

General overview

Corpus	Analytics date	Language
kbp_Latn.jsonl.tsv	12/5/2024	Kabiyè (kbp)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
7,075	46,792	26,333 (56.28 %)	5.3M	24.47 MB	20,864,671

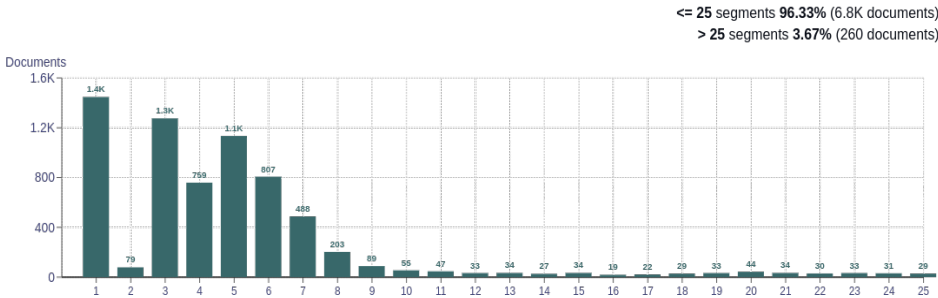
Top 10 domains

Domain	Docs	% of total
wikipedia.org	3.8K	53.87
bible.is	2.4K	34.57
jw.org	488	6.90
ebible.org	89	1.26
bibles.org	84	1.19
breakeveryyoke.com	33	0.47
wikiplanet.click	28	0.40
bywiki.com	12	0.17
revue-gugu.org	9	0.13
unicode.org	6	0.08

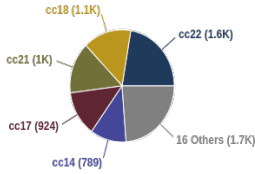
Top 10 TLDs

Domain	Docs	% of total
org	4.5K	63.70
is	2.4K	34.57
com	71	1.00
click	28	0.40
net	13	0.18
cf	6	0.08
ca	2	0.03
fr	1	0.01
nl	1	0.01

Documents size (in segments)

