# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| ltg_Latn.jsonl.tsv | 12/6/2024 | Latgalian (ltg) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 9,209 | 151,382 | 77,506 (51.20 %) | 4.8M | 26,735,255 | 27.41 MB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| lakuga.lv | 3.3K | 36.23% |
| wikipedia.org | 1.8K | 19.43% |
| lgsc.lv | 1.3K | 14.27% |
| naktineica.lv | 352 | 3.82% |
| bonuks.lv | 295 | 3.20% |
| cyxob.lv | 268 | 2.91% |
| lsm.lv | 191 | 2.07% |
| rezeknesbibliot... | 96 | 1.04% |
| sciencegraph.net | 89 | 0.97% |
| jw.org | 80 | 0.87% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| lv | 6.7K | 72.84% |
| org | 1.9K | 20.76% |
| com | 148 | 1.61% |
| eu | 143 | 1.55% |
| net | 108 | 1.17% |
| cz | 86 | 0.93% |
| ru | 35 | 0.38% |
| gov.lv | 16 | 0.17% |
| in | 12 | 0.13% |
| info | 7 | 0.08% |

## Register labels



- HI - 0.1%
- ID - 1.3%
- IN - 26.0%
- IP - 9.3%
- LY - 0.8%
- MIX - 1.8%
- NA - 34.0%
- OP - 3.4%
- SP - 1.9%
- UNK - 21.3%

**MT**:1.5% | 136 Documents

- HI_other - 0.1%
- HI_re - 0.0%
- ID_other - 1.3%
- IN_dtp - 3.2%
- IN_en - 16.3%
- IN_fi - 0.0%
- IN_lt - 0.0%
- IN_other - 6.5%
- IN_ra - 0.0%
- IP_ds - 2.7%
- IP_ed - 0.0%
- IP_other - 6.6%
- LY_other - 0.8%
- MIX - 1.8%
- NA_nb - 4.8%
- NA_ne - 13.9%
- NA_other - 14.9%
- NA_sr - 0.4%
- OP_av - 0.0%
- OP_ob - 0.0%
- OP_other - 0.1%
- OP_rs - 3.1%
- OP_rv - 0.2%
- SP_it - 1.2%
- SP_other - 0.7%
- UNK - 21.3%

## Documents size (in segments)

**<= 25** segments **82.95%** (7.6K documents)
**> 25** segments **17.05%** (1.6K documents)



## Documents by collection

**CC = 75.07%**
**IA = 24.93%**



- cc22 (3K)
- cc18 (2.2K)
- 19 Others (4K)

## Language Distribution

### Number of segments in the Latgalian (ltg) corpus



- Latgalian (ltg) - 55K
- Latvian (lv) - 54K
- Finnish (fi) - 5.6K
- Lithuanian (lt) - 5.5K
- English (en) - 5.2K
- German (de) - 4.4K
- Polish (pl) - 2.2K
- Slovenian (sl) - 1.9K
- Italian (it) - 1.9K
- Dutch (nl) - 1.8K
- 108 Others - 14K

### Percentage of segments in Latgalian (ltg) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (9.2K documents)

## Segment length distribution by token

≤ **49** tokens = **58K** segments | **61K** duplicates

> **50** tokens = **32K** segments | **13K** duplicates



Segments / Number of tokens in the segment

■ Unique segments    ■ Duplicated segments

## Segment noise distribution



| | |
|---|---|
| Too long | 0.97 % |
| Too short | 10.77 % |
| URLs | 3.85 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.48 % |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | i \| 118799    ir \| 55420    nu \| 38709    par \| 38067    kai \| 28443 |
| 2 | labot pirmkodu \| 2730    tys ir \| 2658    latgolys studentu \| 2636    par tū \| 2421    latgalīšu rokstu \| 2093 |
| 3 | latgolys studentu centrs \| 1875    latgalīšu rokstu volūdys \| 1133    latgalīšu kulturys goda \| 706    pi myusim latgolā \| 637    fikys fikys fikys \| 579 |
| 4 | fikys fikys fikys fikys \| 578    latgalīšu kulturys goda bolvys \| 356    nūvoda teritoriskais padaliņs latgolā \| 235    kyskys kys kyskys kys \| 232    kys kyskys kys kyskys \| 232 |
| 5 | fikys fikys fikys fikys fikys \| 577    kys kyskys kys kyskys kys \| 232    kyskys kys kyskys kys kyskys \| 224    latgalīšu kulturys ziņu portals lakuga \| 206    godā solu reorganizej kai pogostu \| 184 |

## 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 | | | | |