

General overview

| Corpus | Analytics date | Language |
|--------------------|----------------|-------------|
| kam_Latn.jsonl.tsv | 12/5/2024 | Kamba (kam) |

Volumes

| Docs | Segments | Unique segments | Tokens | Size | Characters |
|-------|----------|--------------------|--------|--------|------------|
| 1,183 | 14,259 | 9,563 (67.07 %) | 867K | 4.9 MB | 4,631,764 |

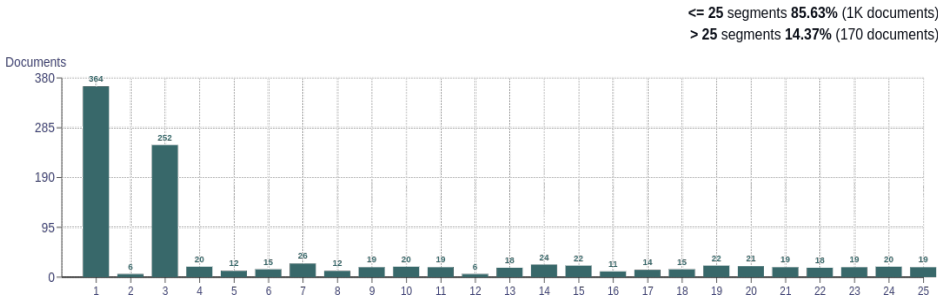
Top 10 domains

| Domain | Docs | % of total |
|------------------|------|------------|
| bible.is | 583 | 49.28 |
| jw.org | 510 | 43.11 |
| sangufm.co.ke | 32 | 2.70 |
| kituonline.com | 9 | 0.76 |
| gospelgo.com | 6 | 0.51 |
| jaladaafrica.org | 5 | 0.42 |
| watchtower.org | 5 | 0.42 |
| graduates.com | 4 | 0.34 |
| rmsradio.co.ke | 4 | 0.34 |
| 4laws.com | 3 | 0.25 |

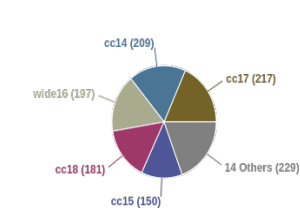
Top 10 TLDs

| Domain | Docs | % of total |
|--------|------|------------|
| is | 583 | 49.28 |
| org | 526 | 44.46 |
| co.ke | 38 | 3.21 |
| com | 31 | 2.62 |
| net | 3 | 0.25 |
| io | 1 | 0.08 |
| web.id | 1 | 0.08 |

Documents size (in segments)

