

General overview

Corpus	Analytics date	Language
HPLT-v2-bos_Latn.tsv	9/22/2024	Bosnian (bs)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
14,613,088	268,156,601			44.01 GB	45,817,404,394

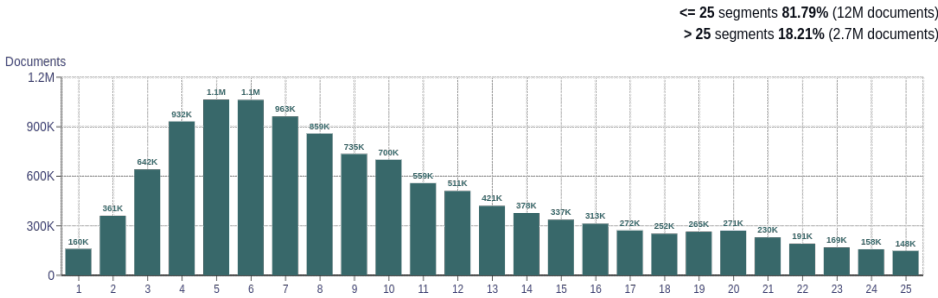
Top 10 domains

Domain	Docs	% of total
klix.ba	372K	2.55
wikipedia.org	273K	1.87
sportske.net	184K	1.26
vesti.rs	181K	1.24
slobodnaevropa.org	174K	1.19
blogspot.com	155K	1.06
blic.rs	140K	0.96
krstarica.com	113K	0.78
b92.net	108K	0.74
mondo.rs	105K	0.72

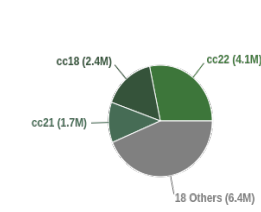
Top 10 TLDs

Domain	Docs	% of total
com	4.4M	29.93
rs	3.9M	26.89
ba	2M	13.45
net	1.3M	8.78
org	891K	6.10
info	501K	3.43
me	461K	3.16
hr	195K	1.33
org.rs	164K	1.12
co.rs	141K	0.96

Documents size (in segments)

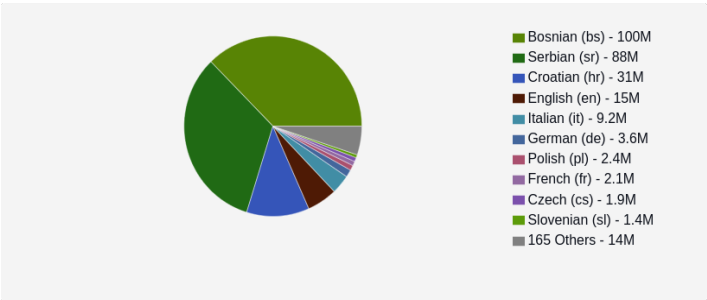


Documents by collection

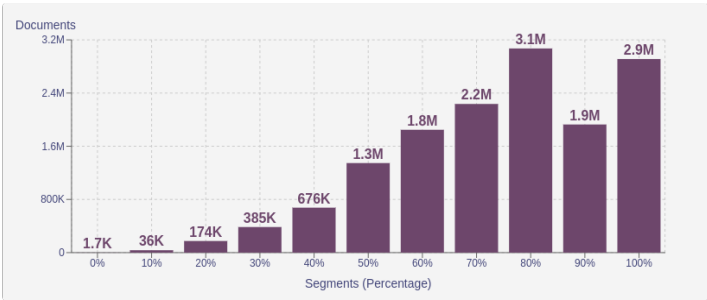


Language Distribution

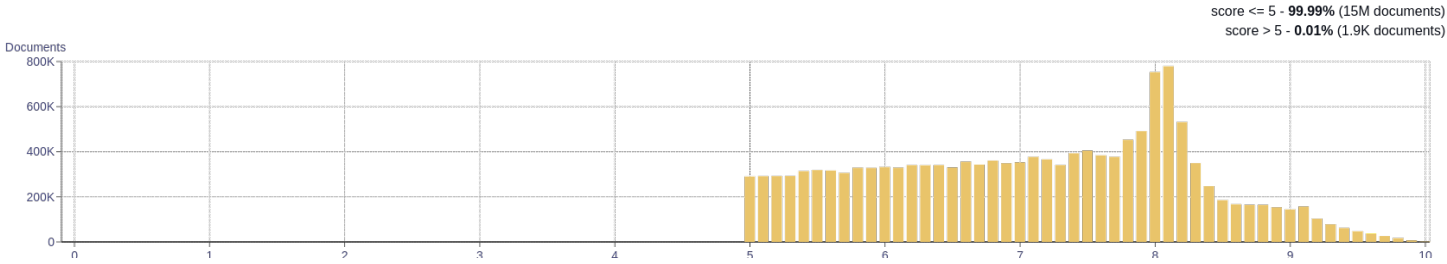
Number of segments



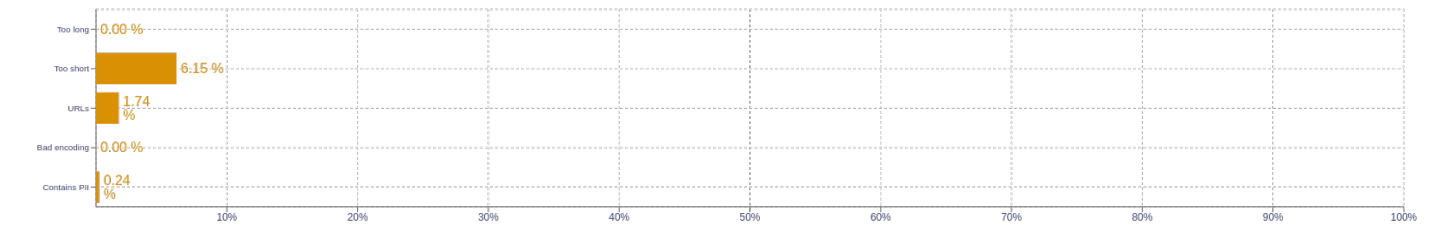
Percentage of segments in Bosnian (bs) inside documents



Distribution of documents by document score



Segment noise distribution



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>