# HPLT Analytics report

## General overview

| Corpus | Analytics date | Language |
|---|---|---|
| HPLT-v2-hrv_Latn.tsv | 9/19/2024 | Croatian (hr) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Size | Characters |
|---|---|---|---|---|---|
| 12,303,820 | 297,130,317 | | | 45.7 GB | 47,710,541,358 |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| wikipedia.org | 264K | 2.15 |
| dnevnik.hr | 170K | 1.38 |
| blogspot.com | 138K | 1.12 |
| tportal.hr | 95K | 0.77 |
| skole.hr | 88K | 0.72 |
| jutarnji.hr | 84K | 0.68 |
| index.hr | 82K | 0.67 |
| metro-portal.hr | 73K | 0.59 |
| 24sata.hr | 66K | 0.54 |
| hrt.hr | 59K | 0.48 |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| hr | 5.5M | 45.05 |
| com | 3.4M | 27.95 |
| org | 592K | 4.81 |
| net | 551K | 4.48 |
| ba | 421K | 3.43 |
| rs | 340K | 2.76 |
| info | 288K | 2.34 |
| com.hr | 247K | 2.01 |
| eu | 219K | 1.78 |
| news | 140K | 1.14 |

## Documents size (in segments)

<= **25** segments **78.21%** (9.6M documents)
> **25** segments **21.79%** (2.7M documents)



## Documents by collection



cc22 (3.9M)
cc18 (1.9M)
cc21 (1.6M)
18 Others (5M)

## Language Distribution

### Number of segments



- Croatian (hr) - 232M
- English (en) - 14M
- Serbian (sr) - 12M
- Italian (it) - 7.5M
- Bosnian (bs) - 5M
- Polish (pl) - 3.1M
- French (fr) - 2.9M
- German (de) - 2.6M
- Czech (cs) - 2.2M
- Portuguese (pt) - 1.3M
- 165 Others - 14M

### Percentage of segments in Croatian (hr) inside documents



## Distribution of documents by document score

score <= 5 - **99.99%** (12M documents)
score > 5 - **0.01%** (1.8K documents)



## Segment noise distribution



- Too long — 0.00 %
- Too short — 13.98 %
- URLs — 1.58 %
- Bad encoding — 0.00 %
- Contains PII — 0.41 %

## 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