_V]  
1 kawa [wa\_V]  
2 kilichokuwa [wa\_V]  
1 kilikuwa [wa\_V]  
1 kilivyokuwa [wa\_V]  
1 kimekuwa [wa\_V]  
1 kinakuwa [wa\_V]  
7 kulikuwa [wa\_V]  
1 kulipokuwa [wa\_V]  
5 kumekuwa [wa\_V]  
1 kunakuwa [wa\_V]  
50 kuwa [wa\_V]  
1 likawa [wa\_V]  
2 likiwa [wa\_V]  
16 lilikuwa [wa\_V]  
1 lilipokuwa [wa\_V]  
4 limekuwa [wa\_V]  
1 litakuwa [wa\_V]  
1 liwalo [wa\_V]  
1 liwe [wa\_V]  
2 mkiwa [wa\_V]  
1 mliokuwa [wa\_V]  
1 mlivyokuwa [wa\_V]  
1 mlizokuwa [wa\_V]  
1 mtakuwa [wa\_V]  
6 mwe [wa\_V]  
7 nalikuwa [wa\_V]  
8 nikawa [wa\_V]  
2 nikiwa [wa\_V]  
3 nilikuwa [wa\_V]  
7 nilipokuwa [wa\_V]  
1 niliyokuwa [wa\_V]  
1 nimekuwa [wa\_V]  
2 nitakuwa [wa\_V]  
2 niwe [wa\_V]  
1 pakawa [wa\_V]  
29 palikuwa [wa\_V]  
3 palipokuwa [wa\_V]  
1 palivyokuwa [wa\_V]  
2 pamekuwa [wa\_V]  
1 pasiwe [wa\_V]  
1 patakuwa [wa\_V]  
2 sijawa [wa\_V]  
1 sikuwa [wa\_V]  
1 tukakuwa [wa\_V]  
40 tukawa [wa\_V]  
6 tukiwa [wa\_V]  
36 tulikuwa [wa\_V]  
2 tuliokuwa [wa\_V]  
14 tulipokuwa [wa\_V]  
2 tulivyokuwa [wa\_V]  
2 tuliyokuwa [wa\_V]  
4 tumekuwa [wa\_V]  
2 tungalikuwa [wa\_V]  
4 tutakuwa [wa\_V]  
9 tuwe [wa\_V]  
11 twalikuwa [wa\_V]  
19 ukawa [wa\_V]  
7 ukiwa [wa\_V]  
37 ulikuwa [wa\_V]  
4 uliokuwa [wa\_V]  
3 ulivyokuwa [wa\_V]  
1 umekuwa [wa\_V]  
1 unakuwa [wa\_V]  
1 utakavyokuwa [wa\_V]  
1 utakuwa [wa\_V]  
10 uwe [wa\_V]  
1 uwepo [wa\_V]  
2 vilikuwa [wa\_V]  
14 wakawa [wa\_V]  
9 wakiwa [wa\_V]  
92 walikuwa [wa\_V]  
25 waliokuwa [wa\_V]  
12 walipokuwa [wa\_V]  
4 walivyokuwa [wa\_V]

1 walizokuwa [wa_V]	2 yamekuwa [wa_V]
12 wamekuwa [wa_V]	2 yamekuwaje [wa_V]
1 wanakuwa [wa_V]	2 yatakavyokuwa [wa_V]
1 watakaokuwa [wa_V]	1 yatakuwa [wa_V]
1 watakiwa [wa_V]	2 yatakuwaje [wa_V]
7 watakuwa [wa_V]	7 yawe [wa_V]
7 wawe [wa_V]	5 zikawa [wa_V]
8 yakawa [wa_V]	1 zikiwa [wa_V]
5 yakuwa [wa_V]	16 zilikuwa [wa_V]
21 yalikuwa [wa_V]	2 zimekuwa [wa_V]
1 yaliokuwa [wa_V]	1 zitakuwa [wa_V]
1 yalipokuwa [wa_V]	1 ziwazo [wa_V]
1 yalivyokuwa [wa_V]	2 ziwe [wa_V]

(7)

kunya - to rain

6 ikanya [nya\_V]  
2 ikinya [nya\_V]  
5 ilikunya [nya\_V]  
2 ilipokunya [nya\_V]  
1 imekunya [nya\_V]  
1 inakunya [nya\_V]  
1 inye [nya\_V]  
1 itakunya [nya\_V]  
2 kunya [nya\_V]  
1 unyayo [nya\_V]  
1 watawanya [nya\_V]  
2 yanya [nya\_V]

(8) kucha - to rise

3 Kakucha [cha\_V]  
1 Kukicha [cha\_V]  
1 Mkacha [cha\_V]  
1 acha [cha\_V]  
2 ache [cha\_V]  
3 kucha [cha\_V]  
2 kulipokucha [cha\_V]  
1 kunakucha [cha\_V]  
1 mwacha [cha\_V]  
4 wacha [cha\_V]  
6 wache [cha\_V]  
1 wacheni [cha\_V]

(9)

kuja - to come

15 Akaja [ja_V]	1 Anakuja [ja_V]
2 Akija [ja_V]	1 Ikaja [ja_V]
7 Alikuja [ja_V]	2 Ilikuja [ja_V]
1 Alipokuja [ja_V]	3 Njoo [ja_V]
1 Amekuja [ja_V]	1 Njoooni [ja_V]



