# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Analytics date | Language |
|---|---|---|
| HPLT-v2-tur_Latn.tsv | 9/27/2024 | Turkish (tr) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Size | Characters |
|---|---|---|---|---|---|
| 116,566,047 | 2,574,897,928 | | | 394.5 GB | 387,178,346,456 |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| blogspot.com | 1.7M | 1.47 |
| blogspot.com.tr | 1.5M | 1.26 |
| hurriyet.com.tr | 1.4M | 1.17 |
| wikipedia.org | 580K | 0.50 |
| docplayer.biz.tr | 529K | 0.45 |
| sabah.com.tr | 485K | 0.42 |
| sikayetvar.com | 454K | 0.39 |
| haberler.com | 422K | 0.36 |
| eksisozluk.com | 396K | 0.34 |
| haberaktuel.com | 393K | 0.34 |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 76M | 64.79 |
| com.tr | 13M | 11.02 |
| net | 10M | 8.89 |
| org | 6.1M | 5.19 |
| org.tr | 1.1M | 0.91 |
| biz.tr | 774K | 0.66 |
| info | 723K | 0.62 |
| gen.tr | 691K | 0.59 |
| gov.tr | 683K | 0.59 |
| edu.tr | 678K | 0.58 |

## Documents size (in segments)

**<= 25** segments **78.06%** (91M documents)
**> 25** segments **21.94%** (26M documents)



## Documents by collection



cc18 (20M)
cc22 (21M)
19 Others (76M)

## Language Distribution

### Number of segments



- Turkish (tr) - 2.2B
- English (en) - 111M
- Italian (it) - 77M
- French (fr) - 24M
- German (de) - 20M
- Dutch (nl) - 11M
- Azerbaijani (az) - 10M
- Spanish (es) - 8.7M
- Swedish (sv) - 7.5M
- Esperanto (eo) - 7.3M
- 165 Others - 78M

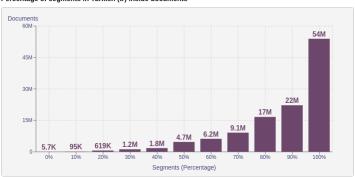### Percentage of segments in Turkish (tr) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (117M documents)



## Segment noise distribution



| | |
|---|---|
| Too long | 1.12 % |
| Too short | 13.98 % |
| URLs | 1.67 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.19 % |

## 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