

General overview

Corpus	Analytics date	Language
HPLT-docslite.hi.tsv	6/9/2024	Hindi (hi)

Volumes

Docs	Segments	Unique segments	Tokens	Size
5,774,861	855,081,850	150,107 (0.02 %)	9.2B	81.68 GB

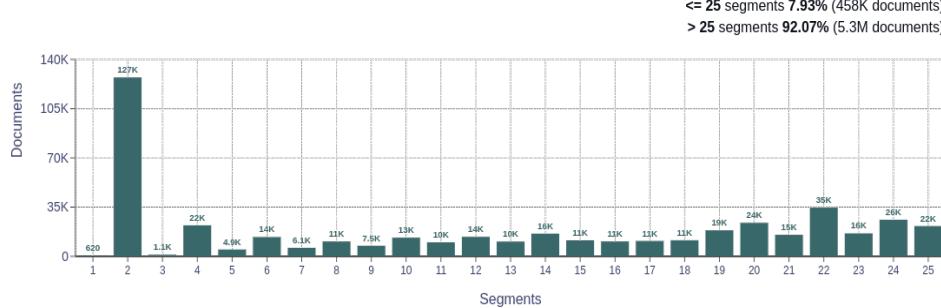
Top 10 domains

Domain	Docs	% of total
alibaba.com	519K	8.99
diebuchsueche.com	404K	6.99
blogspot.in	190K	3.30
lyricsparoles.com	106K	1.84
ju8.me	87K	1.51
blogspot.com	73K	1.26
indiatimes.com	70K	1.21
jagran.com	43K	0.74
zipcodecountry.com	35K	0.60
fanpop.com	28K	0.48

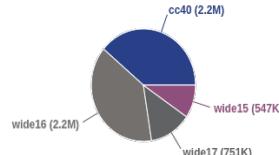
Top 10 TLDs

Domain	Docs	% of total
.com	4.1M	71.51
.in	840K	14.55
.org	195K	3.38
.me	90K	1.57
.co.in	82K	1.42
.page	79K	1.37
.net	74K	1.28
.ae	23K	0.40
.ru	21K	0.37
.co	19K	0.33

Documents size (in segments)

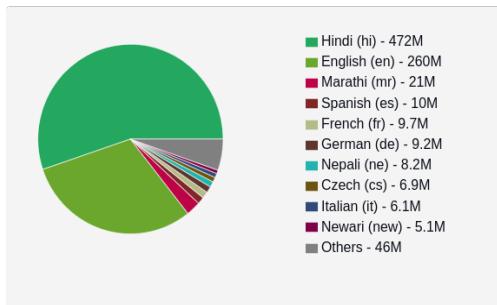


Documents by collection

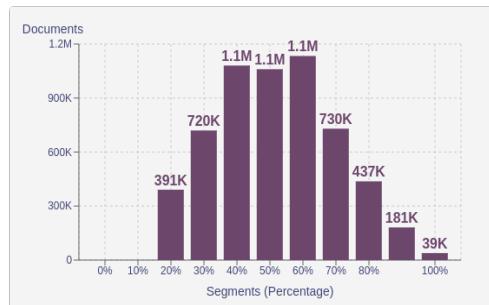


Language Distribution

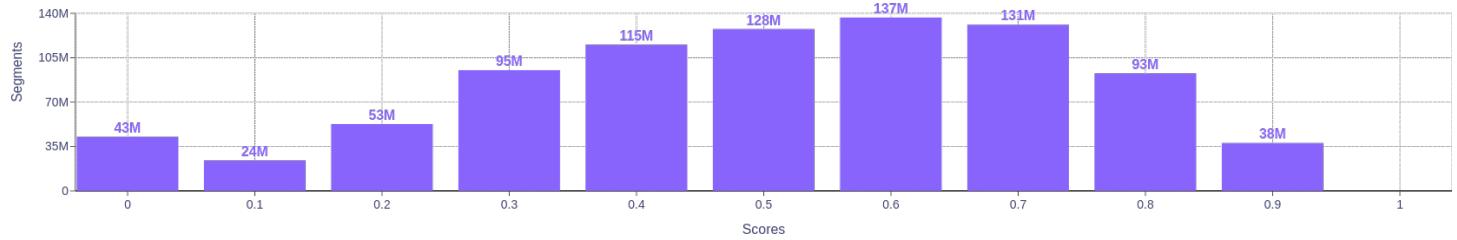
Number of segments



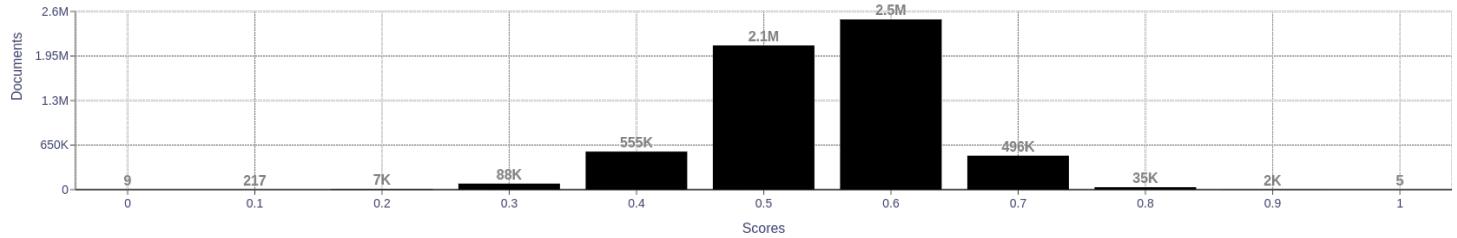
Percentage of segments in Hindi (hi) inside documents



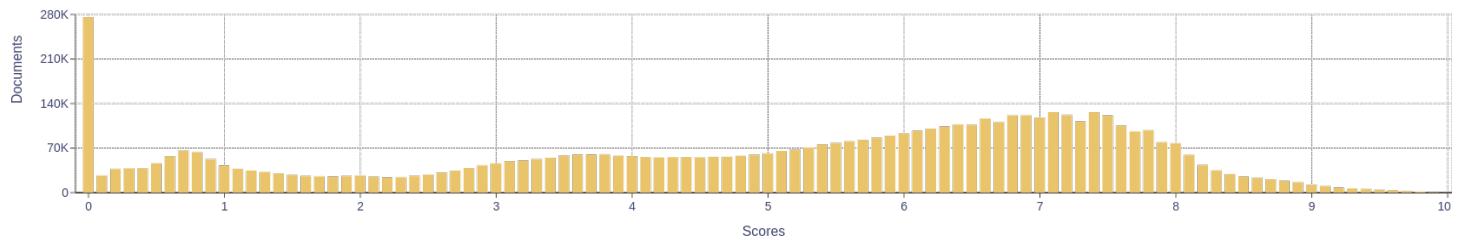
Distribution of segments by fluency score



Distribution of documents by average fluency score



Distribution of documents by document score



Segment length distribution by token

<= 49 tokens = 143M segments | 684M duplicates

> 50 tokens = 29M segments | 7.8M 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>