# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| hplt-v3-snd_Arab | 9/18/2025 | Sindhi |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 363,829 | 6,274,353 | 5,162,379 (82.28 %) | 261M | 1,087,312,087 | 1.78 GB |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| awamiawaz.pk | 42K | 11.47% |
| ktnnews.tv | 17K | 4.60% |
| awamiawaz.com | 16K | 4.42% |
| thetimenews.tv | 16K | 4.32% |
| dailysindhyar.com | 13K | 3.60% |
| pahenjiakhbar.com | 13K | 3.59% |
| sindhexpress.co... | 12K | 3.26% |
| dailysarwan.com | 12K | 3.17% |
| voiceofsindh.co... | 8.4K | 2.30% |
| androidsis.com | 7.3K | 2.00% |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 209K | 57.37% |
| pk | 47K | 12.80% |
| tv | 38K | 10.35% |
| com.pk | 29K | 7.94% |
| org | 14K | 3.75% |
| net | 9.5K | 2.60% |
| zone | 4.6K | 1.27% |
| fr | 1.3K | 0.36% |
| ru | 1.1K | 0.31% |
| co.uk | 776 | 0.21% |

## Register labels



- HI - 0.1%
- ID - 0.0%
- IN - 3.7%
- IP - 1.2%
- LY - 0.2%
- MIX - 0.5%
- NA - 52.0%
- OP - 2.7%
- SP - 0.1%
- UNK - 39.5%

- HI_other - 0.1%
- HI_re - 0.0%
- ID_other - 0.0%
- IN_dtp - 0.7%
- IN_en - 1.5%
- IN_lt - 0.1%
- IN_other - 1.5%
- IP_ds - 0.6%
- IP_other - 0.6%
- LY_other - 0.2%
- MIX - 0.5%
- NA_nb - 0.4%
- NA_ne - 49.0%
- NA_other - 1.1%
- NA_sr - 1.5%
- OP_av - 0.0%
- OP_ob - 0.6%
- OP_other - 1.0%
- OP_rs - 0.9%
- OP_rv - 0.2%
- SP_it - 0.0%
- SP_other - 0.0%
- UNK - 39.5%

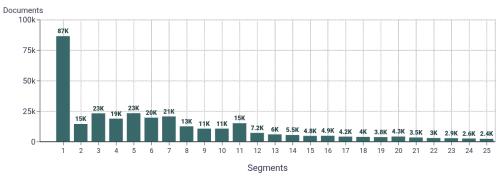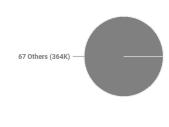🤖 **MT**:37.4% | 136K Documents

## Documents size (in segments) ⓘ

≤ 25 segments **86.73%** (316K documents)
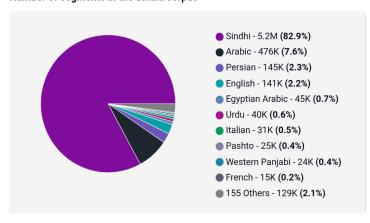\> 25 segments **13.27%** (48K documents)
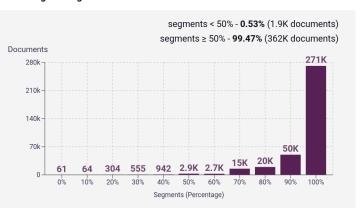


## Document collections

**CC = 97.65%**
**IA = 2.35%**



67 Others (364K)

## Language Distribution

### Number of segments in the Sindhi corpus

- Sindhi - 5.2M **(82.9%)**
- Arabic - 476K **(7.6%)**
- Persian - 145K **(2.3%)**
- English - 141K **(2.2%)**
- Egyptian Arabic - 45K **(0.7%)**
- Urdu - 40K **(0.6%)**
- Italian - 31K **(0.5%)**
- Pashto - 25K **(0.4%)**
- Western Panjabi - 24K **(0.4%)**
- French - 15K **(0.2%)**
- 155 Others - 129K **(2.1%)**

### Percentage of segments in Sindhi inside documents

segments < 50% - **0.53%** (1.9K documents)
segments ≥ 50% - **99.47%** (362K documents)

Documents

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 61 | 64 | 304 | 555 | 942 | 2.9K | 2.7K | 15K | 20K | 50K | 271K |
| 0% | 10% | 20% | 30% | 40% | 50% | 60% | 70% | 80% | 90% | 100% |

Segments (Percentage)

### Distribution of documents by document score

score < 5 - **0%** (0 documents)
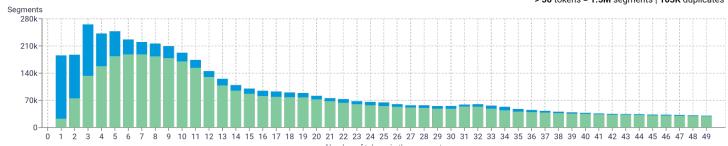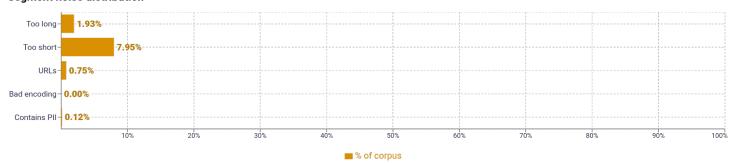score ≥ 5 - **100%** (364K documents)

Documents

### Segment length distribution by token

≤ 49 tokens = **4.7M** segments | **1M** duplicates
> 50 tokens = **1.5M** segments | **103K** duplicates

Segments

### Segment noise distribution

| | |
|---|---|
| Too long | **1.93%** |
| Too short | **7.95%** |
| URLs | **0.75%** |
| Bad encoding | **0.00%** |
| Contains PII | **0.12%** |

■ % of corpus

## Frequent n-grams

| SIZE | N-GRAMS |
|---|---|
| 1 | نه \| 1,354,554    ڪيو \| 1,082,601    جا \| 859,316    طور \| 629,232    ڪئي \| 510,207 |
| 2 | ڪئي وئي \| 110,766    اسلام آباد \| 90,187    ڪيو ويندو \| 85,569    ڪيو وڃي \| 81,426    عمران خان \| 78,105 |
| 3 | پي ٽي آ ء \| 36,139    استعمال ڪيو ويندو \| 20,193    ايس ايس پي \| 16,392    بلاول ڀٽو زرداري \| 15,426    مراد علي شاهه \| 14,641 |
| 4 | سنڌ جي وڏي وزير \| 9,705    سيد مراد علي شاهه \| 7,219    چيئرمين بلاول ڀٽو زرداري \| 6,736    وزير سيد مراد علي \| 6,712    وڏي وزير سيد مراد \| 6,199 |
| 5 | وڏي وزير سيد مراد علي \| 6,105    وزير سيد مراد علي شاهه \| 5,767    ڊبليو ڊبليو ڊبليو ڊبليو \| 4,856    سنڌ جي وڏي وزير سيد \| 4,691    پ چيئرمين بلاول ڀٽو \| 3,830 |

# 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 | | | | |