1 Tukaja [ja\_V]  
1 Tulikuja [ja\_V]  
1 Tutakuja [ja\_V]  
1 Ukija [ja\_V]  
1 Ulikuja [ja\_V]  
1 Vikija [ja\_V]  
1 Wajao [ja\_V]  
1 Waje [ja\_V]  
6 Wakaja [ja\_V]  
1 Wakija [ja\_V]  
7 Walikuja [ja\_V]  
1 Waliokuja [ja\_V]  
2 Walipokuja [ja\_V]  
1 Wamekuja [ja\_V]  
6 aja [ja\_V]  
8 ajapo [ja\_V]  
8 ajaye [ja\_V]  
27 aje [ja\_V]  
172 akaja [ja\_V]  
18 akija [ja\_V]  
58 alikuja [ja\_V]  
2 aliokuja [ja\_V]  
14 alipokuja [ja\_V]  
9 aliyekuja [ja\_V]  
56 amekuja [ja\_V]  
30 anakuja [ja\_V]  
2 anayekuja [ja\_V]  
1 angekuja [ja\_V]  
8 asiye [ja\_V]  
1 asiyekuja [ja\_V]  
1 atakayekuja [ja\_V]  
18 atakuja [ja\_V]  
1 bwanatwaja [ja\_V]  
2 haijaja [ja\_V]  
1 haiji [ja\_V]  
3 haikuja [ja\_V]  
2 haiwaji [ja\_V]  
5 hajaja [ja\_V]  
5 hakuja [ja\_V]  
1 hamji [ja\_V]  
1 hamngekuja [ja\_V]  
1 hatukuja [ja\_V]  
1 havikuja [ja\_V]  
4 hawajaja [ja\_V]  
4 hawaji [ja\_V]  
14 hawakuja [ja\_V]  
1 hawangalikuja [ja\_V]  
1 hazijaja [ja\_V]  
34 huja [ja\_V]  
1 hukuja [ja\_V]  
2 ijapo [ja\_V]  
2 ijayo [ja\_V]  
2 ije [ja\_V]  
17 ikaja [ja\_V]  
3 ikija [ja\_V]  
15 ilikuja [ja\_V]  
5 ilipokuja [ja\_V]  
1 iliyokuja [ja\_V]  
12 imekuja [ja\_V]  
16 inakuja [ja\_V]  
2 isiye [ja\_V]  
1 itakuja [ja\_V]  
13 kaja [ja\_V]  
1 kijacho [ja\_V]  
3 kikaja [ja\_V]  
1 kinakuja [ja\_V]  
128 kuja [ja\_V]  
3 kumekuja [ja\_V]  
1 kunakuja [ja\_V]  
5 likaja [ja\_V]  
2 linakuja [ja\_V]  
1 linapokuja [ja\_V]  
1 litakalokuja [ja\_V]  
4 mje [ja\_V]  
1 mlipokuja [ja\_V]  
1 mlivyokuja [ja\_V]  
2 mmekuja [ja\_V]  
1 msije [ja\_V]  
1 mwaja [ja\_V]  
7 naja [ja\_V]  
2 nakuja [ja\_V]  
1 nalikuja [ja\_V]  
2 niye [ja\_V]  
5 nikaja [ja\_V]  
1 nikiya [ja\_V]  
4 nimekuja [ja\_V]  
2 ninakuja [ja\_V]  
1 ninaokuja [ja\_V]  
2 nitakuja [ja\_V]  
6 njoo [ja\_V]  
4 njooni [ja\_V]  
3 palikuja [ja\_V]  
2 pamekuja [ja\_V]  
1 sikuja [ja\_V]  
2 tuje [ja\_V]  
21 tukaja [ja\_V]  
1 tukaje [ja\_V]  
1 tukija [ja\_V]  
1 tulikuja [ja\_V]  
1 tulipokuja [ja\_V]  
10 tumekuja [ja\_V]  
4 tunakuja [ja\_V]

2	tutakuja [ja_V]	26	waliokuja [ja_V]
1	twaja [ja_V]	5	walipokuja [ja_V]
3	ujao [ja_V]	34	wamekuja [ja_V]
3	uje [ja_V]	26	wanakuja [ja_V]
4	ukaja [ja_V]	6	wanaokuja [ja_V]
3	ukija [ja_V]	1	wanavyokuja [ja_V]
1	ulikuja [ja_V]	1	wangekuja [ja_V]
2	ulipokuja [ja_V]	4	wasije [ja_V]
2	umekuja [ja_V]	1	wasipokuja [ja_V]
1	unakuja [ja_V]	1	watakaokuja [ja_V]
1	unapokuja [ja_V]	9	watakuja [ja_V]
1	ungekuja [ja_V]	3	yaja [ja_V]
1	utakuja [ja_V]	1	yajayo [ja_V]
2	vikaja [ja_V]	1	yaje [ja_V]
1	vikija [ja_V]	1	yakaja [ja_V]
1	vilikuja [ja_V]	1	zije [ja_V]
23	wajao [ja_V]	4	zikaja [ja_V]
3	wajapo [ja_V]	1	zilikuja [ja_V]
22	waje [ja_V]	1	zilipokuja [ja_V]
157	wakaja [ja_V]	3	zimekuja [ja_V]
19	wakija [ja_V]	1	zinakuja [ja_V]
95	walikuja [ja_V]		

Note that the list also includes such exceptional cases as the imperative forms *njoo* (come, sg) and *njooni* (come, pl).

