

General overview

Corpus	Analytics date	Language
HPLT-docslite.ur.tsv	6/8/2024	Urdu (ur)

Volumes

Docs	Segments	Unique segments	Tokens	Size
1,437,830	136,409,453	55,026 (0.04 %)	1.7B	10.61 GB

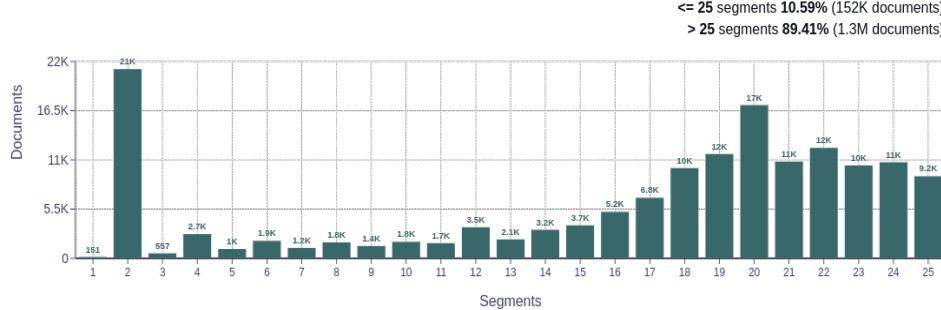
Top 10 domains

Domain	Docs	% of total
diebuchsuche.com	558K	38.83
lyricsparoles.com	119K	8.27
fanpop.com	22K	1.54
wikipedia.org	13K	0.90
dunyapakistan.com	10K	0.70
express.pk	9.7K	0.67
news18.com	9K	0.63
dailypakistan.com.pk	8.9K	0.62
urdumajlis.net	8K	0.55
urduvoa.com	7.5K	0.52

Top 10 TLDs

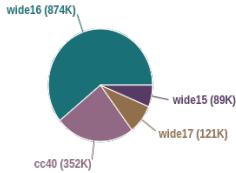
Domain	Docs	% of total
com	1.1M	79.78
pk	62K	4.34
com.pk	52K	3.61
org	45K	3.16
tv	35K	2.46
net	32K	2.25
info	11K	0.75
in	7.5K	0.52
ir	3.7K	0.26
site	2.7K	0.19

Documents size (in segments)



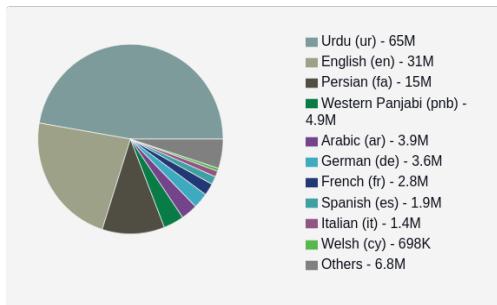
<= 25 segments 10.59% (152K documents)
> 25 segments 89.41% (1.3M documents)

Documents by collection

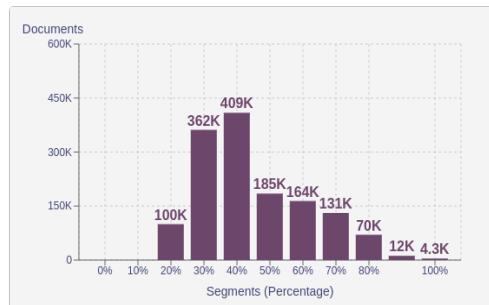


Language Distribution

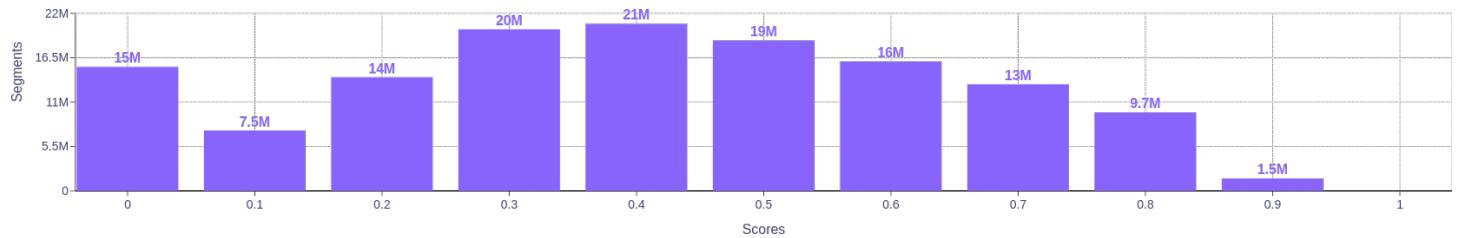
Number of segments



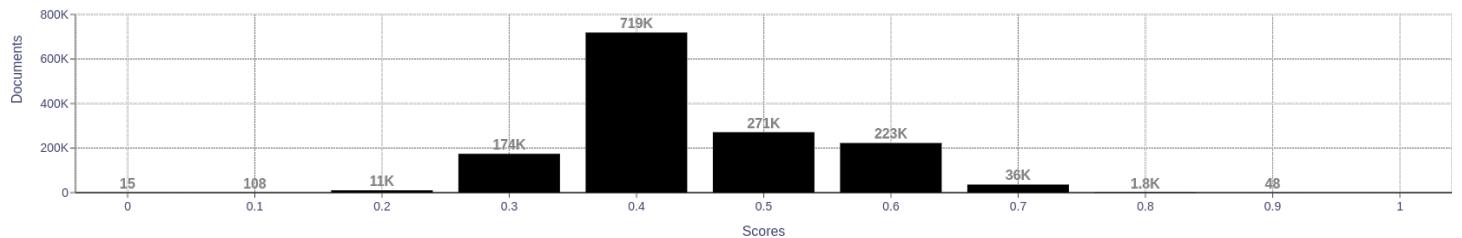
Percentage of segments in Urdu (ur) inside documents



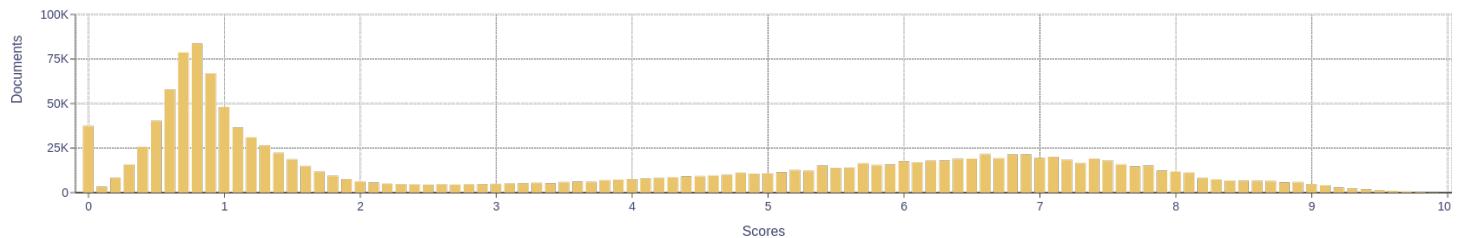
Distribution of segments by fluency score



Distribution of documents by average fluency score



Distribution of documents by document score



Segment length distribution by token

<= 49 tokens = 28M segments | 103M duplicates
> 50 tokens = 5.1M segments | 1.1M 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>