

General overview

Corpus	Analytics date	Language
lug_Latn.jsonl.tsv	9/19/2024	Ganda (lg)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
21,276	407,541	240,013 (58.89 %)	12M	65.91 MB	67,578,510

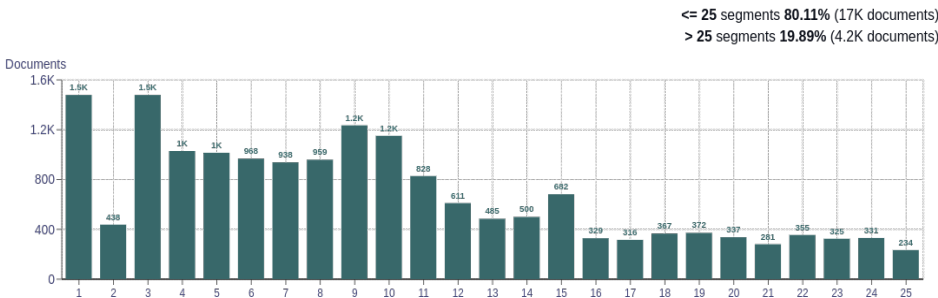
Top 10 domains

Domain	Docs	% of total
bukedde.co.ug	4.2K	19.81
bible.is	2.2K	10.37
dembefm.ug	1.9K	9.01
wordplanet.org	1.9K	8.72
wordproject.org	1.5K	7.13
jw.org	1.4K	6.71
radiosimba.ug	1.1K	5.17
cbsfm.ug	973	4.57
nbs.ug	576	2.71
wikipedia.org	394	1.85

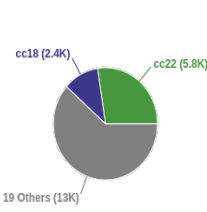
Top 10 TLDs

Domain	Docs	% of total
org	5.6K	26.53
co.ug	5.6K	26.34
ug	4.7K	22.05
com	2.7K	12.75
is	2.2K	10.37
net	165	0.78
info	36	0.17
eu	31	0.15
or.ug	22	0.10
ca	15	0.07

Documents size (in segments)

