

General overview

Corpus	Analytics date	Language
kin_Latn.jsonl.tsv	9/23/2024	Kinyarwanda (rw)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
92,699	1,916,933	1,163,101 (60.68 %)	65M	354.87 MB	365,284,983

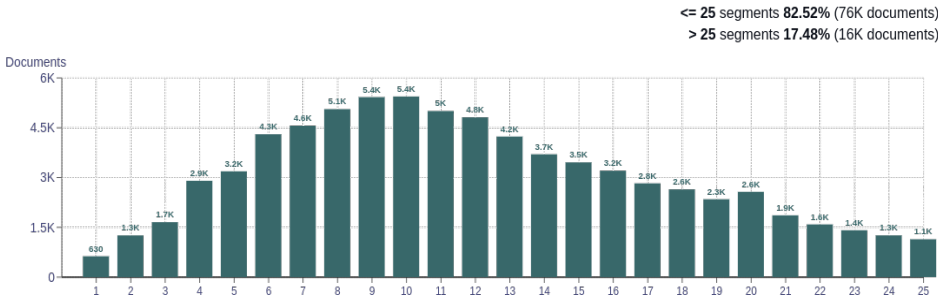
Top 10 domains

Domain	Docs	% of total
igihe.com	11K	11.38
kigalitoday.com	4.5K	4.88
agakiza.org	3.8K	4.09
jw.org	2.7K	2.87
yezu-akuzwe.org	2.6K	2.79
newsorwanda.com	2.4K	2.64
inyarwanda.com	2.3K	2.48
agasaro.com	2.3K	2.44
imirasire.com	2.1K	2.32
umuryango.rw	2K	2.18

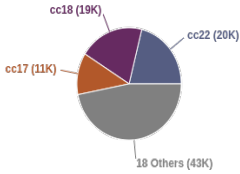
Top 10 TLDs

Domain	Docs	% of total
com	54K	58.54
org	17K	18.60
rw	13K	13.82
fr	2.3K	2.43
co.rw	1.5K	1.64
gov.rw	1.4K	1.48
net	1K	1.11
info	453	0.49
be	232	0.25
ca	215	0.23

Documents size (in segments)

