

General overview

Corpus	Analytics date	Language
HPLT-docsite.hu.tsv	6/13/2024	Hungarian (hu)

Volumes

Docs	Segments	Unique segments	Tokens	Size
11,708,768	1,534,818,521	172,301 (0.01 %)	18B	105.54 GB

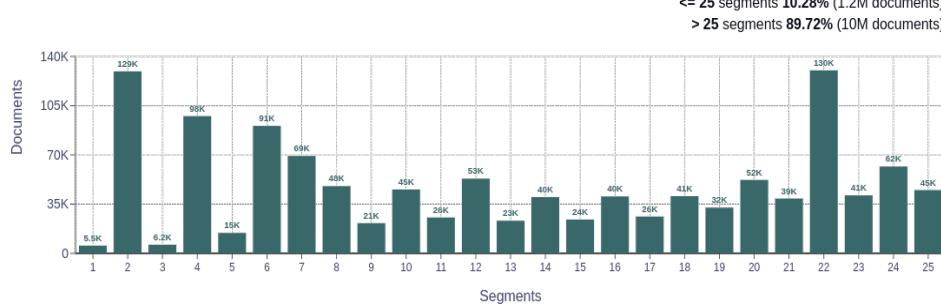
Top 10 domains

Domain	Docs	% of total
blogspot.hu	739K	6.31
blog.hu	235K	2.01
docplayer.hu	142K	1.21
diebuchsuche.com	115K	0.98
blogspot.ro	113K	0.97
blogspot.com	76K	0.65
jú8.me	67K	0.58
lap.hu	59K	0.50
24.hu	49K	0.42
lightinthebox.com	48K	0.41

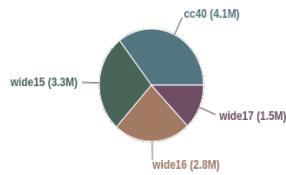
Top 10 TLDs

Domain	Docs	% of total
hu	8.6M	73.31
com	1.5M	12.53
ro	240K	2.05
org	181K	1.55
net	178K	1.52
eu	177K	1.51
info	145K	1.24
me	73K	0.63
sk	73K	0.62
co.hu	57K	0.49

Documents size (in segments)

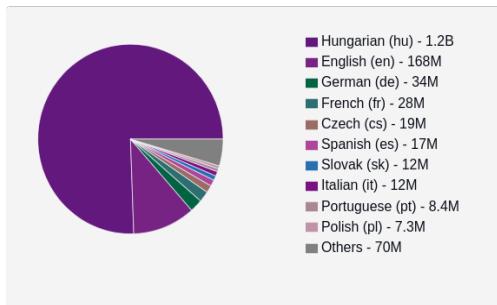


Documents by collection

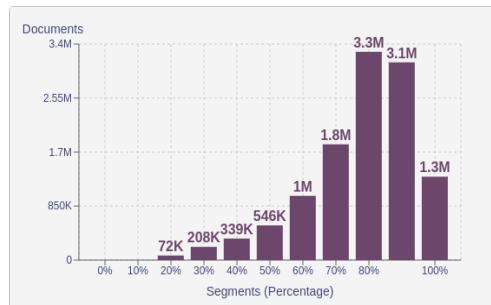


Language Distribution

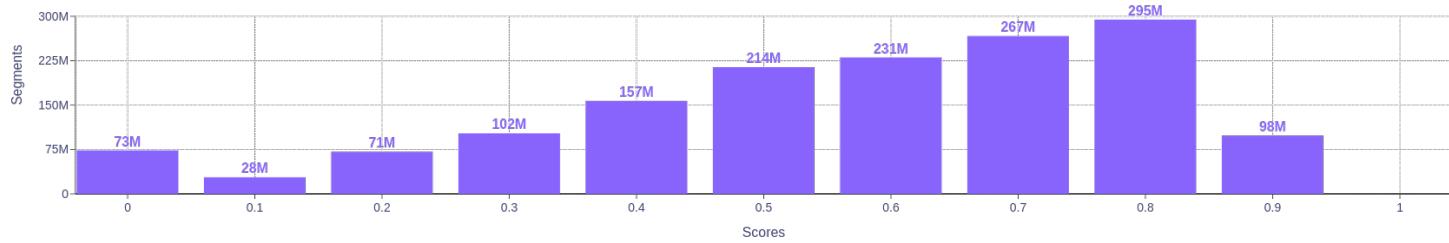
Number of segments



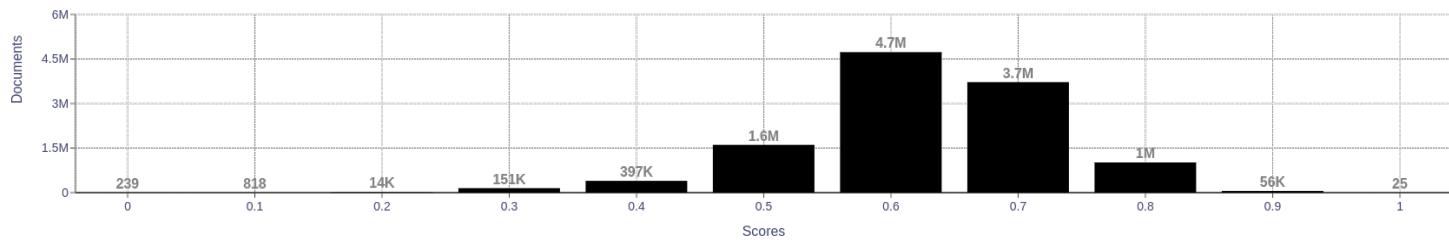
Percentage of segments in Hungarian (hu) inside documents



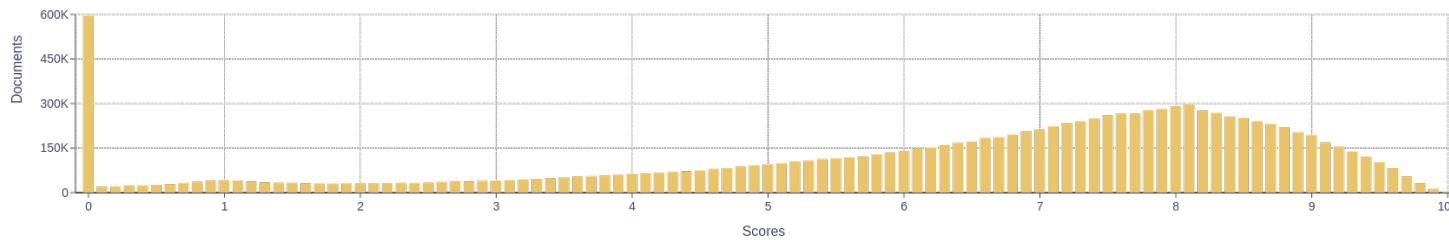
Distribution of segments by fluency score



Distribution of documents by average fluency score



Distribution of documents by document score



Segment length distribution by token

<= 49 tokens = 275M segments | 1.2B duplicates

> 50 tokens = 70M segments | 23M 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>