

General overview

Corpus	Analytics date	Language
HPLT-docslite.zh.tsv	8/19/2024	Chinese (zh)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
1,080,799,436	162,660,044,735			11.26 TB	5,184,542,321,492

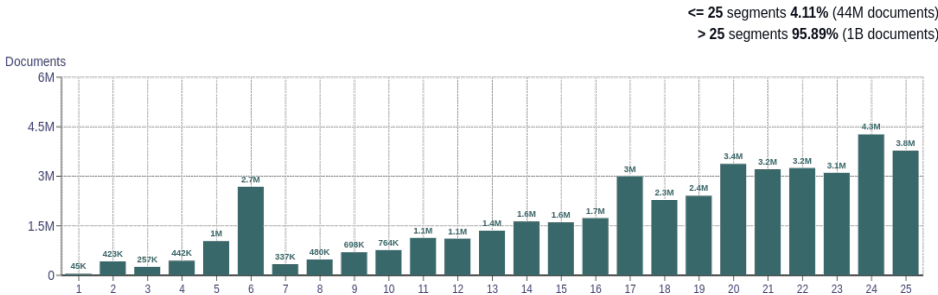
Top 10 domains

Domain	Docs	% of total
58.com	7.6M	0.70
b2b168.com	5.1M	0.47
woaifenxiang.net	4.8M	0.44
68mtv.com	4M	0.37
baixing.com	3.9M	0.36
hc360.com	3.8M	0.35
y40584.cn	3.6M	0.33
baidu.com	3.2M	0.30
ganji.com	3M	0.28
kushubao.com	3M	0.28

Top 10 TLDs

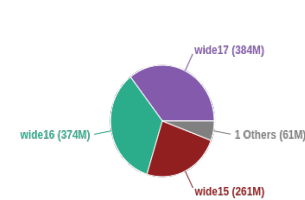
Domain	Docs	% of total
com	691M	63.95
cn	136M	12.59
com.cn	50M	4.61
net	46M	4.23
cc	27M	2.51
club	17M	1.56
org	14M	1.31
win	9.7M	0.90
top	7.1M	0.66
gov.cn	7M	0.65

Documents size (in segments)



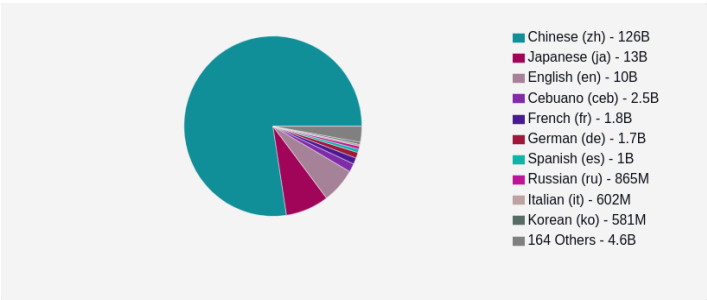
<= 25 segments 4.11% (44M documents)
> 25 segments 95.89% (1B documents)

Documents by collection

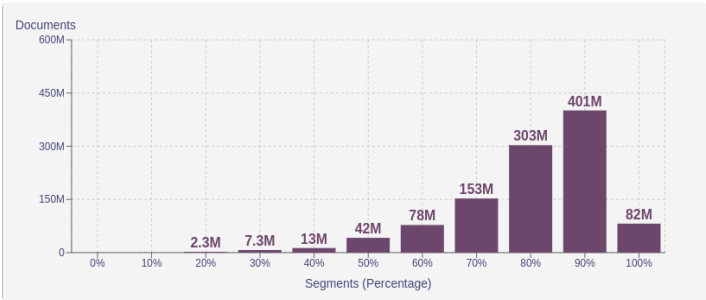


Language Distribution

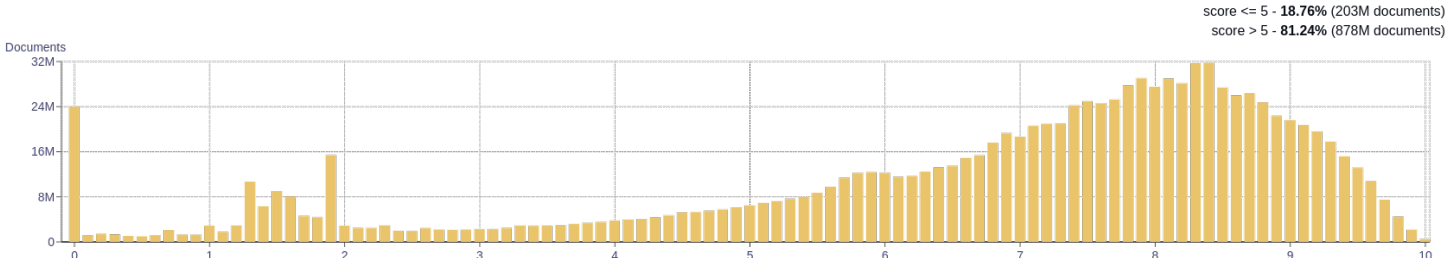
Number of segments



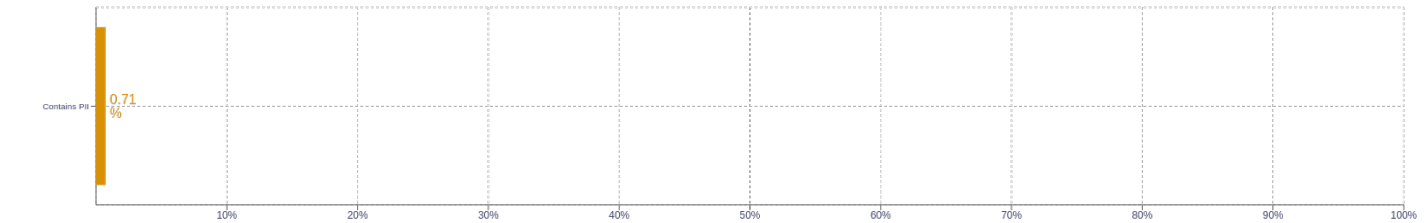
Percentage of segments in Chinese (zh) inside documents



Distribution of documents by document score



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>