

General overview

Corpus	Analytics date	Language
HPLT-v2-slk_Latn.tsv	9/24/2024	Slovak (sk)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
21,827,259	494,276,634			70.65 GB	69,899,657,544

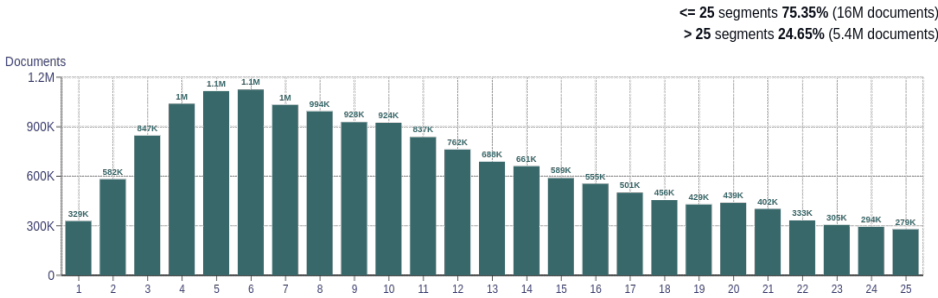
Top 10 domains

Domain	Docs	% of total
sme.sk	1.3M	6.01
wikipedia.org	360K	1.65
firebaseapp.com	291K	1.33
web.app	289K	1.32
pravda.sk	222K	1.02
zdravie.sk	188K	0.86
blogspot.com	161K	0.74
24hod.sk	147K	0.67
aktuality.sk	142K	0.65
netky.sk	118K	0.54

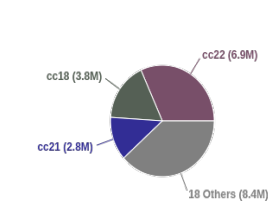
Top 10 TLDs

Domain	Docs	% of total
sk	17M	77.32
com	2.1M	9.44
cz	665K	3.05
eu	628K	2.88
org	526K	2.41
app	289K	1.33
net	212K	0.97
info	122K	0.56
ru	31K	0.14
xyz	26K	0.12

Documents size (in segments)

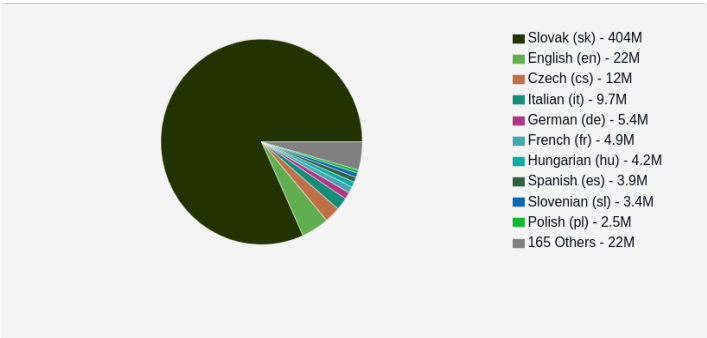


Documents by collection

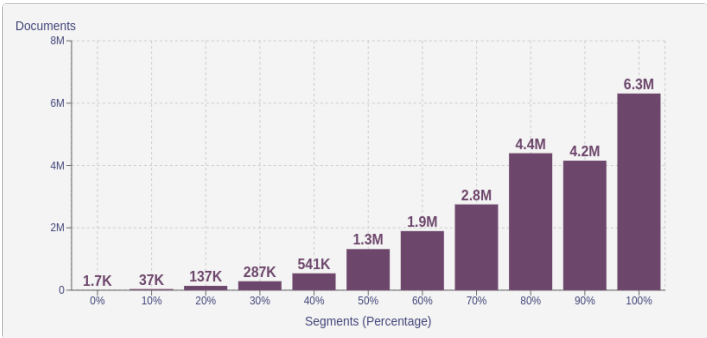


Language Distribution

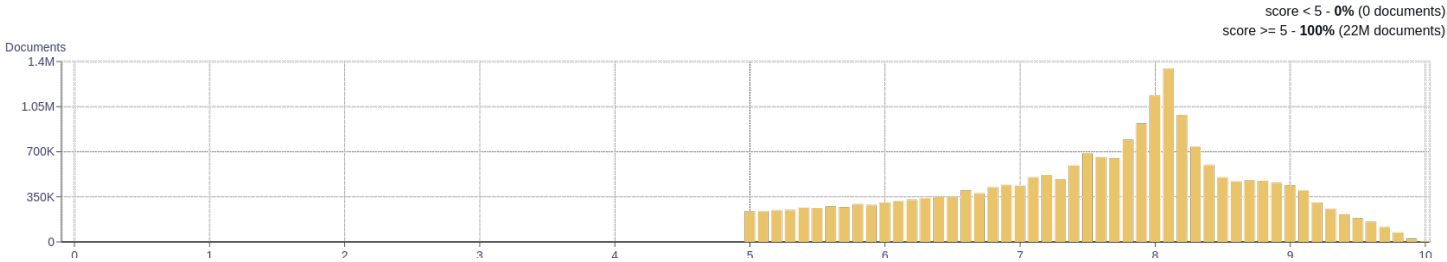
Number of segments



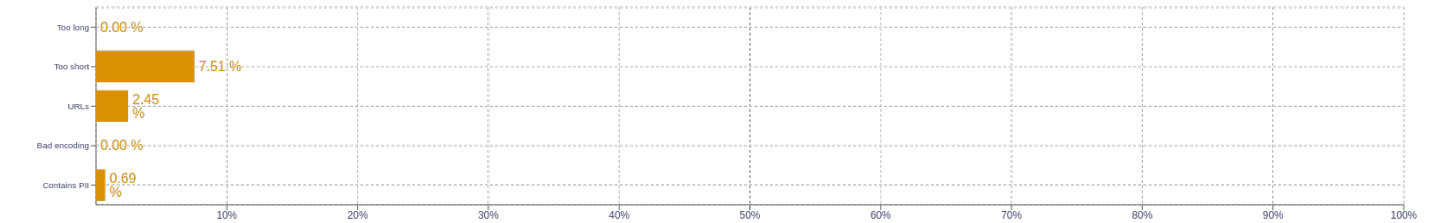
Percentage of segments in Slovak (sk) inside documents



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>