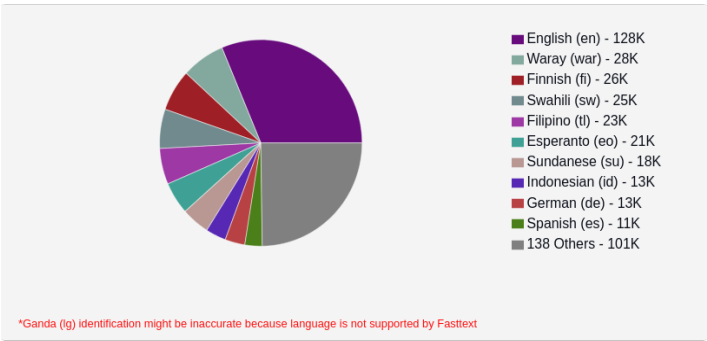


Documents by collection

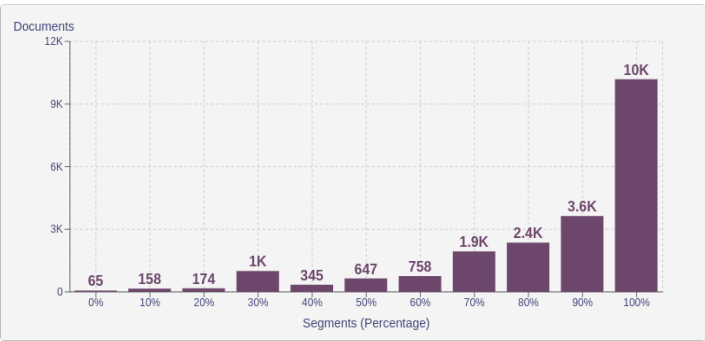


Language Distribution

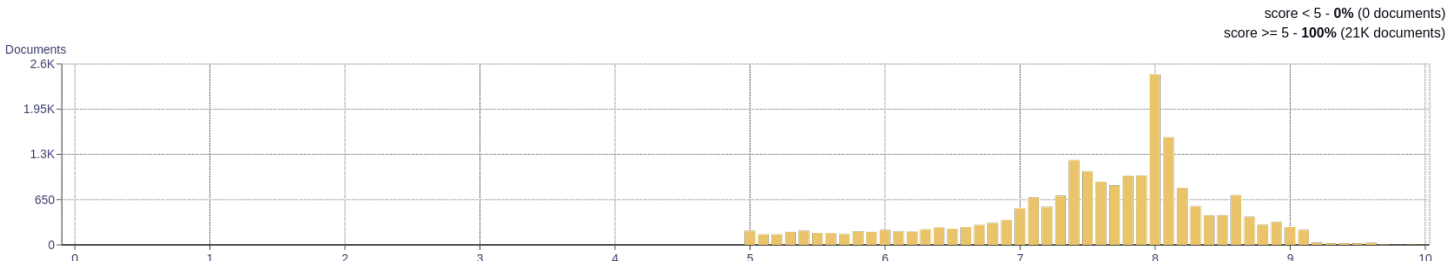
Number of segments



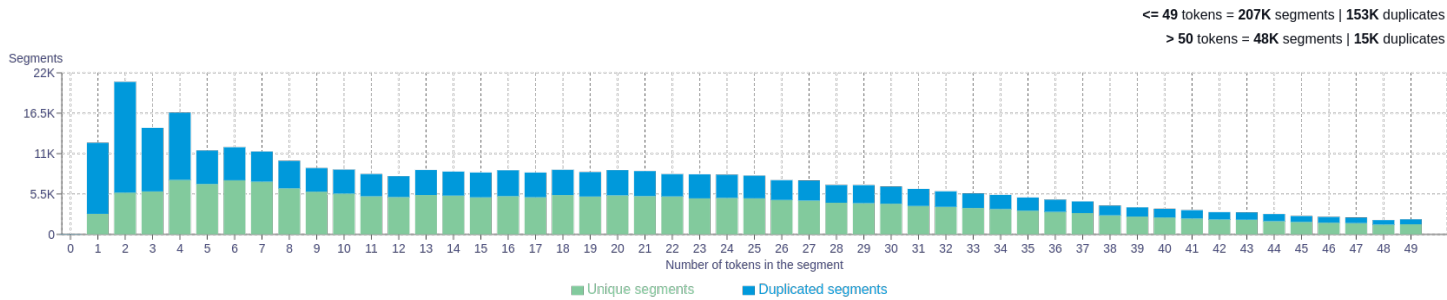
Percentage of segments in Ganda (lg) inside documents



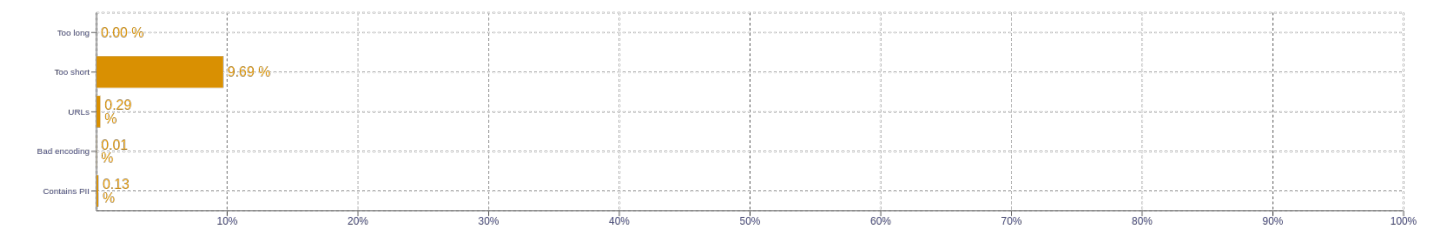
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>n   268953</div> <div>ng   57960</div> <div>omu   50175</div> <div>abantu   38919</div> <div>bw   36845</div>
2	<div>mukama katonda   4401</div> <div>wamu n   3041</div> <div>yesu kristo   2983</div> <div>mukama n   2680</div> <div>mukama waffe   2457</div>
3	<div>abaana ba isiraeri   1757</div> <div>maaso ga mukama   1352</div> <div>mukama katonda wo   1065</div> <div>ekigambo kya mukama   823</div> <div>mukama waffe yesu   679</div>
4	<div>bye bye bye bye   520</div> <div>empapula eziriko enyunzi ezigguka   449</div> <div>mukama waffe yesu kristo   434</div> <div>enyunzi ezigguka ku luno   391</div> <div>endagaano enkadde n'endagaano empya   389</div>
5	<div>bye bye bye bye bye   490</div> <div>eziriko enyunzi ezigguka ku luno   391</div> <div>do you have a story   254</div> <div>your community or an opinion   253</div> <div>you have a story in   253</div>

About HPLT Analytics

Volumes - Segments

Segments correspond to paragraph and list boundaries as defined by HTML elements (<p>, <ul>, <ol>, 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.sitinfo.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>, <ul>, <ol>, 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>