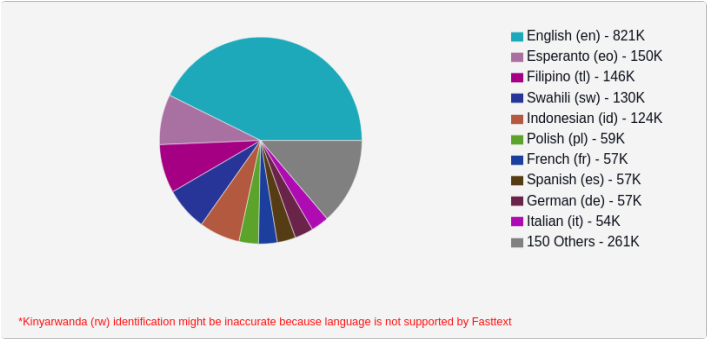


Documents by collection

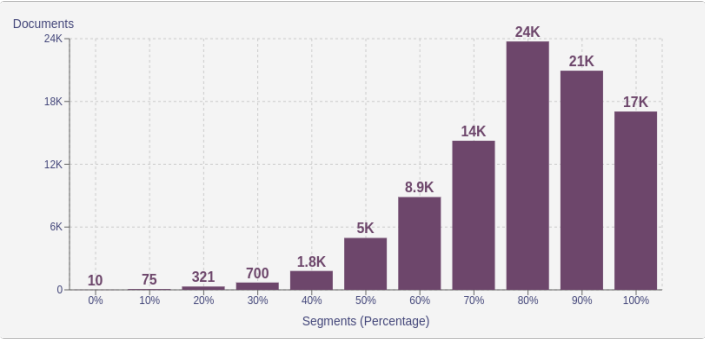


Language Distribution

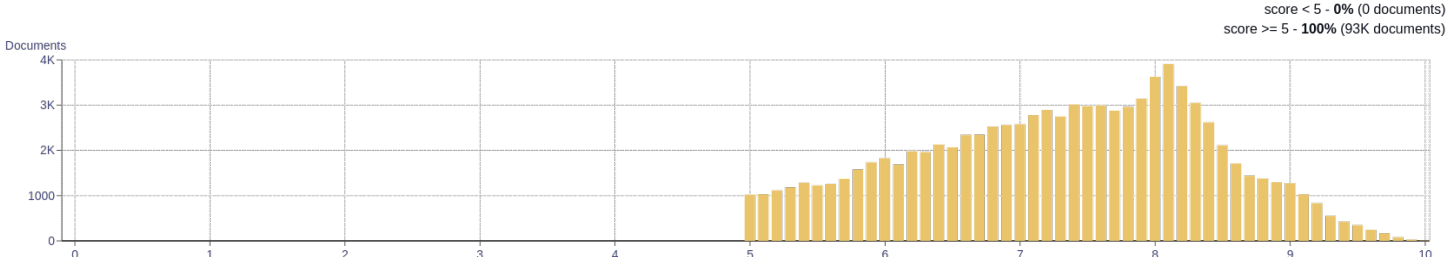
Number of segments



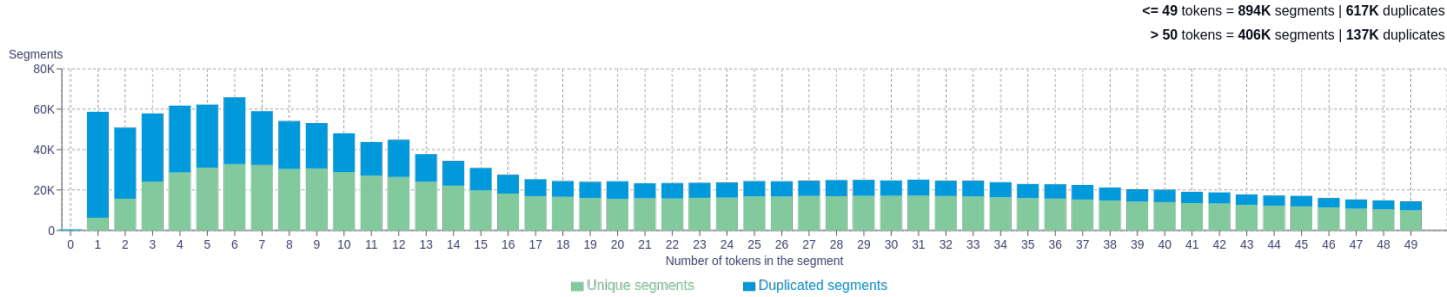
Percentage of segments in Kinyarwanda (rw) inside documents



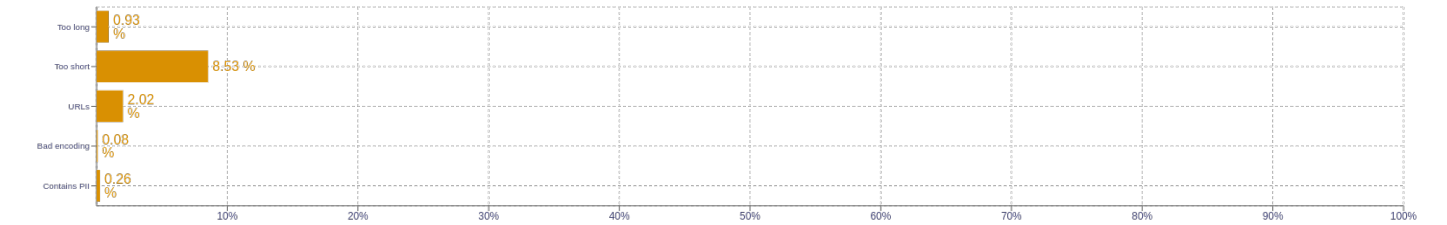
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div><div>n   1103217</div><div>y   455109</div><div>kandi   364482</div><div>imana   319325</div><div>w   249241</div></div>
2	<div><div>u rwanda   56569</div><div>ndetse n   41497</div><div>nyuma y   26384</div><div>cyangwa se   25683</div><div>cyane cyane   21242</div></div>
3	<div><div>kanda hano umusubize   17580</div><div>hirya no hino   8877</div><div>rimwe na rimwe   8137</div><div>jenoside yakorewe abatutsi   6784</div><div>mujyi wa kigali   5388</div></div>
4	<div><div>bite bite bite bite   3498</div><div>imana imuhe amahoro n   2147</div><div>zunze ubumwe za amerika   1721</div><div>allah amuhe amahoro n   1355</div><div>ushinzwe imibereho myiza y   1324</div></div>
5	<div><div>bite bite bite bite bite   3493</div><div>hirya no hino ku isi   1735</div><div>leta zunze ubumwe za amerika   1700</div><div>hirya no hino mu gihugu   1581</div><div>alayihi wa aalih wa sallam   1249</div></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.sltinfo.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>