

General overview

Corpus	Analytics date	Language
ps_1.jsonl.tsv	3/17/2024	Pashto (ps)

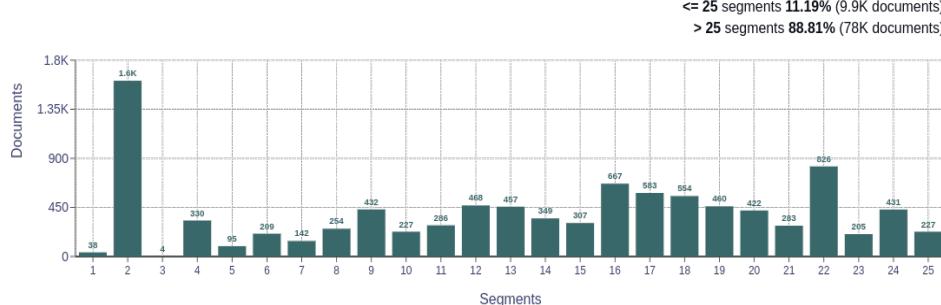
Volumes

Docs	Segments	Unique segments	Tokens	Size
88,212	10,984,513	24,186 (0.22 %)	131M	900.22 MB

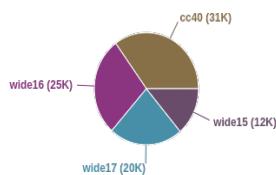
Type-Token Ratio

Pashto (ps)
0.01

Documents size (in segments)