(11)

kupa - to give

2	Akampa [pa_V]	1	atupaye [pa_V]
1	Akanipa [pa_V]	1	avipaye [pa_V]
2	Nikawapa [pa_V]	1	awape [pa_V]
2	Nipe [pa_V]	1	haikumpa [pa_V]
1	Nipeni [pa_V]	1	hakuwapa [pa_V]
2	Tupe [pa_V]	1	hamnipi [pa_V]
1	Tupeni [pa_V]	1	hapi [pa_V]
1	Wakawapa [pa_V]	1	hauwapi [pa_V]
1	Walimpa [pa_V]	1	hawakunipa [pa_V]
17	akampa [pa_V]	1	hawawapi [pa_V]
6	akanipa [pa_V]	1	hukunipa [pa_V]
16	akawapa [pa_V]	4	humpa [pa_V]
3	alimpa [pa_V]	1	hunipi [pa_V]
1	alinipa [pa_V]	2	huwapa [pa_V]
2	aliwapa [pa_V]	1	kampa [pa_V]
1	aliyempa [pa_V]	1	kawapa [pa_V]
1	aliyenipa [pa_V]	2	kukupa [pa_V]
1	aliyewapa [pa_V]	1	kulipa [pa_V]
1	alizompa [pa_V]	10	kumpa [pa_V]
1	amempa [pa_V]	1	kunipa [pa_V]
1	amewapa [pa_V]	1	kutupa [pa_V]
4	anipe [pa_V]	7	kuwapa [pa_V]
1	atawapa [pa_V]	1	kuwapeni [pa_V]

1 mkinipa [pa_V]	1 tumempa [pa_V]
2 mnipe [pa_V]	2 tumpe [pa_V]
1 mpe [pa_V]	1 tupe [pa_V]
1 mpeni [pa_V]	1 tupeni [pa_V]
1 mtampa [pa_V]	1 tutampa [pa_V]
1 napa [pa_V]	1 tutawapa [pa_V]
1 nikakupeni [pa_V]	1 twakupa [pa_V]
5 nikampa [pa_V]	3 unipe [pa_V]
4 nikawapa [pa_V]	1 utupe [pa_V]
1 nikiwapa [pa_V]	8 wakampa [pa_V]
1 nikupe [pa_V]	2 wakanipa [pa_V]
2 nimekupa [pa_V]	2 wakawapa [pa_V]
2 nipe [pa_V]	2 wakimpa [pa_V]
6 nitakupa [pa_V]	2 walimpa [pa_V]
1 nitakupeni [pa_V]	1 walipe [pa_V]
1 niwapazo [pa_V]	1 walizompa [pa_V]
2 niwape [pa_V]	1 wamenipa [pa_V]
2 sikupeni [pa_V]	1 wampa [pa_V]
1 simpi [pa_V]	3 wampe [pa_V]
1 tukajipa [pa_V]	1 watakupa [pa_V]
1 tukampa [pa_V]	1 watanipa [pa_V]
2 tukiwapeni [pa_V]	1 wawape [pa_V]
1 tulipe [pa_V]	

Above we have examples of the monosyllabic verbs in base form. They have also several extended forms, such as applicative, causative, reciprocal and passive forms. Below are some examples of these.

(12)

**kufia - die on behalf of**

2 Alifia [fia\_V]  
1 afie [fia\_V]  
2 akafia [fia\_V]  
2 alikofia [fia\_V]  
1 aliyefia [fia\_V]  
1 atawafia [fia\_V]  
1 hakufia [fia\_V]  
1 kufia [fia\_V]  
1 walifia [fia\_V]  
1 wamewafia [fia\_V]

(13)

**kupeana - to give to each other**

2 kupeana [peana\_V]  
1 wakipeana [peana\_V]  
1 wapeane [peana\_V]

(14)

**kupewa - to be given (that is, to get)**

1 Akapewa [pewa\_V]                      1 Apewa [pewa\_V]

1	Apewaje [pewa_V]	1	lipewe [pewa_V]
5	Tukapewa [pewa_V]	1	mmepewa [pewa_V]
1	Tulipewa [pewa_V]	1	napewa [pewa_V]
1	Wakapewa [pewa_V]	2	nikapewa [pewa_V]
1	Walipewa [pewa_V]	1	nililopewa [pewa_V]
1	Wamepewa [pewa_V]	3	nimepewa [pewa_V]
17	akapewa [pewa_V]	1	sikupewa [pewa_V]
2	akipewa [pewa_V]	34	tukapewa [pewa_V]
1	alilopewa [pewa_V]	2	tuliopewa [pewa_V]
1	aliopewa [pewa_V]	13	tulipewa [pewa_V]
8	alipewa [pewa_V]	1	tulivyopewa [pewa_V]
1	alipopewa [pewa_V]	1	tuliyopewa [pewa_V]
6	amepewa [pewa_V]	2	tumepewa [pewa_V]
1	anapewa [pewa_V]	2	tupewayo [pewa_V]
1	apewa [pewa_V]	1	ukipewa [pewa_V]
1	apewazo [pewa_V]	33	wakapewa [pewa_V]
1	asipopewa [pewa_V]	2	wakipewa [pewa_V]
3	atapewa [pewa_V]	4	waliopewa [pewa_V]
1	hawajapewa [pewa_V]	18	walipewa [pewa_V]
1	hawakupewa [pewa_V]	1	walipopewa [pewa_V]
1	hawapewi [pewa_V]	1	walizopewa [pewa_V]
2	hupewa [pewa_V]	13	wamepewa [pewa_V]
2	kapewa [pewa_V]	2	wapewe [pewa_V]
16	kupewa [pewa_V]	1	wasipewe [pewa_V]
1	likapewa [pewa_V]		