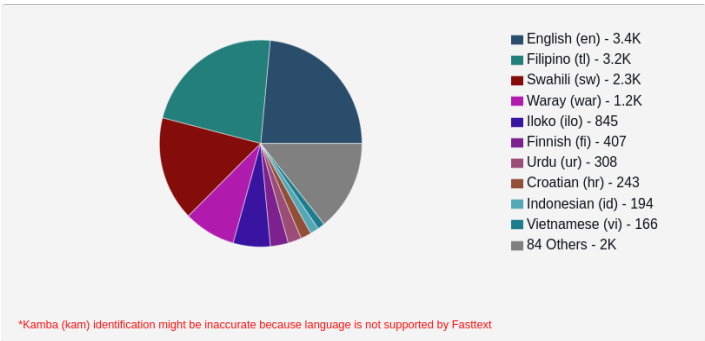


Documents by collection

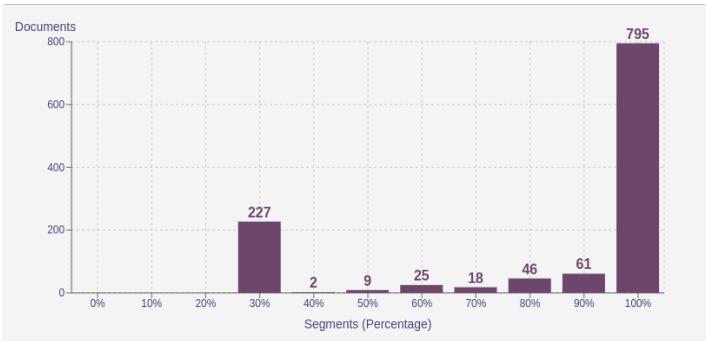


Language Distribution

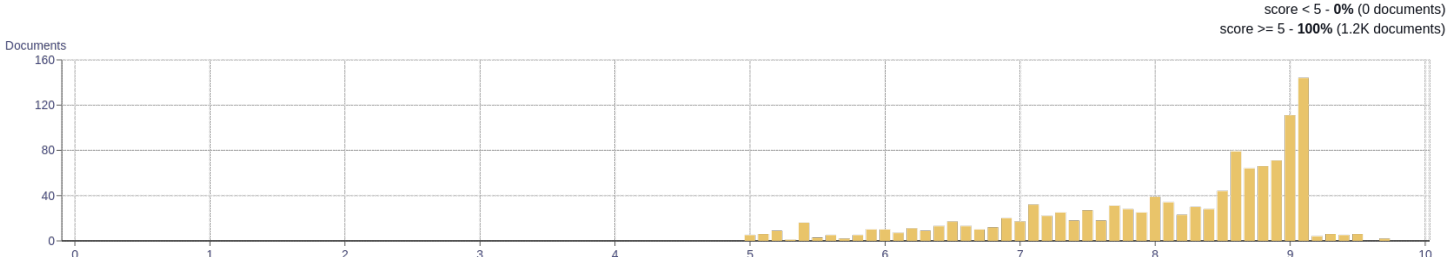
Number of segments



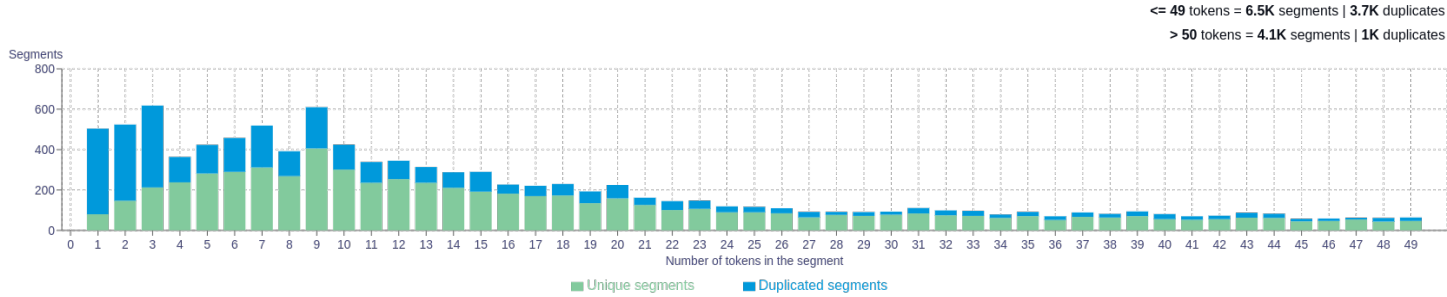
Percentage of segments in Kamba (kam) inside documents



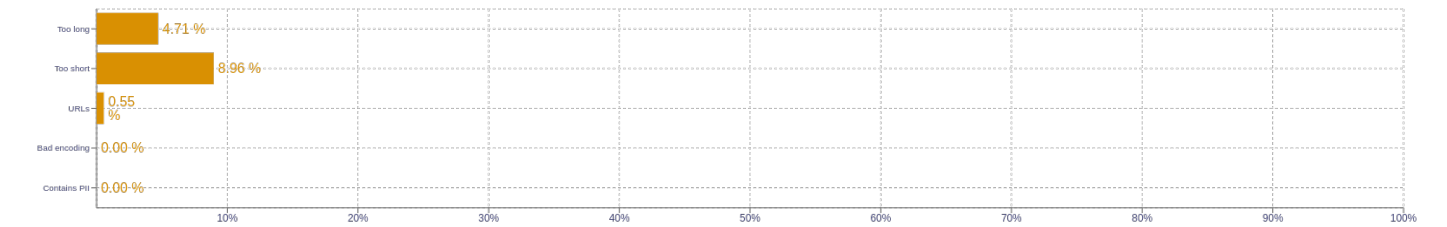
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

| Size | n-grams |
|------|--|
| 1 | <div>kwa 4012</div> <div>ali 3916</div> <div>andũ 3849</div> <div>yeova 3785</div> <div>sika 3780</div> |
| 2 | <div>pĩna ali 424</div> <div>sya yeova 408</div> <div>ngũsĩ sya 397</div> <div>kũũ nthĩ 372</div> <div>sika ali 362</div> |
| 3 | <div>ngũsĩ sya yeova 393</div> <div>iũlũ wa nthĩ 320</div> <div>tene na tene 261</div> <div>atamu na eva 181</div> <div>munumba ya kĩnzua 177</div> |
| 4 | <div>kishekuļu wĩĩtu yesu kilisito 93</div> <div>maũndũ ma vata kuma 56</div> <div>asu asu asu asu 49</div> <div>ũu wĩ o vo 46</div> <div>ũvoo mũseo wa ũsumbĩ 42</div> |
| 5 | <div>maũndũ ma vata kuma ndetonĩ 53</div> <div>asu asu asu asu asu 36</div> <div>ĩndĩ o na ũu wĩ 35</div> <div>soma maelesyo ma mũthya namba 35</div> <div>ũsumbĩ wa ngai nĩ kyaũ 34</div> |

About HPLT Analytics

Volumes - Segments

Segments correspond to paragraph and list boundaries as defined by HTML elements (<p>, , , etc.) replaced by newlines.

Volumes - Tokens

Tokenized with <https://github.com/hplt-project/data-analytics-tool/blob/main/tokenizers-info.md>

Type-Token Ratio

Lexical variety computed as "number of types (uniques)/number of tokens", after removing punctuation (<https://www.slinfo.com/wp-content/uploads/2014/01/type-token-ratio.pdf>).

Document size (in segments)

Segments correspond to paragraph and list boundaries as defined by HTML elements (<p>, , , etc.) replaced by newlines.

Language distribution

Language identified with FastSpell (<https://github.com/mbanon/fastspell>).

Distribution of segments by fluency score

Obtained with Monocleaner (<https://github.com/bitextor/monocleaner>).

Distribution of documents by average fluency score

Obtained with Monocleaner (<https://github.com/bitextor/monocleaner>).

Distribution of documents by document score

Obtained with Web Docs Scorer (<https://github.com/pablop16n/web-docs-scorer/>).

Segment length distribution by token

Tokenized with <https://github.com/hplt-project/data-analytics-tool/blob/main/tokenizers-info.md>

Segment noise distribution

Obtained with Bicleaner Hardrules (<https://github.com/bitextor/bicleaner-hardrules/>).

Frequent n-grams

Tokenized with <https://github.com/hplt-project/data-analytics-tool/blob/main/tokenizers-info.md>, after removing n-grams starting or ending in a stopword. Stopwords from <https://github.com/hplt-project/data-analytics-tool/blob/main/scripts/resources/README.txt>