

General overview

Corpus	Analytics date	Language
kas_Arab.jsonl.tsv	11/4/2024	Kashmiri (ks)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
949	27,108	11,770 (43.42 %)	708K	5.87 MB	3,441,826

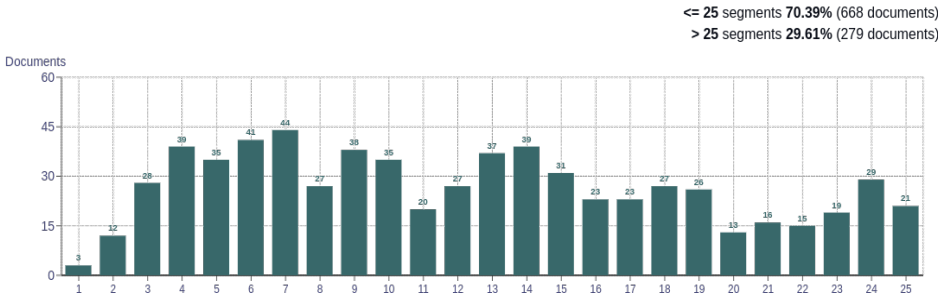
Top 10 domains

Domain	Docs	% of total
muneeburrahman.com	302	31.82
neabmagazine.com	254	26.77
wikipedia.org	157	16.54
neabinternational.org	52	5.48
wiktionary.org	45	4.74
newschecker.in	29	3.06
gotquestions.org	24	2.53
vikaspedia.in	19	2.00
aminkamil.blog	17	1.79
koshurakhbar.com	13	1.37

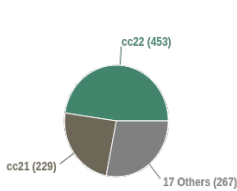
Top 10 TLDs

Domain	Docs	% of total
com	594	62.59
org	280	29.50
in	50	5.27
blog	17	1.79
ir	3	0.32
net	2	0.21
ae	1	0.11
ru	1	0.11
io	1	0.11

Documents size (in segments)

