

General overview

Corpus	Analytics date	Language
HPLT-docslite.sl.tsv	6/8/2024	Slovenian (sl)

Volumes

Docs	Segments	Unique segments	Tokens	Size
2,196,149	285,527,540	72,459 (0.03 %)	3B	15.65 GB

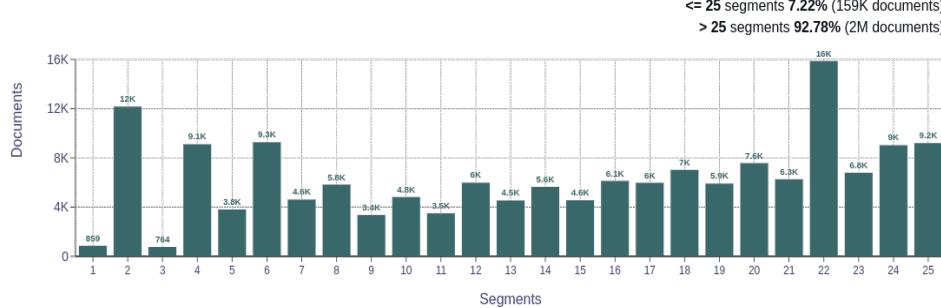
Top 10 domains

Domain	Docs	% of total
diebuchsuche.com	61K	2.78
blogspot.si	34K	1.54
metropolitan.si	30K	1.37
europages.si	26K	1.20
delo.si	22K	0.98
wikipedia.org	19K	0.85
agoda.com	17K	0.79
blogspot.com	16K	0.73
sta.si	14K	0.64
siol.net	14K	0.64

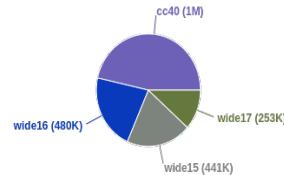
Top 10 TLDs

Domain	Docs	% of total
si	1.3M	61.31
com	484K	22.05
net	105K	4.76
org	85K	3.89
eu	61K	2.77
info	26K	1.16
tv	5.5K	0.25
de	5.3K	0.24
at	5.3K	0.24
news	4.7K	0.21

Documents size (in segments)

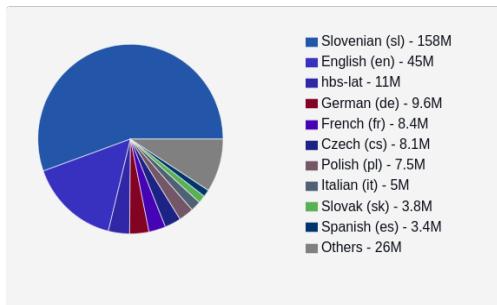


Documents by collection

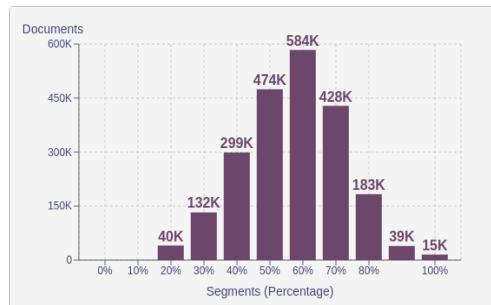


Language Distribution

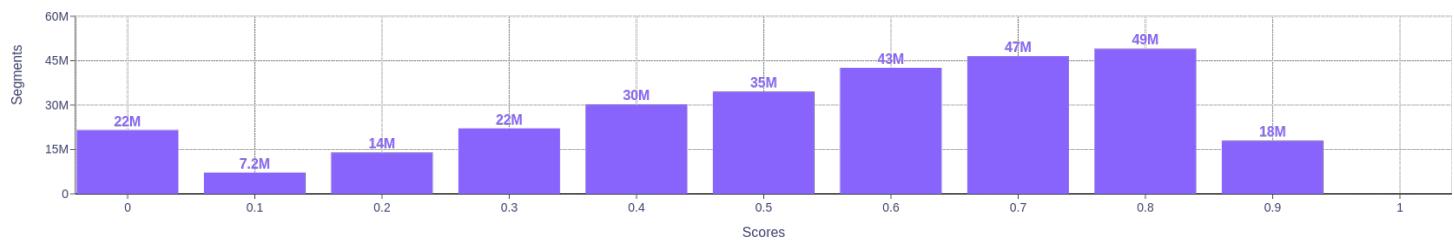
Number of segments



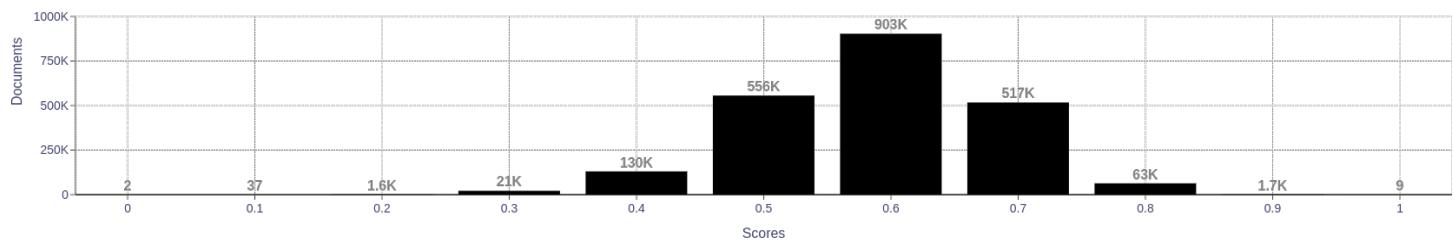
Percentage of segments in Slovenian (sl) inside documents



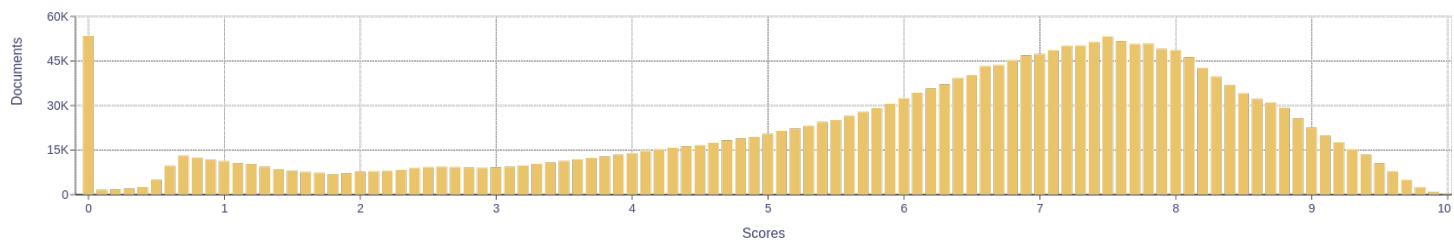
Distribution of segments by fluency score



Distribution of documents by average fluency score



Distribution of documents by document score



Segment length distribution by token

<= 49 tokens = 52M segments | 221M duplicates

> 50 tokens = 13M segments | 3.3M duplicates



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>