# HPLT Analytics report

**HPLT**Analytics

## General overview

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

### Volumes

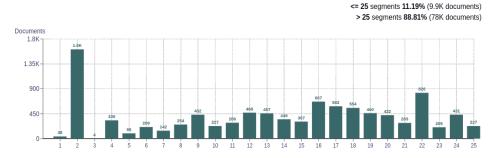| Docs | Segments | Unique segments | Tokens | Size | Characters |
|---|---|---|---|---|---|
| 88,212 | 10,984,513 | 2,078,407 (18.92 %) | 131M | 900.22 MB | |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| taand.com | 3.8K | 4.27 |
| spogmairadio.af | 3.5K | 4.00 |
| larawbar.net | 2.7K | 3.04 |
| qamosona.com | 2.5K | 2.87 |
| pashtovoa.com | 2.2K | 2.52 |
| khabarial.com | 2.2K | 2.46 |
| kandahar-tv.com | 2.1K | 2.42 |
| cri.cn | 1.9K | 2.16 |
| nunn.asia | 1.8K | 2.04 |
| wikipedia.org | 1.7K | 1.92 |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 55K | 62.31 |
| af | 8K | 9.03 |
| net | 6.1K | 6.97 |
| org | 5K | 5.70 |
| gov.af | 2.6K | 2.99 |
| cn | 1.9K | 2.16 |
| asia | 1.9K | 2.11 |
| info | 1.5K | 1.70 |
| org.af | 570 | 0.65 |
| tv | 515 | 0.58 |

## Documents size (in segments)

**<= 25** segments **11.19%** (9.9K documents)
**> 25** segments **88.81%** (78K documents)



## Documents by collection



cc40 (31K)
wide16 (25K)
wide15 (12K)
wide17 (20K)

## Language Distribution

### Number of segments



- Pashto (ps) - 5.3M
- Persian (fa) - 1.7M
- English (en) - 1.4M
- Arabic (ar) - 794K
- Urdu (ur) - 250K
- South Azerbaijani (azb) - 149K
- German (de) - 144K
- French (fr) - 141K
- Interlingue (ie) - 118K
- Mazanderani (mzn) - 112K
- 155 Others - 830K

### Percentage of segments in Pashto (ps) inside documents



## Distribution of documents by document score

score < 5 - **37.17%** (33K documents)
score >= 5 - **62.83%** (55K documents)



## Segment length distribution by token

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



- Unique segments
- Duplicated segments

## Segment noise distribution



Too long 0,20 %
Too short 39.48 %
URLs 0.90 %
Bad encoding 0.00 %

**Frequent n-grams**

| Size | n-grams |
|---|---|
| 1 | کی \| 1652289   چې \| 1388817   دی \| 501104   افغانستان \| 464528   دې \| 428709 |
| 2 | افغانستان کې \| 71463   مهم خبرونه \| 51888   اونۍ مهم \| 46617   days ago \| 43757   hours ago \| 42234 |
| 3 | اونۍ مهم خبرونه \| 46614   all rights reserved \| 21826   داس هم شته \| 20739   صلی الله علیه \| 13008   ساینس او ټکنالوژی \| 11728 |
| 4 | صلی الله علیه وسلم \| 11980   دنور لارویان رغیزه خبرونه \| 9078   چې په افغانستان کې \| 6580   from twitter for iphone \| 6146   opens in new window \| 5917 |
| 5 | رسول الله صلی الله علیه \| 5450   مونږ سره په تماس کې \| 4343   a password will be e \| 4498   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>, <ul>, <ol>, 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 or 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>, <ul>, <ol>, 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