# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| hplt-v3-pbt_Arab | 9/24/2025 | Southern Pashto |

## Volumes

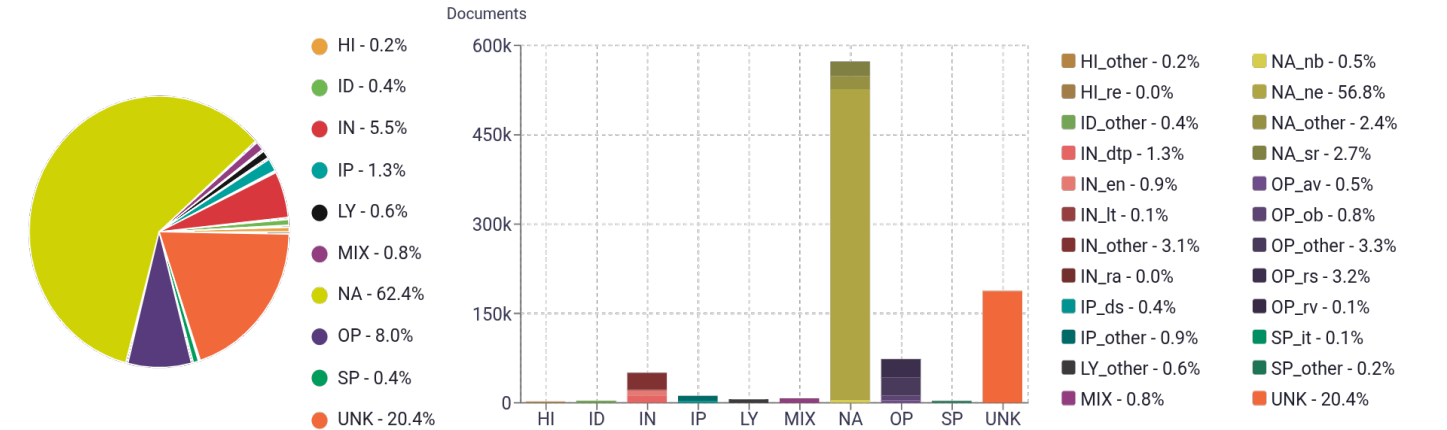| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 918,708 | 15,745,595 | 12,584,332 (79.92 %) | 577M | 2,431,797,440 | 3.99 GB |

## Top 10 domains

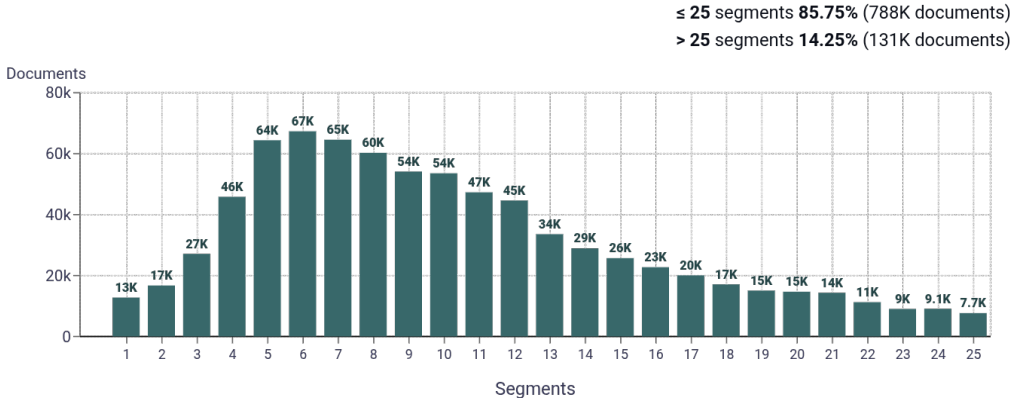| Domain | Docs | % of total |
|---|---|---|
| nunn.asia | 61K | 6.66% |
| pashtovoa.com | 44K | 4.74% |
| mashaalradio.com | 40K | 4.34% |
| azadiradio.com | 29K | 3.19% |
| voadeewanews.com | 26K | 2.84% |
| taand.com | 23K | 2.51% |
| bbc.com | 21K | 2.31% |
| bakhtarnews.af | 19K | 2.07% |
| tolafghan.com | 18K | 2.00% |
| dw.com | 18K | 1.95% |

## Top 10 TLDs

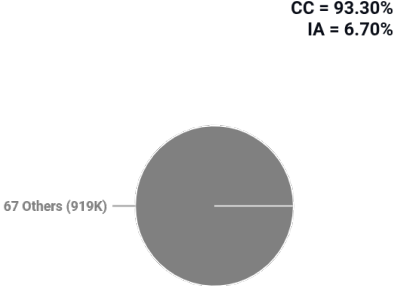| Domain | Docs | % of total |
|---|---|---|
| com | 583K | 63.47% |
| af | 66K | 7.14% |
| asia | 62K | 6.72% |
| net | 45K | 4.91% |
| org | 28K | 3.04% |
| gov.af | 27K | 2.89% |
| cn | 15K | 1.61% |
| tv | 11K | 1.24% |
| com.af | 7.8K | 0.85% |
| co | 6.3K | 0.69% |

## Register labels



Pie chart legend:
- HI - 0.2%
- ID - 0.4%
- IN - 5.5%
- IP - 1.3%
- LY - 0.6%
- MIX - 0.8%
- NA - 62.4%
- OP - 8.0%
- SP - 0.4%
- UNK - 20.4%