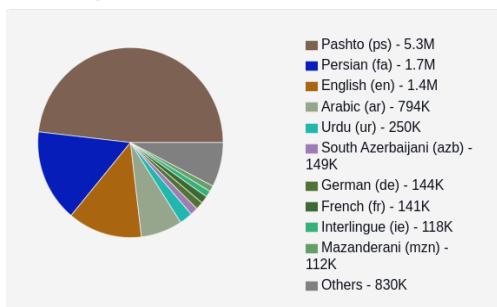


Documents by collection

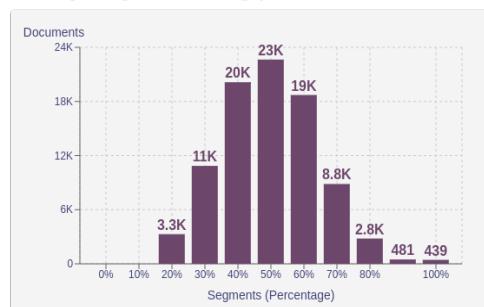


Language Distribution

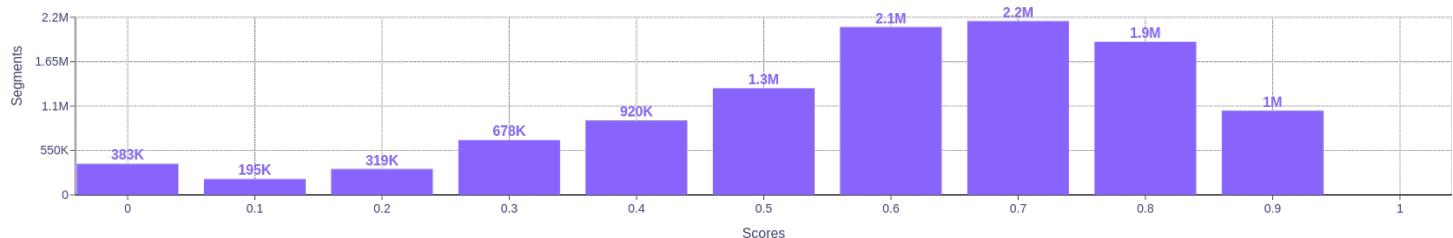
Number of segments



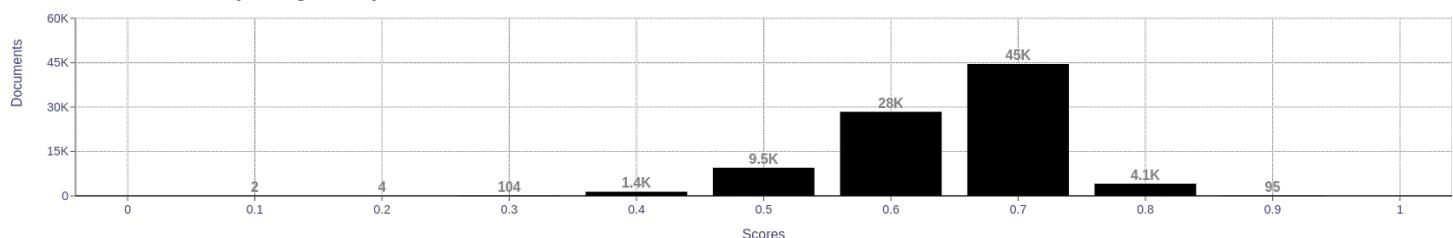
Percentage of segments in Pashto (ps) inside documents



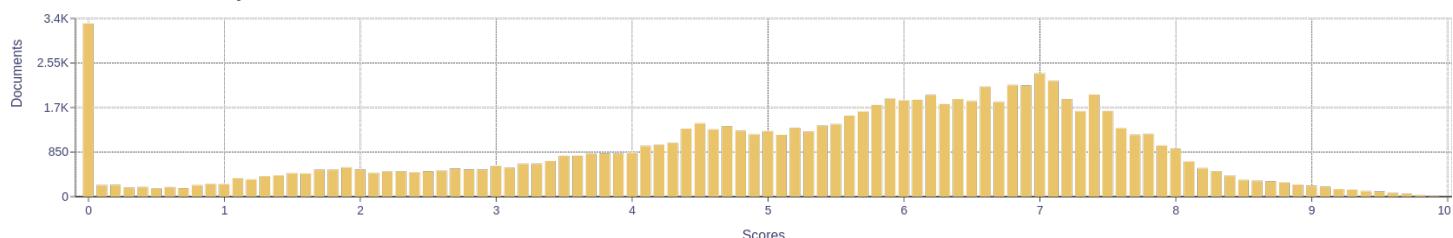
Distribution of segments by fluency score



Distribution of documents by average fluency score

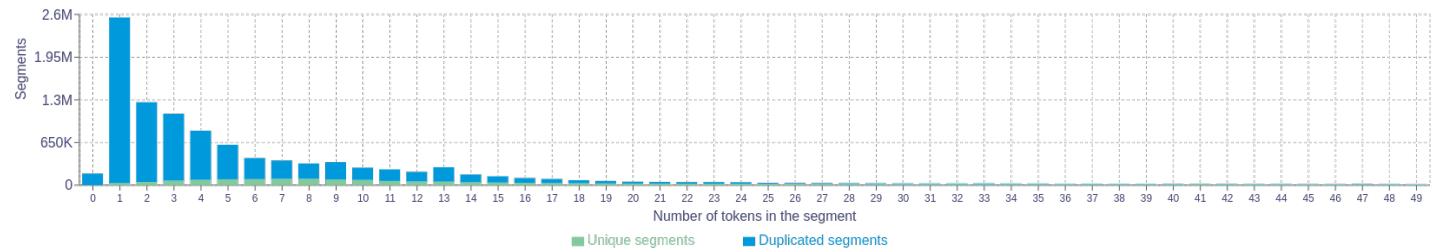


Distribution of documents by document score



Segment length distribution by token

<= 49 tokens = 1.7M segments | 8.8M duplicates
 > 50 tokens = 509K segments | 111K duplicates



Segment noise distribution



Frequent n-grams

Size	n-grams
1	1652289 کب 1388817 جب 501104 دی 464528 افغانستان 428709 دی
2	71463 افغانستان کب 51888 مهم خبرو نه 46617 اون، مهم days ago 43757 hours ago 42234
3	46614 اون، مهم خبرو نت all rights reserved 21826 20739 داس هم شته 13008 11728 ساینس او بکالوری مل الله عليه
4	11980 دنور لا رویان رغبہ خبرو نه 9078 جب یہ افغانستان کب 6580 6146 from twitter for iphone opens in new window 5917
5	5450 رسول الله مل الله عليه a password will be e 4498 4343 موبائل سرہ به تماس کن your email address will not 3816 email address will not be 3816

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.slinfo.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>