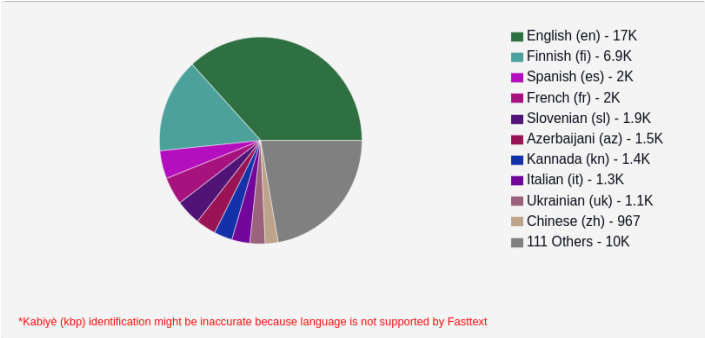


Documents by collection

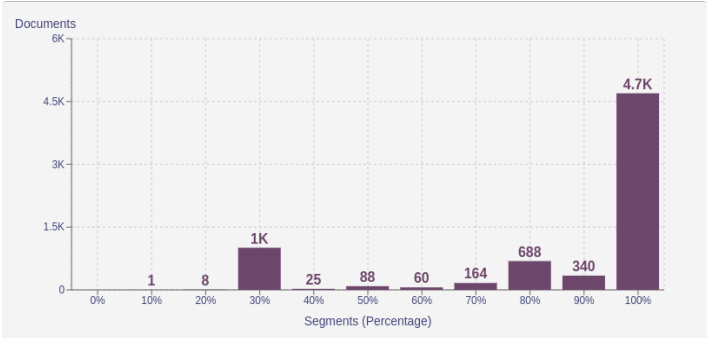


Language Distribution

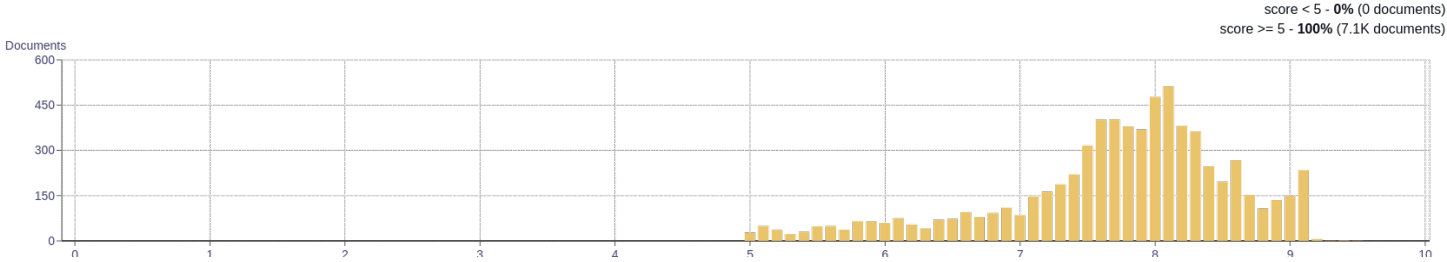
Number of segments



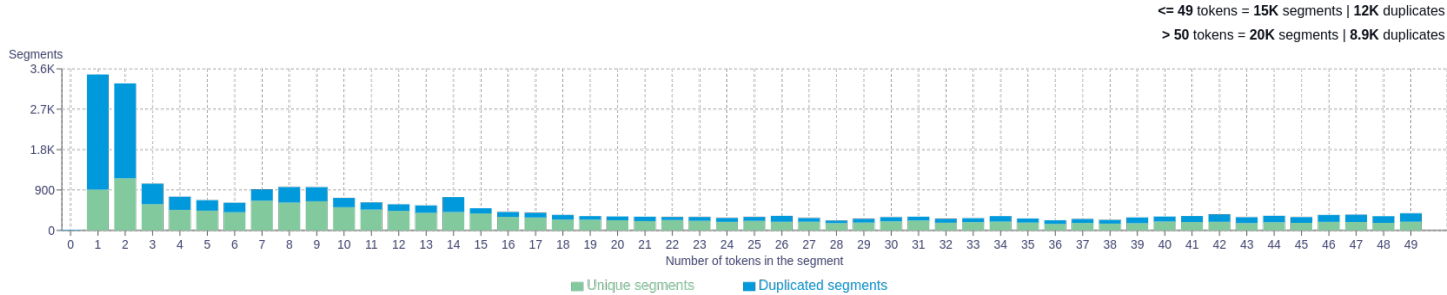
Percentage of segments in Kabiyè (kbp) inside documents



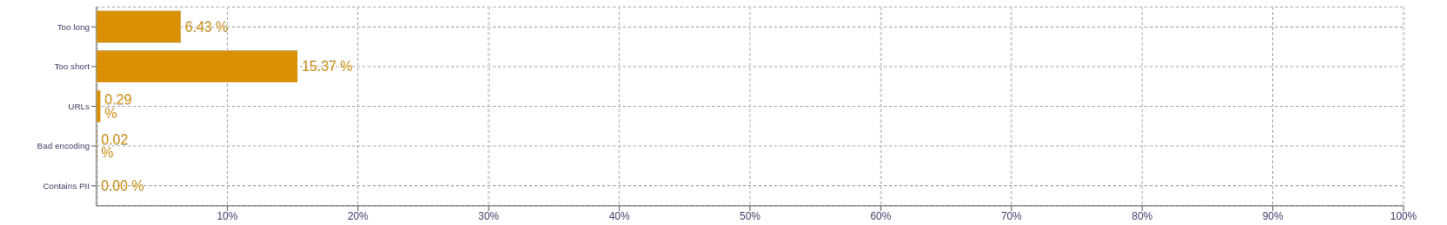
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div><div>ı 91417</div><div>taa 78909</div><div>ba 64121</div><div>sı 53301</div><div>ma 53182</div></div>
2	<div><div>hu aa 5303</div><div>dı ba 4766</div><div>dı v 3980</div><div>εjaε taa 3948</div><div>te tın 3921</div></div>
3	<div><div>nala hu aa 1545</div><div>pə taɣa pʊlʊ 1492</div><div>εjaε dıne dı 1242</div><div>hu buloɣ aa 1062</div><div>v a baa 956</div></div>
4	<div><div>tıya ba a baa 687</div><div>pə taɣa pʊlʊ tɔɔ 687</div><div>tıya v a baa 608</div><div>yee yee yee yee 570</div><div>re yesu bası tıya 432</div></div>
5	<div><div>bası tıya ba a baa 639</div><div>bası tıya v a baa 604</div><div>yee yee yee yee yee 568</div><div>kalʊ na ba tı tın 327</div><div>avr mai juın juıl aou̇t 311</div></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.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>, , , 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>