

General overview

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

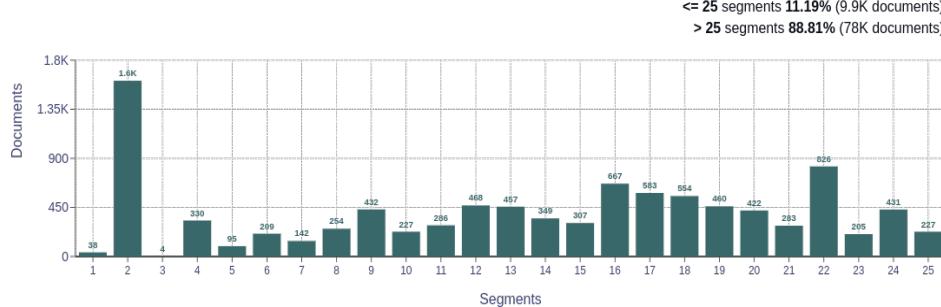
Volumes

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

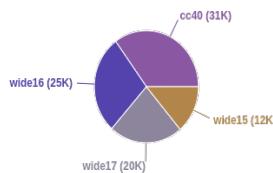
Type-Token Ratio

| Pashto (ps) |
|-------------|
| 0.01 |

Documents size (in segments)

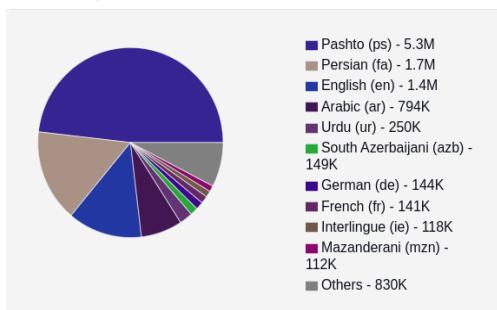


Documents by collection

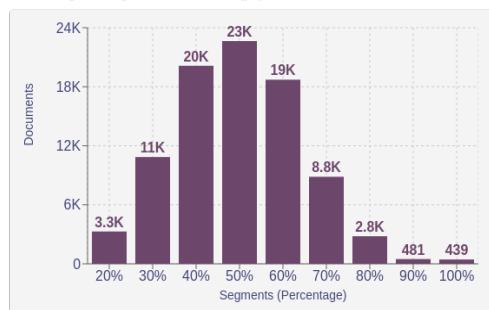


Language Distribution

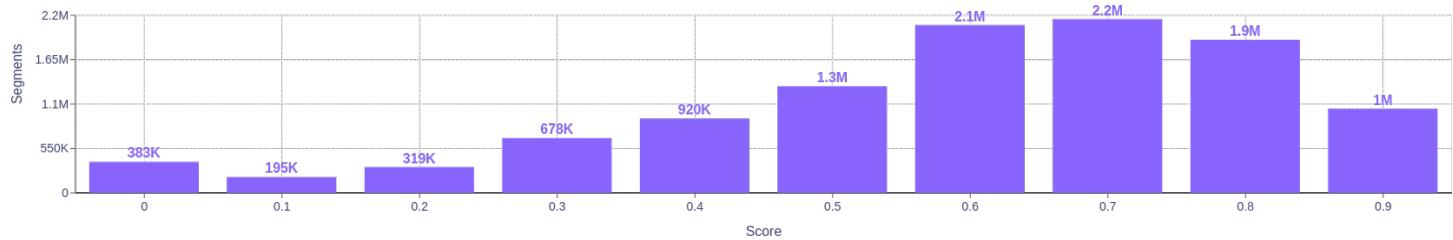
Number of segments



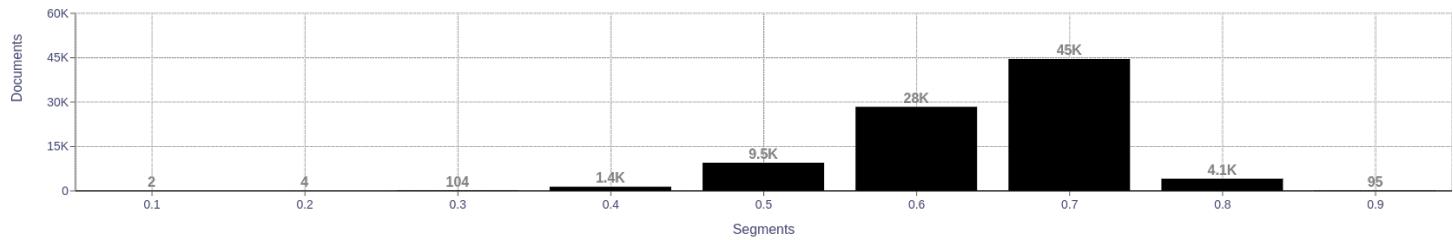
Percentage of segments in Pashto (ps) inside documents



