# HPLT Analytics report

HPLT Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| lit_Latn.jsonl.tsv | 6/10/2025 | Lithuanian (lt) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 13,338,275 | 322,101,556 | 137,631,836 (42.73 %) | 8.1B | 50,084,920,486 | 49.59 GB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| wikipedia.org | 401K | 3.00% |
| delfi.lt | 253K | 1.90% |
| 15min.lt | 241K | 1.81% |
| lzinios.lt | 231K | 1.73% |
| diena.lt | 207K | 1.55% |
| mokslobaze.lt | 142K | 1.06% |
| lrt.lt | 128K | 0.96% |
| alkas.lt | 125K | 0.93% |
| blogspot.com | 116K | 0.87% |
| hotels.com | 91K | 0.68% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| lt | 11M | 79.01% |
| com | 1.4M | 10.44% |
| org | 534K | 4.01% |
| eu | 278K | 2.09% |
| net | 165K | 1.24% |
| info | 94K | 0.70% |
| ru | 22K | 0.17% |
| co.uk | 21K | 0.16% |
| today | 16K | 0.12% |
| pl | 14K | 0.11% |

## Register labels



- HI - 3.2%
- ID - 1.9%
- IN - 15.8%
- IP - 26.6%
- LY - 0.2%
- MIX - 4.0%
- NA - 29.1%
- OP - 5.9%
- SP - 0.9%
- UNK - 12.5%

🤖 **MT**:9.4% | 1.3M Documents



- HI_other - 1.5%
- HI_re - 1.7%
- ID_other - 1.9%
- IN_dtp - 4.6%
- IN_en - 3.4%
- IN_fi - 0.0%
- IN_lt - 1.3%
- IN_other - 6.4%
- IN_ra - 0.1%
- IP_ds - 23.9%
- IP_ed - 0.0%
- IP_other - 2.7%
- LY_other - 0.2%
- MIX - 4.0%
- NA_nb - 4.4%
- NA_ne - 19.0%
- NA_other - 3.4%
- NA_sr - 2.3%
- OP_av - 1.4%
- OP_ob - 1.4%
- OP_other - 1.3%
- OP_rs - 0.9%
- OP_rv - 0.9%
- SP_it - 0.6%
- SP_other - 0.2%
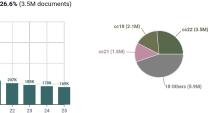- UNK - 12.5%

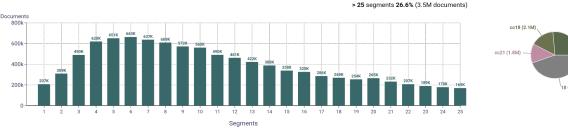## Documents size (in segments)

<= 25 segments **73.4%** (9.8M documents)
> 25 segments **26.6%** (3.5M documents)



## Documents by collection

CC = 66.62%
IA = 33.38%



cc18 (2.1M), cc22 (3.5M), cc21 (1.8M), 18 Others (5.9M)
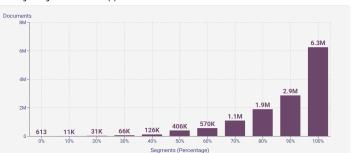
## Language Distribution

### Number of segments in the Lithuanian (lt) corpus



- Lithuanian (lt) - 275M
- English (en) - 11M
- Italian (it) - 9.5M
- Spanish (es) - 2.4M
- Esperanto (eo) - 2.2M
- German (de) - 2.2M
- French (fr) - 2.2M
- Polish (pl) - 1.9M
- Portuguese (pt) - 1.6M
- Finnish (fi) - 1.4M
- 165 Others - 12M

### Percentage of segments in Lithuanian (lt) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (13M documents)

## Segment length distribution by token

Segments
20M
15M
10M
5M
0

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49

Number of tokens in the segment

■ Unique segments ■ Duplicated segments

## Segment noise distribution



| | |
|---|---|
| Too long | 0.72 % |
| Too short | 14.77 % |
| URLs | 1.90 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.46 % |

10% 20% 30% 40% 50% 60% 70% 80% 90% 100%

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | yra \| 52706592    buvo \| 21554146    gali \| 19860198    m. \| 14828048    lietuvos \| 13546087 |
| 2 | gali būti \| 5800176    turi būti \| 2411956    šiuo metu \| 2050120    šiek tiek \| 1681239    lietuvos respublikos \| 1633727 |
| 3 | kartus per dieną \| 329067    tuo pačiu metu \| 311230    akcijų pasirinkimo sandoriai \| 280489    širdies ir kraujagyslių \| 267274    automobilių stovėjimo aikštelė \| 262772 |
| 4 | šį viešbutį per paskutiniąją \| 251796    viešbutį per paskutiniąją valandą \| 251796    vilniaus visuomenės sveikatos biuras \| 160519    socialinės apsaugos ir darbo \| 125022    visuomenės sveikatos biuras kviečia \| 106680 |
| 5 | šį viešbutį per paskutiniąją valandą \| 251796    peržiūrėjo šį viešbutį per paskutiniąją \| 251796    visuomenės informavimo priemonėse bei interneto \| 76992    informavimo priemonėse bei interneto tinklalapiuose \| 76906    visuomenės sveikatos biuras vykdo programą \| 75924 |

## 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 | | | | |