

General overview

Corpus	Analytics date	Language
HPLT-docslite.el.tsv	6/30/2024	Greek (el)

Volumes

Docs	Segments	Unique segments	Tokens	Size
15,833,090	4,557,166,094	242,639 (0.01 %)	43B	311.27 GB

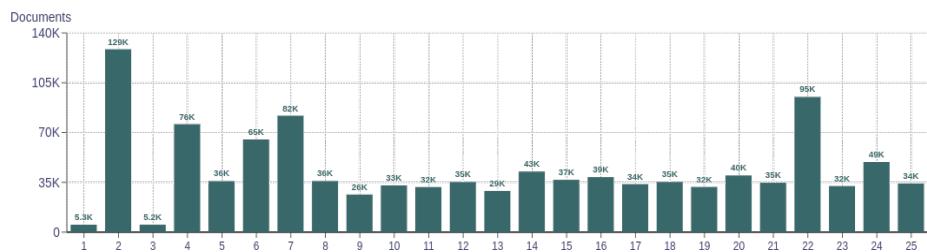
Top 10 domains

Domain	Docs	% of total
blogspot.gr	2.4M	15.37
blogspot.com	785K	4.96
wordpress.com	216K	1.36
blogspot.be	183K	1.15
docplayer.gr	156K	0.99
inewsg.com	137K	0.86
blogspot.nl	121K	0.77
diebuchsue.com	118K	0.75
blogspot.ch	89K	0.56
rotise.gr	86K	0.54

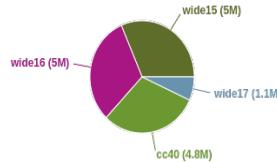
Top 10 TLDs

Domain	Docs	% of total
gr	10M	64.02
com	3.1M	19.30
org	212K	1.34
be	186K	1.18
eu	170K	1.08
net	165K	1.04
com.cy	164K	1.04
nl	125K	0.79
de	97K	0.61
ch	96K	0.60

Documents size (in segments)

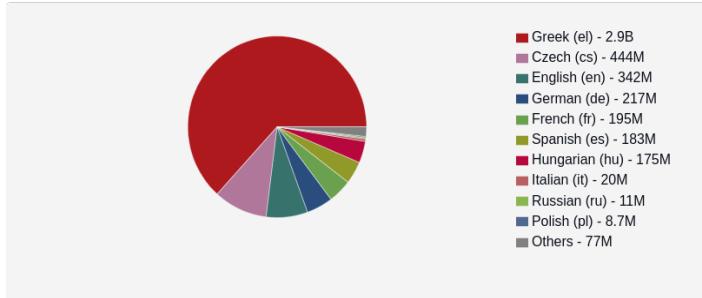


Documents by collection

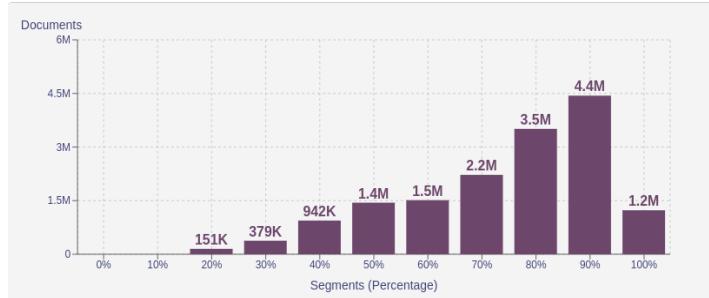


Language Distribution

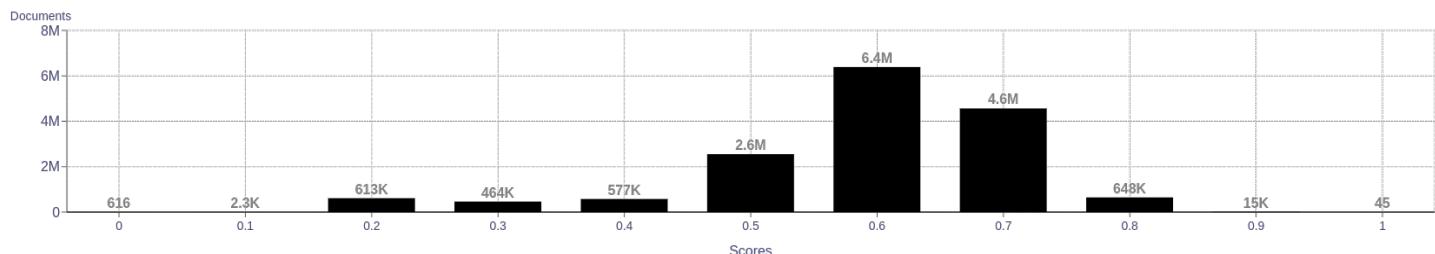
Number of segments



Percentage of segments in Greek (el) inside documents

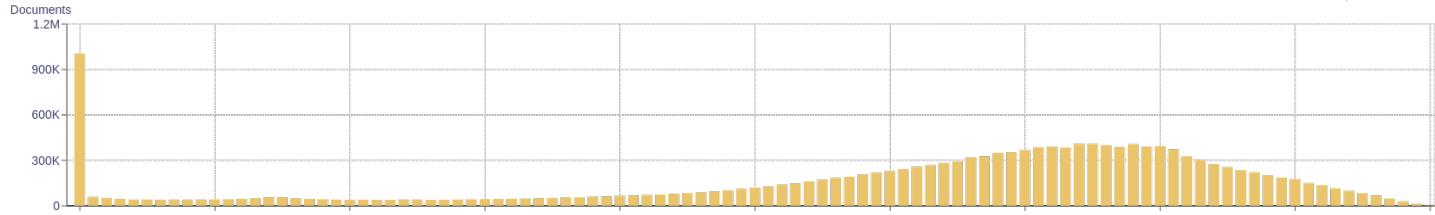


Distribution of documents by average fluency score



Distribution of documents by document score

score <= 5 - 22.88% (3.6M documents)
score > 5 - 77.12% (12M documents)



Segment length distribution by token

<= 49 tokens = 319M segments | 4.1B duplicates
> 50 tokens = 114M segments | 52M 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/mabanon/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>