(15)

**kulisha** - to make eat, that is: to feed

1	Kalisha [lisha_V]
1	hatukulishi [lisha_V]
2	kuwalisha [lisha_V]
1	wakatulisha [lisha_V]
2	wanaonilisha [lisha_V]
1	watanilisha [lisha_V]

## 5 Searching on the basis of the surface form

Traditionally, the surface string search was the only method in information retrieval. In isolating languages such as English, it still is a viable method. The situation becomes much more complex with heavily inflecting languages such as Swahili or Finnish. The above examples demonstrate this. The direct string search is viable also in Swahili, if we want to search on the basis of surface strings, be they whole words or part of words.

With the search method demonstrated here, the search key may be the beginning or end part of the word, or any part inside the word, or the whole word. The system finds the matches and surrounds the found key section of the word with square brackets. Examples are below.

(16)

Key word is: aliwa

NEW-96- Dec. 21 , Mwalimu G. Yohana Mpalila [aliwa]sili hapa toka Mwiti .  
NEW-96- Dec. 22 , Mwalimu Augustino [aliwa]sili akitoka Miwa .  
NEW-96- Assubuhi ya Uz[aliwa] tuliimba Carols tukizunguka mjini tukaingia  
Kanisani saa a kwanza mpaka saa 4 kasoro dakika kumi .  
NEW-96- Siku ile ya Uz[aliwa] watoto waliambiwa marks V zao wakapokea  
haki zao .  
NEW-96- Dec. 26 , Mwalimu Filipino [aliwa]leta watoto wake wapokee zawadi  
zao walioshinda na wasioshinda pia .  
NEW-96- Dec. 28 , Archdeacon Farler [aliwa]sili hapa pamoja\_na Mwalimu  
James Chigulu kwani yee ni mgonjwa wa\_macho .  
NEW-96- Jioni yake watu walikusanyika nyumbani kuonyeshwa sanamu za taa  
Takatifu za Uz[aliwa] .  
MAS-96- Bassi barua ilipofika pwani ndipo akatokea Bwana mkubwa wa Lindi  
kuja kutaka mali yote waliotwaa kwa watu wa safari , na mali  
yaliyotw[aliwa] katika vita hii ndio pembe , watumwa , na tumbako .  
MAS-96- Dec. 23 , Padre Hugh na watu wake walio Wamasihiya na wanafunzi  
wakaja kujiweka tayari kwa\_siku kuu ya Uz[aliwa] .  
MAS-96- Penyi Nangoo nikapotea njia , kwani wenzangu n[aliwa]acha mitini  
wakila matunda , nikafika mpaka ndani huko kwenyi mashamba ya akina Satia  
.  
NYA-96- Bassi tukajaribu kuvunja rupia , ache tusumbuke , killa tuendapo  
inakat[aliwa] kwa\_sababu rupia zitumikazo huku zimepigwa\_chapa , bassi  
tukajaribu killa duka mwisho kwa\_bahati tukavunja .  
NYA-96- Tena karibu wazungu wetu wote walikata\_tamaa walitaka kurudi kwao  
Ulaya , illa Bwana Faulkner [aliwa]tuliza , akisema "Haithuru twendeni  
kuko\_huko Nyasa .

The result can be arranged also so that the key word is aligned, and a chosen amount of context on each side can be printed. In example (17), the context size is defined on the basis of the number of character slots. In this case the left context is 20 slots and the right context 30 slots.

(17)

NEW-96- u G. Yohana Mpalila [aliwa]sili hapa toka Mwiti .  
NEW-96- , Mwalimu Augustino [aliwa]sili akitoka Miwa .  
NEW-96- Assubuhi ya Uz[aliwa] tuliimba Carols tukizu  
NEW-96- Siku ile ya Uz[aliwa] watoto waliambiwa mark  
NEW-96- 26 , Mwalimu Filipino [aliwa]leta watoto wake wapoke  
NEW-96- , Archdeacon Farler [aliwa]sili hapa pamoja\_na Mwa  
NEW-96- a taa Takatifu za Uz[aliwa] .  
MAS-96- i , na mali yaliyotw[aliwa] katika vita hii ndio p  
MAS-96- i kwa\_siku kuu ya Uz[aliwa] .  
MAS-96- a , kwani wenzangu n[aliwa]acha mitini wakila matu  
NYA-96- illa tuendapo inakat[aliwa] kwa\_sababu rupia zitum  
NYA-96- illa Bwana Faulkner [aliwa]tuliza , akisema "Haith

The context can also be defined on the basis of the number of words on each side. In (18) the context is three words on each size. Also, separate punctuation marks are treated as words.

