

General overview

Corpus	Analytics date	Language
HPLT-v2-ita_Latn.tsv	10/6/2024	Italian (it)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
221,752,424	5,127,273,785			769.93 GB	815,696,448,525

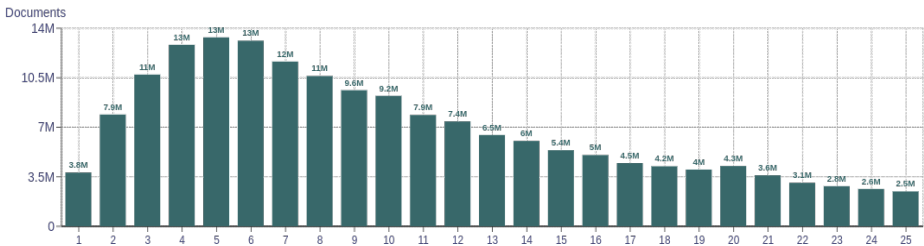
Top 10 domains

Domain	Docs	% of total
blogspot.com	6.2M	2.78
blogspot.it	3.8M	1.71
wordpress.com	2.3M	1.02
wikipedia.org	2M	0.90
kijiji.it	1.6M	0.70
repubblica.it	956K	0.43
alternativa.org	814K	0.37
corriere.it	759K	0.34
tripadvisor.it	694K	0.31
docplayer.it	637K	0.29

Top 10 TLDs

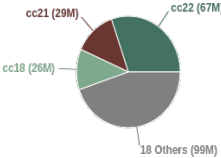
Domain	Docs	% of total
it	125M	56.38
com	59M	26.76
org	9.9M	4.48
net	7.5M	3.36
eu	3.4M	1.52
info	2.5M	1.13
ch	2M	0.90
tv	816K	0.37
biz	424K	0.19
de	412K	0.19

Documents size (in segments)



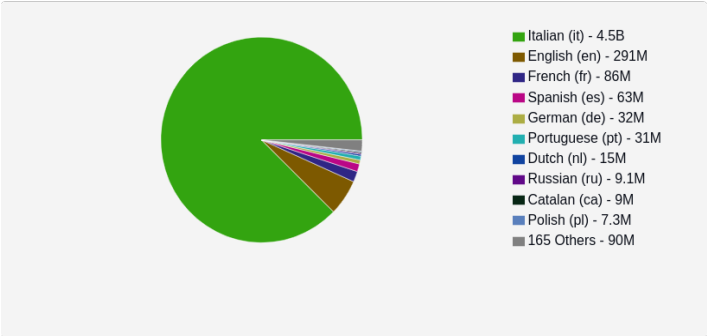
<= 25 segments 77.86% (173M documents)
> 25 segments 22.14% (49M documents)

Documents by collection

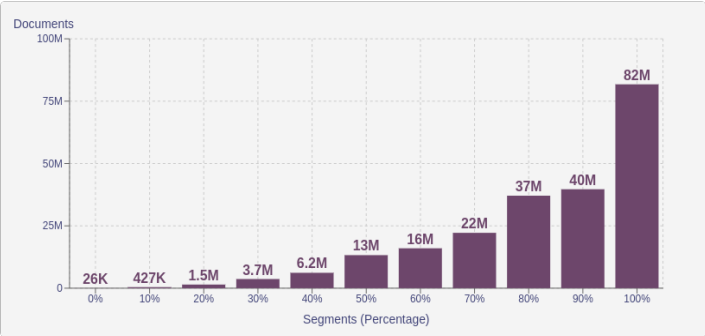


Language Distribution

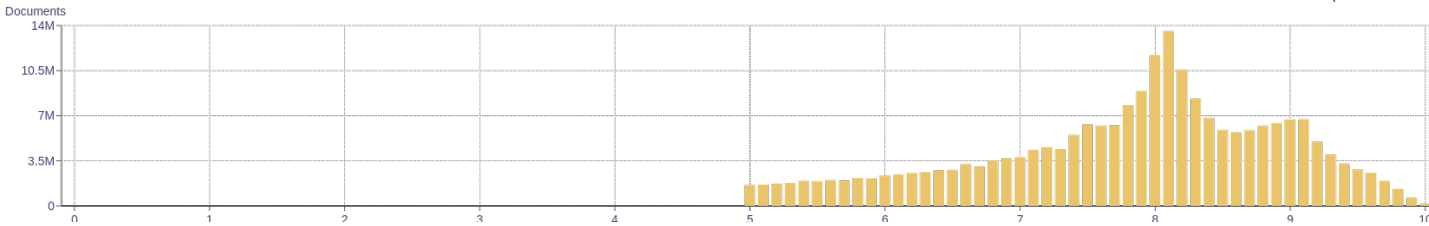
Number of segments



Percentage of segments in Italian (it) inside documents

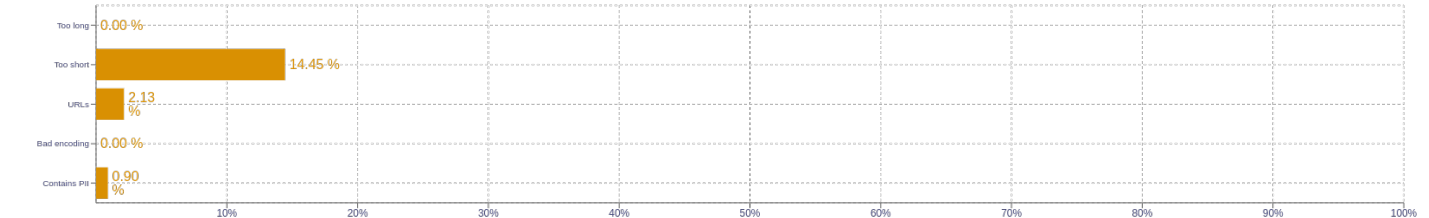


Distribution of documents by document score



score < 5 - 0% (0 documents)
score >= 5 - 100% (222M documents)

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>