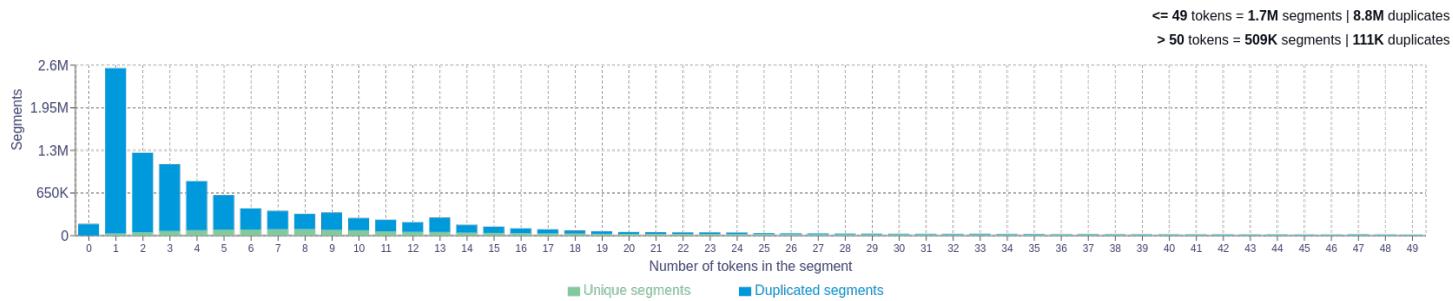
Distribution of segments by fluency score



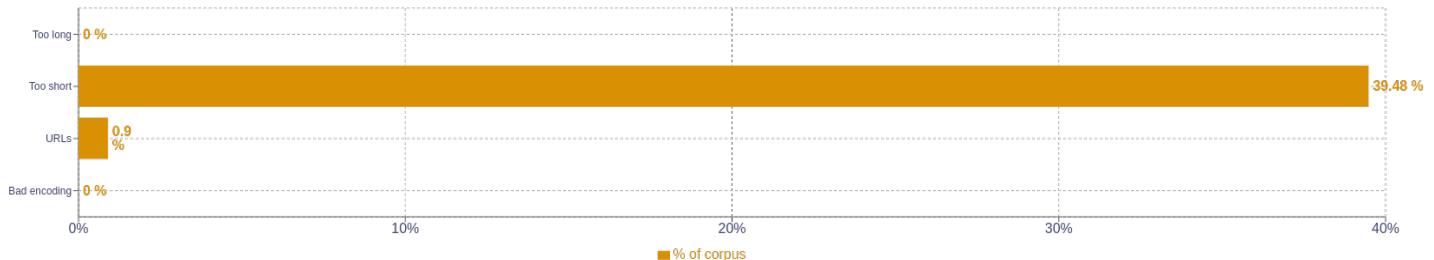
Distribution of documents by average fluency score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

| Size | n-grams |
|------|---|
| 1 | (1652289 کب) (1388817 جب) (501104 دی) (464528 افغانستان) (428709 دی) (428709 افغانستان) (428709 دی) (428709 دی) |
| 2 | (71463 افغانستان کی) (51888 مهم خبرونه) (46617 اون مهم) (43757 days ago) (42234 hours ago) |
| 3 | (46614 اون مهم خبرونه) (all rights reserved 21826) (20739 داسی هم شد) (13008 ملی اللہ علیہ) (11728 سپس او بکنالوری) (6580 ملی اللہ علیہ وسلم) |
| 4 | (11980 دنور لاڑویان رعنیہ خبرونہ) (9078 چن یہ افغانستان کی) (6146 from twitter for iphone) (5917 opens in new window) |
| 5 | (5450 رسول اللہ ملی اللہ علیہ) (4498 a password will be e) (4343 مونر سرہ یہ تماس کی) (3816 your email address will not) (3816 email address will not be) |

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>).

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.

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>