(18)

NEW-96- G. Yohana Mpalila [aliwa]sili hapa toka Mwiti  
NEW-96- , Mwalimu Augustino [aliwa]sili akitoka Miwa .  
NEW-96- Assubuhi ya Uz[aliwa] tuliimba Carols tukizunguka  
NEW-96- Siku ile ya Uz[aliwa] watoto waliambiwa marks  
NEW-96- , Mwalimu Filipino [aliwa]leta watoto wake wapokee  
NEW-96- , Archdeacon Farler [aliwa]sili hapa pamoja\_na  
NEW-96- taa Takatifu za Uz[aliwa] .  
MAS-96- , na mali yaliyotw[aliwa] katika vita hii  
MAS-96- kwa\_siku kuu ya Uz[aliwa] .  
MAS-96- , kwani wenzangu n[aliwa]acha mitini wakila matu  
NYA-96- , killa tuendapo inakat[aliwa] kwa\_sababu rupia  
NYA-96- illa Bwana Faulkner [aliwa]tuliza , akisema "Haithuru

We see in the examples above that the direct string search method is not ideal in this case, because hits include verb prefix combinations, verb stems or part of them, and parts of proper names. The examples serve merely as demonstration on how results can be presented in various ways.

We get more precise results when we use the analysed forms as search key.

## 6 Searching on the basis of analysis result

Keeping the context on both sides of the hit as three words, we test with two stems, *uzaliwa\_N* and *zaliwa\_V* (19-20).

(19)

Search key: *uzaliwa\_N* (birth, that is Christmas)

CHU-89- Alhamisi baada\_ya Uzaliwa [uzaliwa\_N] sisi Waalimu wa  
MAS-89- kuingia tunangojea Uzaliwa [uzaliwa\_N] .  
SHA-89- kuu za Uzaliwa [uzaliwa\_N] tunaona hapana pah  
YER-89-S i kuu ya Uzaliwa [uzaliwa\_N] , 1888 .  
YER-89-N ana zote za Uzaliwa [uzaliwa\_N] mnazoimba leo Mbwe  
CHI-90- hii ya Uzaliwa [uzaliwa\_N] hatukupata kustare  
CHI-90- ndio karibu\_na Uzaliwa [uzaliwa\_N] naona Dec. 24  
MAG-91- yetu ya Uzaliwa [uzaliwa\_N] ilistawi sana .  
MIS-91- kuu ya Uzaliwa [uzaliwa\_N] tulikuwapo wote La  
MAG-92- yetu ya Uzaliwa [uzaliwa\_N] ilitufurahisha san  
MAG-92- kuu ya Uzaliwa [uzaliwa\_N] .  
KOR-92- kuu ya Uzaliwa [uzaliwa\_N] wa Bwana wetu  
KIU-93- mosi mbele\_ya Uzaliwa [uzaliwa\_N] ikawa shughuli yet  
CHI-94- kuu ya Uzaliwa [uzaliwa\_N] vyakula vimeisha m  
CHI-94- kuu ya Uzaliwa [uzaliwa\_N] .  
KIU-94- kuu ya Uzaliwa [uzaliwa\_N] tukapamba nyumba y  
KIU-94- nyimbo za Uzaliwa [uzaliwa\_N] ( Christmas Carols  
KIU-94- Siku ya Uzaliwa [uzaliwa\_N] alitusayidia kwa u  
NEW-95- kuu ya Uzaliwa [uzaliwa\_N] ilikuwa sala kubwa  
mku-95- huku mpaka Uzaliwa [uzaliwa\_N] , nami nalipata  
MBW-96- nyimbo za Uzaliwa [uzaliwa\_N] .  
MBW-96- siku ya Uzaliwa [uzaliwa\_N] wa Bwana wetu  
KOR-96- sikukuu ya Uzaliwa [uzaliwa\_N] wa Bwana tukahamis  
MWI-96- 'holidays' za Uzaliwa [uzaliwa\_N] .  
NEW-96- Assubuhi ya Uzaliwa [uzaliwa\_N] tuliimba Carols tu

NEW-96-                    ile ya Uzaliwa [uzaliwa\_N] watoto waliambiwa  
NEW-96-                    Takatifu za Uzaliwa [uzaliwa\_N] .  
MAS-96-                    kuu ya Uzaliwa [uzaliwa\_N] .

