# HPLT Analytics report

**HPLT** Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| slv_Latn.jsonl.tsv | 6/7/2025 | Slovenian (sl) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 10,277,172 | 238,597,185 | 107,466,653 (45.04 %) | 6.4B | 35,029,805,659 | 33.47 GB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| wikipedia.org | 276K | 2.69% |
| sta.si | 226K | 2.20% |
| siol.net | 141K | 1.37% |
| delo.si | 130K | 1.26% |
| blogspot.com | 120K | 1.16% |
| dnevnik.si | 83K | 0.81% |
| metropolitan.si | 77K | 0.75% |
| slo-tech.com | 73K | 0.71% |
| rtvslo.si | 72K | 0.70% |
| agoda.com | 71K | 0.69% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| si | 6.3M | 61.21% |
| com | 2.2M | 21.01% |
| net | 558K | 5.43% |
| org | 555K | 5.40% |
| eu | 252K | 2.45% |
| info | 120K | 1.17% |
| se | 28K | 0.27% |
| je | 24K | 0.24% |
| tv | 22K | 0.21% |
| at | 19K | 0.19% |

## Register labels



- HI - 3.5%
- ID - 2.2%
- IN - 15.4%
- IP - 29.3%
- LY - 0.1%
- MIX - 5.2%
- NA - 29.7%
- OP - 6.2%
- SP - 0.7%
- UNK - 7.8%

**MT**: 5.0% | 514K Documents

- HI_other - 2.0%
- HI_re - 1.5%
- ID_other - 2.2%
- IN_dtp - 5.4%
- IN_en - 2.7%
- IN_fi - 0.0%
- IN_lt - 1.7%
- IN_other - 5.6%
- IN_ra - 0.0%
- IP_ds - 25.4%
- IP_ed - 0.0%
- IP_other - 3.9%
- LY_other - 0.1%
- MIX - 5.2%
- NA_nb - 7.9%
- NA_ne - 14.4%
- NA_other - 3.8%
- NA_sr - 3.5%
- OP_av - 1.3%
- OP_ob - 1.4%
- OP_other - 1.3%
- OP_rs - 0.9%
- OP_rv - 1.3%
- SP_it - 0.6%
- SP_other - 0.1%
- UNK - 7.8%

## Documents size (in segments)

<= 25 segments **78.6%** (8.1M documents)
> 25 segments **21.4%** (2.2M documents)



## Documents by collection

CC = 70.36%
IA = 29.64%



## Language Distribution

### Number of segments in the Slovenian (sl) corpus



- Slovenian (sl) - 179M
- English (en) - 12M
- Serbian (sr) - 11M
- Italian (it) - 7M
- Czech (cs) - 4.5M
- Polish (pl) - 3.6M
- German (de) - 3.4M
- French (fr) - 2.3M
- Slovak (sk) - 2M
- Esperanto (eo) - 1.6M
- 165 Others - 12M

### Percentage of segments in Slovenian (sl) inside documents



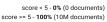## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (10M documents)

## Segment length distribution by token

Segments

- Unique segments
- Duplicated segments

Number of tokens in the segment

## Segment noise distribution



| | |
|---|---|
| Too long | 0.88 % |
| Too short | 15.68 % |
| URLs | 1.86 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.69 % |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | lahko \| 25982227   zelo \| 6691629   a \| 6631583   leta \| 6410881   zato \| 6406868 |
| 2 | uredi kodo \| 974504   spletni strani \| 834745   republike slovenije \| 684698   spletne strani \| 655489   osebnih podatkov \| 540481 |
| 3 | uradni list rs \| 388741   člena tega zakona \| 187426   državna revizijska komisija \| 166992   dodaj v košarico \| 159247   odstavka tega člena \| 157674 |
| 4 | celotna novica je dostopna \| 136900   evropskega parlamenta in sveta \| 64408   e-poštni naslov je zaščiten \| 55487   obvezno pokojninsko in invalidsko \| 48545   uradnem listu republike slovenije \| 47721 |
| 5 | novica je dostopna le naročnikom \| 136880   naslov je zaščiten proti smetenju \| 55436   kazensko odgovoren za javno spodbujanje \| 44940   posameznik kazensko odgovoren za javno \| 44926   odgovoren za javno spodbujanje sovraštva \| 44788 |

## 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 | | | | |