# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| mkd_Cyrl.jsonl.tsv | 9/24/2024 | Macedonian (mk) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 3,565,647 | 57,008,331 | 27,635,192 (48.48 %) | 1.7B | 9,386,182,083 | 15.64 GB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| wikipedia.org | 180K | 5.05% |
| slobodnaevropa.mk | 79K | 2.23% |
| voanews.com | 57K | 1.61% |
| kurir.mk | 55K | 1.53% |
| daily.mk | 53K | 1.48% |
| netpress.com.mk | 42K | 1.18% |
| rbth.com | 39K | 1.10% |
| hitportal.com.mk | 39K | 1.08% |
| republika.mk | 36K | 1.02% |
| kafepauza.mk | 34K | 0.94% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| mk | 2M | 55.09% |
| com | 544K | 15.27% |
| com.mk | 433K | 12.14% |
| org | 263K | 7.37% |
| gov.mk | 86K | 2.40% |
| org.mk | 80K | 2.25% |
| edu.mk | 39K | 1.08% |
| net | 32K | 0.90% |
| info | 14K | 0.40% |
| news | 13K | 0.36% |

## Register labels



- HI - 2.7%
- ID - 0.4%
- IN - 12.7%
- IP - 9.6%
- LY - 0.2%
- MIX - 3.2%
- NA - 54.9%
- OP - 6.6%
- SP - 0.9%
- UNK - 8.8%

**MT**:5.8% | 206K Documents

- HI_other - 1.3%
- HI_re - 1.4%
- ID_other - 0.4%
- IN_dtp - 3.0%
- IN_en - 5.1%
- IN_fi - 0.0%
- IN_lt - 0.7%
- IN_other - 3.8%
- IN_ra - 0.1%
- IP_ds - 7.8%
- IP_ed - 0.0%
- IP_other - 1.8%
- LY_other - 0.2%
- MIX - 3.2%
- NA_nb - 1.5%
- NA_ne - 42.6%
- NA_other - 4.1%
- NA_sr - 6.7%
- OP_av - 2.2%
- OP_ob - 1.5%
- OP_other - 1.3%
- OP_rs - 1.1%
- OP_rv - 0.6%
- SP_it - 0.6%
- SP_other - 0.2%
- UNK - 8.8%

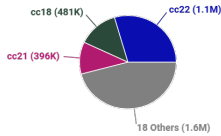## Documents size (in segments)

**<= 25** segments **85.34%** (3M documents)
**> 25** segments **14.66%** (523K documents)



## Documents by collection

CC = 63.48%
IA = 36.52%



- cc18 (481K)
- cc22 (1.1M)
- cc21 (396K)
- 18 Others (1.6M)

## Language Distribution

### Number of segments in the Macedonian (mk) corpus



- Macedonian (mk) - 50M
- English (en) - 1.6M
- Italian (it) - 1.3M
- Serbian (sr) - 1.2M
- Russian (ru) - 829K
- Bulgarian (bg) - 468K
- Ukrainian (uk) - 326K
- German (de) - 243K
- French (fr) - 219K
- Spanish (es) - 69K
- 162 Others - 726K

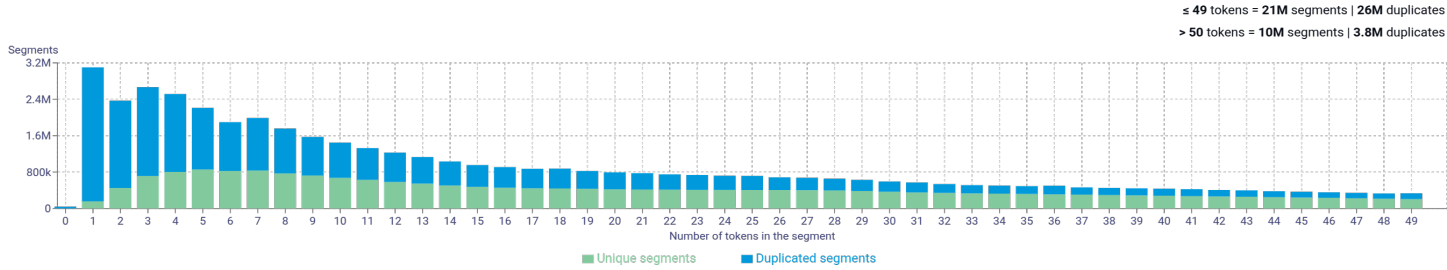### Percentage of segments in Macedonian (mk) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (3.6M documents)

## Segment length distribution by token

≤ 49 tokens = **21M** segments | **26M** duplicates
> **50** tokens = **10M** segments | **3.8M** duplicates



Segments

■ Unique segments   ■ Duplicated segments

Number of tokens in the segment

## Segment noise distribution



| | |
|---|---|
| Too long | 0.94 % |
| Too short | 12.21 % |
| URLs | 1.19 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.23 % |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | година \| 3930388    македонија \| 2807325    време \| 2074517    години \| 1767968    дел \| 1730470 |
| 2 | република македонија \| 539687    уреди извор \| 506624    станува збор \| 277149    ве молиме \| 245983    голем број \| 245318 |
| 3 | република северна македонија \| 88737    можат да бидат \| 85258    лигата на шампионите \| 60021    втората светска војна \| 59704    владата на република \| 57927 |
| 4 | владата на република македонија \| 41668    министерството за внатрешни работи \| 40710    труд и социјална политика \| 37252    собранието на република македонија \| 35902    авторски текстови е казниво \| 30930 |
| 5 | текстови е казниво со закон \| 30976    неделата директно во вашето сандаче \| 28528    најдобрите стории на неделата директно \| 28528    дистрибуира во каква било форма \| 25188    писмена дозвола од македонската информативна \| 25171 |

## 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 | | | | |