(20)

Search key: zaliwa\_V (to be born)

YER-89-                    usiku ule alipozaliwa [zaliwa\_V] Mwana wa Daudi  
YER-89-                    huonyesha mahali alipozaliwa [zaliwa\_V] Isa Masiya na  
MBW-89-                    , mtoto amezaliwa [zaliwa\_V] Sept. 10 ,  
MKU-90-                    wa watumwa watakozaliwa [zaliwa\_V] sasa watakuwa huru  
HAD-90-                    watu wengine waliozaliwa [zaliwa\_V] katika ustaarabu ,  
SHA-91-                    Alhamisi alizaliwa [zaliwa\_V] mtoto kwa Alfred  
CHI-92-                    ndiyo siku aliyozaliwa [zaliwa\_V] Sultani wao Madachi  
NAN-92-                    mtoto mchanga aliyezaliwa [zaliwa\_V] kama siku themintas  
MSI-93-                    yakawa kama amezaliwa [zaliwa\_V] nayo .  
CHI-94-                    27 , walizaliwa [zaliwa\_V] watoto wawili Yusuf  
NEW-95-                    nzige hawa waliozaliwa [zaliwa\_V] hapa , kwani  
MBW-95-                    Denys , alizaliwa [zaliwa\_V] Feb. 18 ,

Note that when we use the base form as search key, the hit will be produced immediately after the surface form of the word.

## 7 Choosing source text

The identification code at the beginning of each sentence gives a possibility to define the search material. For example, we can extract all the sentences that deal with war in 1888 (21).

(21)

VIT-88-                    Wadoicha {mdoicha\_N} wanaona {ona\_V} fetheha {fetheha\_N} sana  
{sana\_AD-ADJ} kama {kama\_ADV} mashauri {shauri\_N} yao {ao\_PRON} yote  
{ote\_PRON} yamefahitika V .  
VIT-88-                    Sasa {sasa\_ADV} wameshawishi {shawishi\_V} Waingreza  
{Waingreza\_PROPN} huko {huko\_ADV} Ulaya {Ulaya\_N} kushariki {shariki\_V}  
kazi {kazi\_N} ya {ya\_GEN-CON} kuzuia {zuia\_N} na {na\_CC} kuwaathibu  
{athibu\_N} watu {mtu\_N} wote {ote\_PRON} wanaochukua {chukua\_V} watumwa  
{tumwa\_V} baharini {bahari\_N} , {,\_COMMA} na {na\_CC} wanaochukua  
{chukua\_V} bunduki {bunduki\_N} ao {ao\_N} baruti {baruti\_N} .  
VIT-88-                    Bassi {bassi\_ADV} sasa {sasa\_ADV} wanaanza {anza\_V} kuvizia  
{zia\_V} pwani {pwani\_ADV} yote {ote\_PRON} , {,\_COMMA} Wadoicha  
{mdoicha\_N} wanavizia {zia\_V} tangu {tangu\_PREP} Tanga {Tanga\_N} hatta  
{hatta\_ADV} kisiwa {kisiwa\_N} Mafya {jifya\_N} na {na\_CC} Waingreza  
{waingreza\_N} tangu {tangu\_PREP} Mafya {jifya\_N} hatta {hatta\_ADV} Rovuma  
{Rovuma\_PROPN} , {,\_COMMA} na {na\_CC} tangu {tangu\_PREP} Tanga {Tanga\_N}  
hatta {hatta\_ADV} Lamu {Lamu\_N} .  
VIT-88-                    Merikebu {merikebu\_N} nyingi {ingi\_PRON} zinavizia {zia\_V} ,  
{,\_COMMA} na {na\_CC} merikebu {merikebu\_N} moja {moja\_NUM} kubwa  
{kubwa\_ADJ} sana {sana\_AD-ADJ} za {za\_GEN-CON} chuma {chuma\_N} tupu  
{tupu\_ADJ} , {,\_COMMA} nene {nene\_ADJ} kabisa {kabisa\_ADV} , {,\_COMMA}  
inakaa {kaa\_V} sikuzote {sikuzote\_ADV} bandarini {bandari\_N} hapa  
{hapa\_ADV} Unguja {Unguja\_N} ina {ina\_V} ngome {ngome\_N} mbili  
{mbili\_NUM} za {za\_GEN-CON} chuma {chuma\_N} tupu {tupu\_ADJ} juu\_yake  
{juu\_ake\_ADV} , {,\_COMMA} mnamo {mnamo\_PREP} mizinga {mzinga\_N} ya  
{ya\_GEN-CON} ajabu {ajabu\_N} mno {mno\_AD-ADJ} ya {ya\_GEN-CON} kupigia

{pigia\_N} adui {adui\_N} akikaa {kaa\_V} mbali {mbali\_ADV} mayili  
{mayili\_N} tano {tano\_NUM} .

The number of reports on war varies greatly between years. In (22) are statistics on the number of sentences each year.

(22)  
1888 - 4  
1889 - 93  
1890 - 3  
1891 - 0  
1892 - 22  
1893 - 0  
1894 - 0  
1895 - 0  
1896 - 0

However, these statistics do not give the full picture of the situation, because the reports take into account only those cases, where a specific title on war is present. War reports are also in many other reports. For this reason, we must take another approach. We search for information using keys based on actual words in text. Obvious key words are *vita* (war), *kupigana* (to fight), *bunduki* (gun), *mkuki* (spear), *mshale* (arrow). In order to save space, I will give only statistics on each word in *Msimulizi* (23).

(23)  
Key word: *vita\_N* (war) - 273  
          *pigana\_V* (fight) - 93  
          *bunduki\_N* (gun) - 170  
          *mkuki\_N* (spear) - 9  
          *mshale\_N* (arrow) - 7  
          *uthia\_N* (threat) - 120

The words *vita*, *pigana* and *uthia* appear quite often each year. This gives a more reliable picture of the situation than the picture given by the titles alone. The noun *bunduki* is not a very reliable indicator of war, because *bunduki* was very much used as salute to welcome an honoured guest. *Buduki* was so recklessly used that death accidents in welcome ceremonies were reported.