Bar chart legend:
- HI_other - 0.2%
- HI_re - 0.0%
- ID_other - 0.4%
- IN_dtp - 1.3%
- IN_en - 0.9%
- IN_lt - 0.1%
- IN_other - 3.1%
- IN_ra - 0.0%
- IP_ds - 0.4%
- IP_other - 0.9%
- LY_other - 0.6%
- MIX - 0.8%
- NA_nb - 0.5%
- NA_ne - 56.8%
- NA_other - 2.4%
- NA_sr - 2.7%
- OP_av - 0.5%
- OP_ob - 0.8%
- OP_other - 3.3%
- OP_rs - 3.2%
- OP_rv - 0.1%
- SP_it - 0.1%
- SP_other - 0.2%
- UNK - 20.4%

🤖 **MT**:15.8% | 145K Documents

## Documents size (in segments) ⓘ

≤ **25** segments **85.75%** (788K documents)
> **25** segments **14.25%** (131K documents)



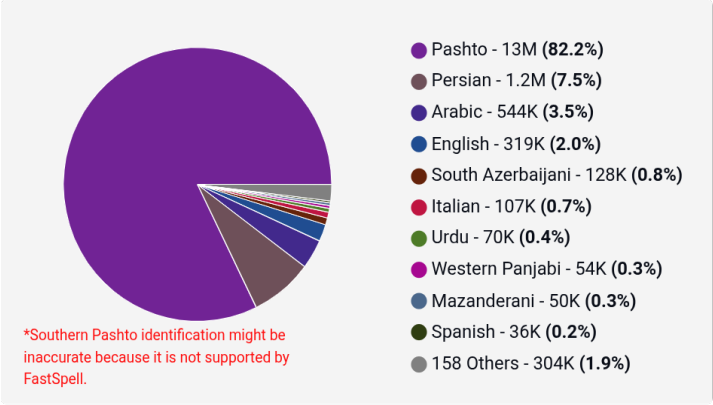## Document collections

**CC = 93.30%**
**IA = 6.70%**



67 Others (919K)

# Language Distribution

## Number of segments in the Southern Pashto corpus



- Pashto - 13M **(82.2%)**
- Persian - 1.2M **(7.5%)**
- Arabic - 544K **(3.5%)**
- English - 319K **(2.0%)**
- South Azerbaijani - 128K **(0.8%)**
- Italian - 107K **(0.7%)**
- Urdu - 70K **(0.4%)**
- Western Panjabi - 54K **(0.3%)**
- Mazanderani - 50K **(0.3%)**
- Spanish - 36K **(0.2%)**
- 158 Others - 304K **(1.9%)**

*Southern Pashto identification might be inaccurate because it is not supported by FastSpell.

## Percentage of segments in Southern Pashto inside documents

segments < 50% - **1.19%** (11K documents)
segments ≥ 50% - **98.81%** (908K documents)



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score ≥ 5 - **100%** (919K documents)



## Segment length distribution by token

≤ 49 tokens = **12M** segments | **2.8M** duplicates
> 50 tokens = **3.5M** segments | **325K** duplicates



## Segment noise distribution



- Too long: **0.54%**
- Too short: **6.86%**
- URLs: **0.54%**
- Bad encoding: **0.00%**
- Contains PII: **0.09%**

■ % of corpus

## Frequent n-grams

| SIZE | N-GRAMS |
|------|---------|
| 1 | چي \| 11,437,703   کې \| 10,860,973   دی \| 3,154,710   د \| 3,086,039   یې \| 2,723,042 |
| 2 | افغانستان کې \| 365,914   دی چي \| 356,822   چي تاسو \| 310,129   حال کې \| 245,306   کې چي \| 240,102 |
| 3 | تاسو کولی شئ \| 121,664   حال کې چي \| 102,433   چي په دی \| 100,235   چي د افغانستان \| 90,824   چي د دی \| 75,747 |
| 4 | چي په افغانستان کې \| 54,826   پداسي حال کې چي \| 49,411   صلی الله علیه وسلم \| 46,164   حال کې ده چي \| 27,568   صلی الله علیه وسلم \| 19,260 |
| 5 | رسول الله صلی الله علیه \| 26,219   رسول الله صلی الله علیه \| 11,407   آژانس د خبر له مخي \| 10,201   اره خه نه دي ویلي \| 8,548   جمهور خبري آژانس د خبر \| 7,540 |

# 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

**Register labels**

| Name | Abbr. | Name | Abbr. | Name | Abbr. |
|---|---|---|---|---|---|
| **Machine-translated** | MT | **How-to or instructions** | **HI** | Description of a thing or person | dtp |
| **Lyrical** | LY | Recipe | re | FAQ | fi |
| **Spoken** | SP | **Informational persuasion** | **IP** | Legal terms & conditions | lt |
| Interview | it | Description with intent to sell | ds | **Opinion** | **OP** |
| **Interactive discussion** | ID | News & opinion blog or editorial | ed | Review | rv |
| **Narrative** | NA | **Informational description** | **IN** | Opinion blog | ob |
| News report | ne | Enciclopedia article | en | Denominational religious blog or sermon | rs |
| Sports report | sr | Research article | ra | Advice | av |
| Narrative blog | nb | | | | |