# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| als_Latn.jsonl.tsv | 9/23/2024 | Albanian (als) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 5,385,262 | 95,101,632 | 48,574,292 (51.08 %) | 3.2B | 16,005,838,206 | 16.0 GB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| evropaelire.org | 134K | 2.48% |
| zeriamerikes.com | 110K | 2.05% |
| wikipedia.org | 100K | 1.86% |
| botasot.info | 71K | 1.32% |
| albeu.com | 48K | 0.90% |
| blogspot.com | 43K | 0.80% |
| teksteshqip.com | 42K | 0.78% |
| telegrafi.com | 41K | 0.77% |
| koha.net | 37K | 0.69% |
| shqiperia.com | 33K | 0.61% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 2.3M | 41.79% |
| al | 1.1M | 20.44% |
| org | 474K | 8.81% |
| net | 411K | 7.62% |
| info | 305K | 5.66% |
| mk | 162K | 3.02% |
| tv | 153K | 2.84% |
| gov.al | 62K | 1.15% |
| com.al | 54K | 1.01% |
| ch | 48K | 0.90% |

## Register labels



- HI - 1.5%
- ID - 0.9%
- IN - 9.0%
- IP - 6.1%
- LY - 1.1%
- MIX - 2.0%
- NA - 61.7%
- OP - 8.1%
- SP - 1.2%
- UNK - 8.3%

🤖 **MT**:4.4% | 235K Documents

- HI_other - 0.8%
- HI_re - 0.7%
- ID_other - 0.9%
- IN_dtp - 3.2%
- IN_en - 1.9%
- IN_fi - 0.0%
- IN_lt - 0.5%
- IN_other - 3.4%
- IN_ra - 0.0%
- IP_ds - 5.0%
- IP_ed - 0.0%
- IP_other - 1.1%
- LY_other - 1.1%
- MIX - 2.0%
- NA_nb - 0.9%
- NA_ne - 49.8%
- NA_other - 4.2%
- NA_sr - 6.8%
- OP_av - 1.3%
- OP_ob - 2.4%
- OP_other - 2.1%
- OP_rs - 2.0%
- OP_rv - 0.3%
- SP_it - 0.8%
- SP_other - 0.4%
- UNK - 8.3%

## Documents size (in segments)

**<= 25** segments **84.81%** (4.6M documents)
**> 25** segments **15.19%** (818K documents)



## Documents by collection

**CC = 69.11%**
**IA = 30.89%**



- cc18 (845K)
- cc22 (1.6M)
- cc21 (578K)
- 18 Others (2.3M)

## Language Distribution

### Number of segments in the Albanian (als) corpus



- Albanian (sq) - 71M
- English (en) - 6.4M
- Italian (it) - 3.7M
- Lithuanian (lt) - 2.5M
- French (fr) - 1M
- Esperanto (eo) - 881K
- Serbian (sr) - 777K
- German (de) - 729K
- Spanish (es) - 613K
- Turkish (tr) - 605K
- 164 Others - 6.6M
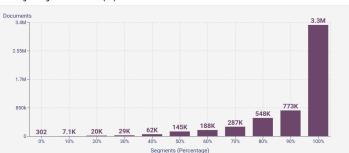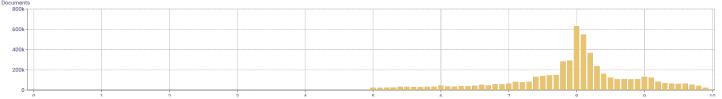
### Percentage of segments in Albanian (als) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (5.4M documents)

## Segment length distribution by token

Segments

6M
4.5M
3M
1.5M
0

Number of tokens in the segment
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

■ Unique segments  ■ Duplicated segments

## Segment noise distribution



| | |
|---|---|
| Too long | 1.34 % |
| Too short | 11.82 % |
| URLs | 1.76 % |
| Bad encoding | 0.01 % |
| Contains PII | 0.11 % |

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

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | është \| 21759174   shumë \| 8812607   kanë \| 6473161   të \| 6048064   duhet \| 5137758 |
| 2 | është shumë \| 404479   është bërë \| 375825   kanë qenë \| 364525   read more \| 318634   edi rama \| 293209 |
| 3 | herë të parë \| 274560   shtetet e bashkuara \| 260659   redakto tekstin burimor \| 237717   duhet të jetë \| 225196   republikës së kosovës \| 189623 |
| 4 | gjithnjë e më shumë \| 59240   luftës së dytë botërore \| 57162   luaj online flash lojë \| 47112   është shumë e rëndësishme \| 38496   sot e kësaj dite \| 32286 |
| 5 | miqtë tuaj më të mirë \| 68169   ndajnë këtë lojë me miqtë \| 67888   harroni të vlerësoni këtë game \| 55687   shtetet e bashkuara të amerikës \| 54702   shteteve të bashkuara të amerikës \| 45287 |

## 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 | | | | |