We can direct the search to certain material. For example, we can retrieve the war-related material from Uganda (24). For saving space, context is contracted.

(24)  
UGA-89- Bassi Uganda sasa vita [vita\_N] tupu wala\_mwisho  
UGA-90- shua ngine , walete vita [vita\_N] huko Sese , laki  
UGA-94- Hapana hofu ya vita [vita\_N] siku\_hizi .  
UGA-90- u , wakautwaa wakapigana [pigana\_V] sana wakatwaa  
UGA-90- saidia sana , tumepigana [pigana\_V] na wale adui z  
UGA-94- wa jamaa zetu tutapigana [pigana\_V] marra ; kwani  
UGA-94- ka inchi yetu , tupigane [pigana\_V] tena wapate ku



UGA-89- enda barazani na bunduki [bunduki\_N] zao tayari ,  
 UGA-89- alimpigia mfalme bunduki [bunduki\_N] asimpate , la  
 UGA-90- karibia ikapigwa bunduki [bunduki\_N] ya pwani ikap  
 UGA-89- ua wale wawili kwa mkuki [mkuki\_N] wake , lakini m  
 UGA-89- awili aliowaua kwa mkuki [mkuki\_N] wake , mmoja nd  
 UGA-89- Bassi marra uthia [uthia\_N] mtupu na machaf  
 UGA-90- esi Uganda , kwani uthia [uthia\_N] uko wa Kamega ,

## 8 Use of non-standard words in Msimulizi

The analysed version of Msimulizi makes it also possible to study the use of non-standard words. A word is considered non-standard, when it is no longer used in standard Swahili, and the word is intentionally written in the form as it was printed. The non-standard element may be in the form of the word stem or in inflection. *Msimulizi* contains also clear typing errors, and such errors were corrected in editing phase. If the word was not in standard form and it was repeatedly written in the same way, it was considered a non-standard form of the word. All such non-standard words were included into the morphological analyzer, so that they could be retrieved in any word form.

In (25) is a list of non-standard words of *Msimulizi*. Only such words are listed, where the non-standard element is in the word stem.

(25)

48 "adhuhuri" N	41 "burre" ADV
20 "afathali" ADV	2 "bwato" N
1 "ahadia" V	10 "chayi" N
10 "alasili" N	2 "chukuwa" V
2 "angawa" CONJ	9 "dayima" ADV
29 "anjili" N	2 "doktari" N
17 "asharini" NUM	1 "edaashara" NUM
1 "ashirini" NUM	361 "enyi" POSS-PRON
44 "asikari" N	2 "faraga" N
2 "athabu" N	4 "fathaika" V
3 "athibiwa" V	5 "fathili" N
3 "athibu" V	2 "fathili" V
1 "athimishwa" V	1 "fathiliwa" V
1 "athimiwa" V	16 "fayida" N
28 "balli" CONJ	5 "fayidi" V
88 "balyozi" N	2 "fayidia" V
3 "bandera" ADV	2 "fayidika" V
43 "bandera" N	2 "fayidiwa" V
29 "barozi" N	29 "fetha" N
49 "barra" N	2 "fetheha" N
34 "bass" ADV	1 "fithuli" N
1680 "bassi" ADV	6 "forotha" N
23 "billa" PREP	14 "fortha" N
1 "biskwiti" N	27 "frasi" N
4 "bithaa" N	9 "fullani" ADJ
16 "boonde" N	2 "fungasa" V
2 "borra" ADJ	2 "ghasiya" N
67 "buddi" N	6 "ghathabu" N

125	"ginsi" N	2	"mkohani" N
3	"hassa" ADV	6	"mngine" ADJ
3	"hatiya" N	3	"mosikiti" N
1474	"hatta" ADV	4	"moskiti" N
6	"hayi" ADJ	47	"mpila" N
3	"hifathi" V	2	"msayidia" N
4	"ilimu" N	6	"mshujaa" N
327	"illa" CONJ	6	"msikiaji" N
8	"illakini" CONJ	1	"mudda" N
55	"illi" CONJ	2	"mwaanafunzi" N
137	"ingine" ADJ	14	"nahotha" N
1	"itlafu" N	16	"nakhotha" N
2	"jaisha" V	1	"nassi" ADV
4	"jihathari" V	78	"naswi" ADV
3	"juwa" V	2	"nduwi" N
2	"kabithi" V	9	"nuss" ADJ-PRE
2	"kabithiwa" V	1	"nussra" ADJ-PRE
1	"kasikazini" ADV	102	"nussu" ADJ-PRE
2	"kassa" ADV	1	"ogolea" V
9	"kassi" ADV	4	"pilao" N
5	"kassorobo" ADV	1	"polomoka" V
2	"kassrobo" ADV	4	"raiya" N
1	"katha" ADJ	2	"rassi" N
19	"kathalika" ADV	37	"rathi" N
14	"kathawakatha" ADJ	6	"rayia" N
4	"kathi" N	1	"rayiya" N
6	"khutubu" V	2	"rhotuba" N
5	"kidoicha" N	1	"ridio" N
438	"killa" ADJ-PRE	3	"rithi" V
8	"kinwa" N	1	"rithiwa" V
34	"kitwa" N	4	"saalamu" N
2	"kuaheri" ADV	2	"sabwini" NUM
2	"kuaherini" ADV	3	"salaama" N
5	"kwiba" V	6	"salam" N
16	"kwitwa" V	1	"sayidia" N
6	"magaribi" ADV	142	"sayidia" V
4	"malkiya" N	3	"sayidiana" V
4	"mangaribi" N	12	"sayidiwa" V
32	"mangine" N	24	"selaha" N
23	"manovari" N	15	"sermala" N
527	"marra" ADV	6	"settini" NUM
10	"mathabahu" N	48	"shangwi" N
18	"mathbahu" N	5	"sharbati" N
18	"mayili" N	1	"shariki" N
2	"mayiti" N	37	"shariki" V
59	"mdoicha" N	4	"sharikisha" V
1	"mdoichi" N	1	"sharikishwa" V
17	"mgine" PRON	6	"shauko" N
25	"mhadi" N	2	"sherbeti" N
1	"mharamia" N	63	"shidda" N
1	"mjermani" N	1	"sigileti" N

