

General overview

Corpus	Analytics date	Language
HPLT-docslite.cs.tsv	6/16/2024	Czech (cs)

Volumes

Docs	Segments	Unique segments	Tokens	Size
16,987,093	2,096,798,338	163,737 (0.01 %)	23B	127.38 GB

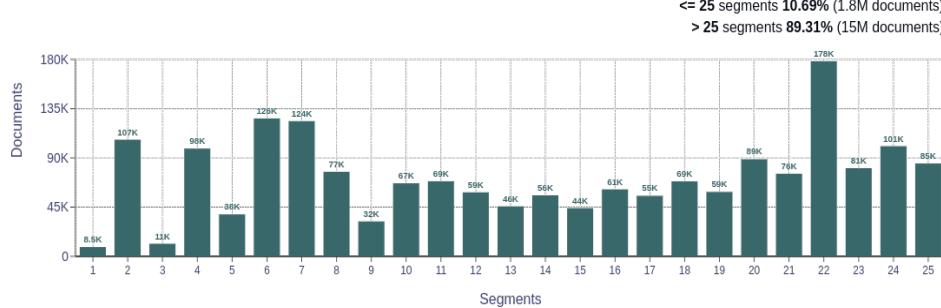
Top 10 domains

Domain	Docs	% of total
idnes.cz	360K	2.12
docplayer.cz	259K	1.53
blogspot.cz	251K	1.48
sleviste.cz	163K	0.96
edb.cz	153K	0.90
karaoketexty.cz	149K	0.88
web.app	122K	0.72
firebaseapp.com	116K	0.68
diebuchsueche.com	114K	0.67
ju8.me	95K	0.56

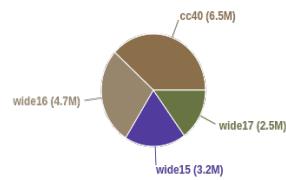
Top 10 TLDs

Domain	Docs	% of total
cz	14M	79.84
com	1.5M	8.87
eu	487K	2.87
net	308K	1.81
org	273K	1.60
info	150K	0.88
app	122K	0.72
sk	122K	0.72
me	98K	0.58
de	30K	0.17

Documents size (in segments)

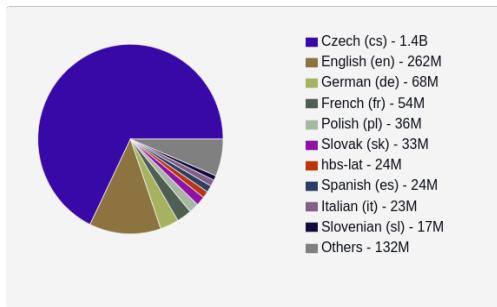


Documents by collection

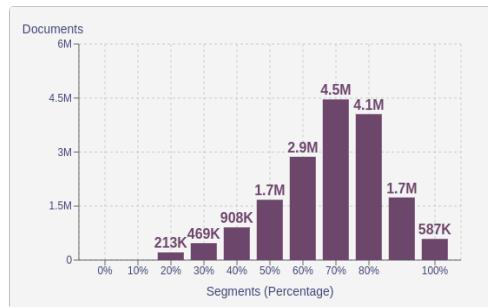


Language Distribution

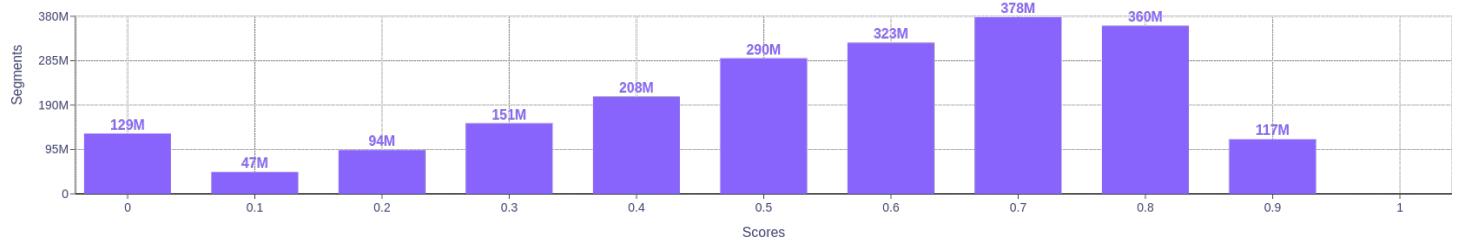
Number of segments



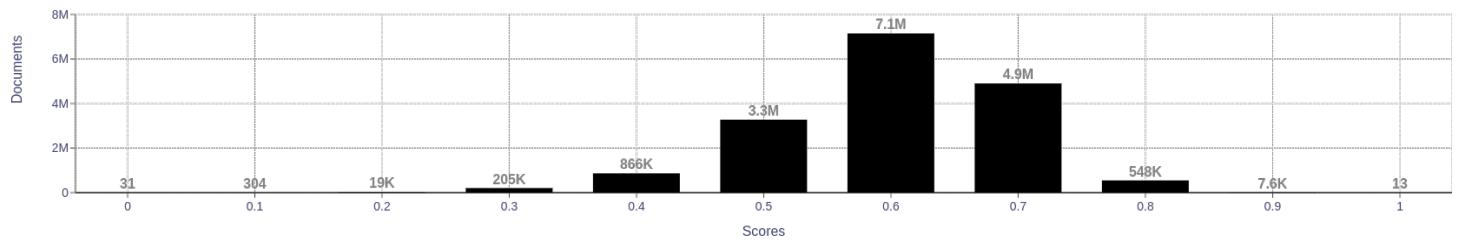
Percentage of segments in Czech (cs) inside documents



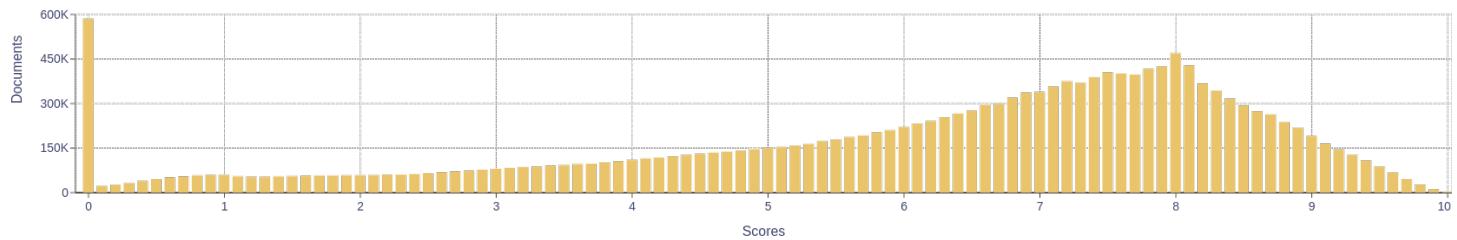
Distribution of segments by fluency score



Distribution of documents by average fluency score



Distribution of documents by document score



Segment length distribution by token

<= 49 tokens = 360M segments | 1.6B duplicates

> 50 tokens = 90M segments | 27M 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>