# HPLT Analytics report

## General overview

| Corpus | Analytics date | Language |
|---|---|---|
| kas_Arab.jsonl.tsv | 12/3/2024 | 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)

<= 25 segments **70.39%** (668 documents)
> 25 segments **29.61%** (281 documents)



## Documents by collection

cc22 (453)
cc21 (229)
17 Others (267)



## Language Distribution

### Number of segments

- Urdu (ur) - 17K
- Western Panjabi (pnb) - 4.4K
- Persian (fa) - 3K
- English (en) - 953
- Arabic (ar) - 631
- Interlingua (ia) - 152
- Italian (it) - 127
- Mazanderani (mzn) - 100
- Spanish (es) - 90
- Korean (ko) - 84
- 62 Others - 874



*Kashmiri (ks) identification might be inaccurate because language is not supported by Fasttext

### Percentage of segments in Kashmiri (ks) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (949 documents)



## Segment length distribution by token

<= 49 tokens = 10K segments | 13K duplicates
> 50 tokens = 3.7K segments | 2K duplicates



- Unique segments
- Duplicated segments

## Segment noise distribution

| | |
|---|---|
| Too long | 1.29 % |
| Too short | 11.09 % |
| URLs | 0.21 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.00 % |

**Frequent n-grams**

| Size | n-grams |
|---|---|
| 1 | نَ \| 14742   چھ \| 14677   منز \| 12340   پ \| 7439   نَ \| 7036 |
| 2 | فلكس بنر \| 5864   بنر فلكس \| 5842   چھ نَ \| 1824   منز چھ \| 1306   چھِ نَ \| 699   منز چھ نَ \| 176   تم تم تم \| 176 |
| 3 | بنر فلكس بنر \| 5840   فلكس بنر فلكس \| 5838   آؤلی پَند پؤرَپی \| 235   منز چھ نَ \| 194 |
| 4 | بنر فلكس بنر فلكس \| 5836   فلكس بنر فلكس بنر \| 5828   تم تم تم تم \| 173   پَند پؤرَپی پرٯٹھ آمٹن \| 127   آؤلی پَند پؤرَپی پرٯٹھ \| 121 |
| 5 | بنر فلكس بنر فلكس بنر \| 5826   فلكس بنر فلكس بنر فلكس \| 5824   تم تم تم تم تم \| 171   آؤلی پَند پؤرَپی پرٯٹھ آمٹن \| 121   كرتس منز حد۔ روس خوش۔ \| 65 |

## 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