1	"siswi" PRON	26	"tissa" NUM
1	"sitaashara" NUM	1	"tissini" NUM
1	"sittini" NUM	1	"torumpeta" N
11	"sumulia" V	1	"towa" V
2	"suruwali" N	5	"tufano" N
7	"swafi" ADJ	1	"tufanu" N
4	"tafathali" ADV	1	"uhayi" N
3	"tafathalini" ADV	11	"unuss" ADJ-PRE
6	"taraja" V	123	"unussu" ADJ-PRE
12	"thahabu" N	1	"uswafi" N
34	"thaifu" ADJ	3	"uthaifu" N
1	"thalathini" NUM	1	"uthalimu" N
6	"thambi" N	10	"uthi" V
1	"thaminia" V	120	"uthia" V
120	"thani" V	13	"uthika" V
15	"thania" V	1	"uthiwa" V
6	"thaniwa" V	2	"uwa" V
5	"tharau" V	1	"vunjilia" V
3	"tharuba" N	6	"walla" ADV
11	"thawabu" N	83	"wallakini" CONJ
2	"thehebu" N	419	"wagine" N
4	"thelitashara" NUM	8	"wagineo" N
4	"thoofika" V	7	"yayi" N
2	"thoofisha" V	138	"zayidi" AD-ADJ
2	"thoruba" N	1	"zowea" V
1	"thulumiwa" V	1	"zowelea" V
4	"thulumu" V	3	"zuwia" V
25	"thuru" V		

Many words are classified as non-standard because of the orthography. Typical examples are such consonants that are written with digraphs, such as *dh* > *th* (e.g. *thoruba*). Between two vowels is often added the half-consonant *w* (e.g. *ua* > *uwa*, *saidia* > *sayidia*). Doubling of the consonant occurs after syllables with stress (*basi* > *bassi*, *bali* > *balli*, *ila* > *illa*, *nusu* > *nussu*, *tisa* > *tissa*).

Non-standard morphemes, which do not feature in the above list, include such prefixes of verbs as *tuli* > *twali*, *mli* > *mwali*, *zili* > *zali*, *nita* > *nta*, *na* > *ha*. The prefixes include the subject marker and the TAM marker. Examples of verb forms are in (26).

(26)

[twali]penda  
 [twali]shinda  
 [twali]wapokea  
 [twali]watazamia  
 [twali]kaa  
 [twali]kusanyika  
 [twali]cheza  
 [twali]karibishwa  
 [mwali]nitendeaje  
 [zali]ungua  
 [zali]fanya

[zali]batilika  
[zali]tujia  
[zali]tazamiwa  
[nta]jaribu  
[nta]sahau  
[nta]uvunjilia  
[nta]chagua

## 9 Discussion

Archives contain lots of printed material, which is immensely valuable for many kinds of purposes. Computers offer means for permanent preservation of this material. The material is normally preserved in bitmap form, which itself is merely another way to make the material readable for human beings. Digital form would allow quick, accurate and comprehensive information retrieval, if text would be machine-readable. The OCR is a step forward, but unfortunately it works properly only, if the original text is clean. Therefore, in practice the OCR-read text must be corrected manually. Even this task is not easy, if the original text is difficult to read.

It is possible to construct post-editing programs, which do corrections automatically. However, only part of corrections can be done in this way without access to context. The only fully reliable method is to go the whole text through manually.

Manual correcting is often considered expensive, which is why cheaper methods have been sought for. We must consider the benefits of proper editing and weight them against costs. If the text to be edited, such as *Msimulizi*, has permanent value as source material, manual editing is not too expensive. Therefore, I have done it for *Msimulizi*. There are many more historical materials, which would deserve similar treatment as *Msimulizi*.

When the text is edited, it can be converted into the form suitable for computational treatment without much effort. This text form can then be processed into rich text format. Searches can then be made to this text format. The method enhances vastly more accurate and comprehensive search methods than traditional surface text search.