# HPLT Analytics report

**HPLT** Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| est_Latn.jsonl.tsv | 6/6/2025 | Estonian (et) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 8,449,258 | 264,352,085 | 102,779,262 (38.88 %) | 5.6B | 35,759,962,808 | 34.35 GB |

### 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% |

## Register labels



- HI - 3.4%
- ID - 3.7%
- IN - 14.1%
- IP - 23.0%
- LY - 0.1%
- MIX - 3.9%
- NA - 33.9%
- OP - 5.4%
- SP - 0.9%
- UNK - 11.6%

🤖 **MT**:8.5% | 722K Documents

Documents

- HI_other - 2.0%
- HI_re - 1.4%
- ID_other - 3.7%
- IN_dtp - 3.9%
- IN_en - 3.7%
- IN_fi - 0.0%
- IN_lt - 1.2%
- IN_other - 5.3%
- IN_ra - 0.1%
- IP_ds - 20.1%
- IP_ed - 0.0%
- IP_other - 2.9%
- LY_other - 0.1%
- MIX - 3.9%
- NA_nb - 8.5%
- NA_ne - 18.7%
- NA_other - 3.7%
- NA_sr - 3.0%
- OP_av - 0.6%
- OP_ob - 1.7%
- OP_other - 1.3%
- OP_rs - 0.6%
- OP_rv - 1.2%
- SP_it - 0.6%
- SP_other - 0.3%
- UNK - 11.6%

## Documents size (in segments)

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



## Documents by collection

**CC = 73.12%**
**IA = 26.88%**



cc22 (2.5M)
cc18 (1.5M)
cc21 (1.3M)
18 Others (3.2M)

## Language Distribution

### Number of segments in the Estonian (et) corpus



- 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

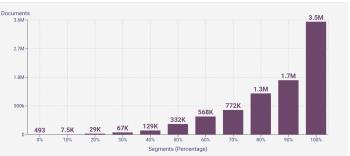### Percentage of segments in Estonian (et) inside documents



## Distribution of documents by document score

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

## Segment length distribution by token

≤ **49** tokens = **86M** segments | **149M** duplicates
> **50** tokens = **29M** segments | **12M** duplicates



Segments

| Segments | |
| --- | --- |
| 28M | |
| 21M | |
| 14M | |
| 7M | |
| 0 | |

Number of tokens in the segment

■ Unique segments  ■ Duplicated segments

## Segment noise distribution



| | |
| --- | --- |
| Too long | 0.48 % |
| Too short | 18.34 % |
| URLs | 1.55 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.78 % |

## Frequent n-grams

| Size | n-grams |
| --- | --- |
| 1 | eesti \| 14056503   võib \| 10394869   kes \| 9821756   kuid \| 9028429   ole \| 8052598 |
| 2 | võib olla \| 1415584   samal ajal \| 918895   mitte ainult \| 798046   ettevõttega ühendust \| 739660   euroopa liidu \| 664639 |
| 3 | viimase tunni jooksul \| 332931   vaadanud seda hotelli \| 332641   külastajat on vaadanud \| 332641   hotelli viimase tunni \| 332641   olulisemate uudiste kokkuvõte \| 276237 |
| 4 | vaadanud seda hotelli viimase \| 332641   hotelli viimase tunni jooksul \| 332641   vastas dr jüri ennet \| 93500   president toomas hendrik ilves \| 87180   isamaa ja res publica \| 73437 |
| 5 | vaadanud seda hotelli viimase tunni \| 332641   külastajat on vaadanud seda hotelli \| 332641   vana ning kuulub väljaande digitaalsesse \| 48927   aastat vana ning kuulub väljaande \| 48927   võib olla vajalik kaasaegsete allikatega \| 48926 |

## 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 | | | | |