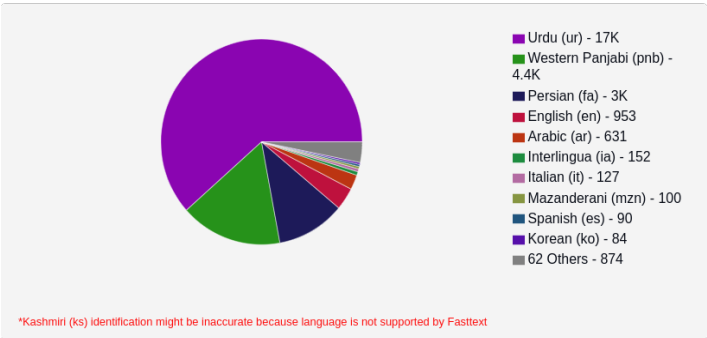


Documents by collection

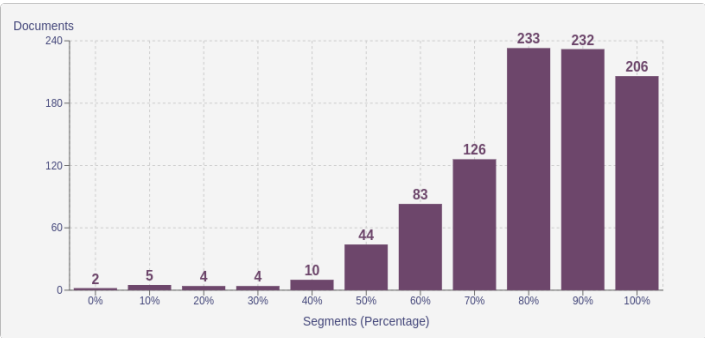


Language Distribution

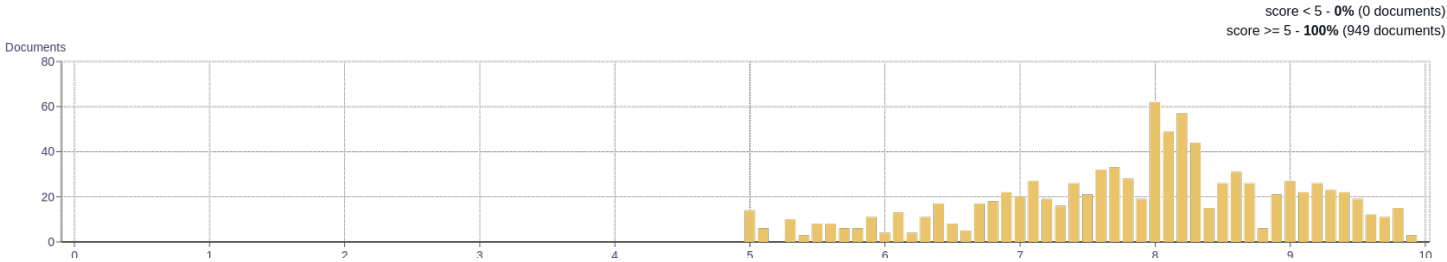
Number of segments



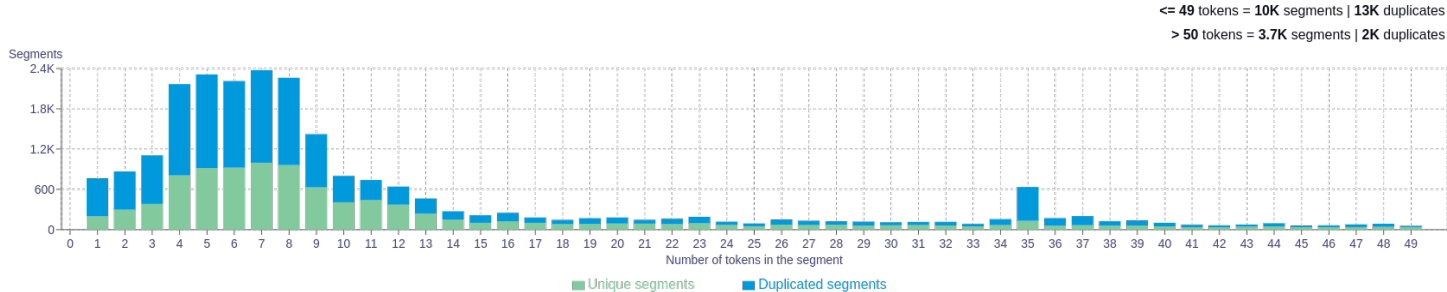
Percentage of segments in Kashmiri (ks) inside documents



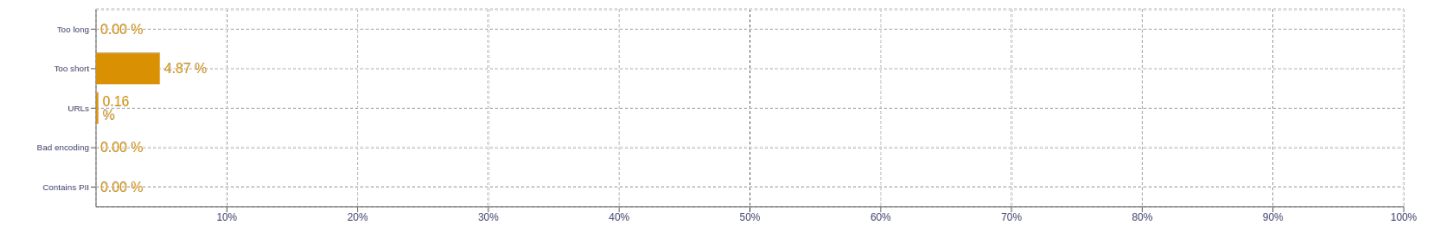
Distribution of documents by document score



Segment length distribution by token



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>