# HPLT Analytics report

## General overview

| Corpus | Analytics date | Language |
|---|---|---|
| HPLT-v2-est_Latn.tsv | 9/20/2024 | Estonian (et) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Size | Characters |
|---|---|---|---|---|---|
| 8,449,320 | 264,400,562 | | | 34.35 GB | 35,759,962,808 |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| wikipedia.org | 293K | 3.47 |
| blogspot.com | 290K | 3.43 |
| postimees.ee | 288K | 3.40 |
| err.ee | 249K | 2.95 |
| delfi.ee | 225K | 2.67 |
| aripaev.ee | 158K | 1.87 |
| pilguheit.com | 149K | 1.76 |
| ohtuleht.ee | 117K | 1.39 |
| kliinik.ee | 93K | 1.10 |
| wordpress.com | 85K | 1.00 |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| ee | 5.5M | 65.36 |
| com | 1.7M | 19.75 |
| org | 436K | 5.16 |
| eu | 222K | 2.63 |
| net | 116K | 1.37 |
| fi | 76K | 0.90 |
| com.ee | 53K | 0.62 |
| edu.ee | 49K | 0.58 |
| info | 42K | 0.50 |
| pt | 21K | 0.25 |

## Documents size (in segments)

<= 25 segments **75.41%** (6.4M documents)
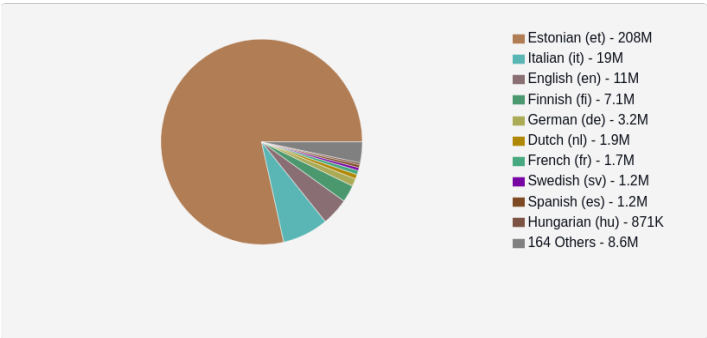> 25 segments **24.59%** (2.1M documents)



## Documents by collection



## Language Distribution

### Number of segments



- Estonian (et) - 208M
- Italian (it) - 19M
- English (en) - 11M
- Finnish (fi) - 7.1M
- German (de) - 3.2M
- Dutch (nl) - 1.9M
- French (fr) - 1.7M
- Swedish (sv) - 1.2M
- Spanish (es) - 1.2M
- Hungarian (hu) - 871K
- 164 Others - 8.6M
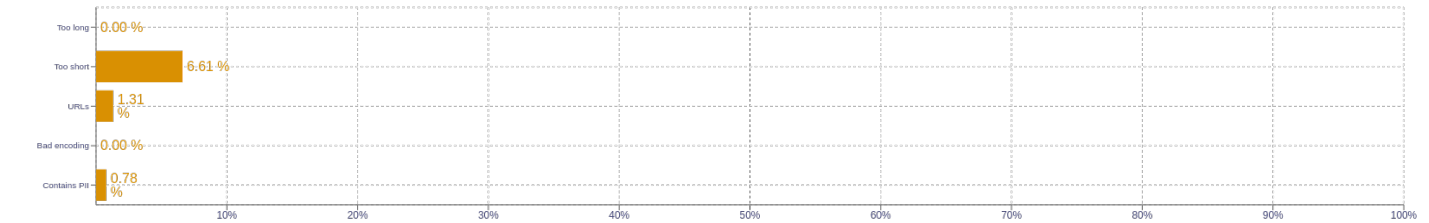
### Percentage of segments in Estonian (et) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (8.4M documents)



## Segment noise distribution

## 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