# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| kir_Cyrl.jsonl.tsv | 9/22/2024 | Kyrgyz (ky) |

## Volumes

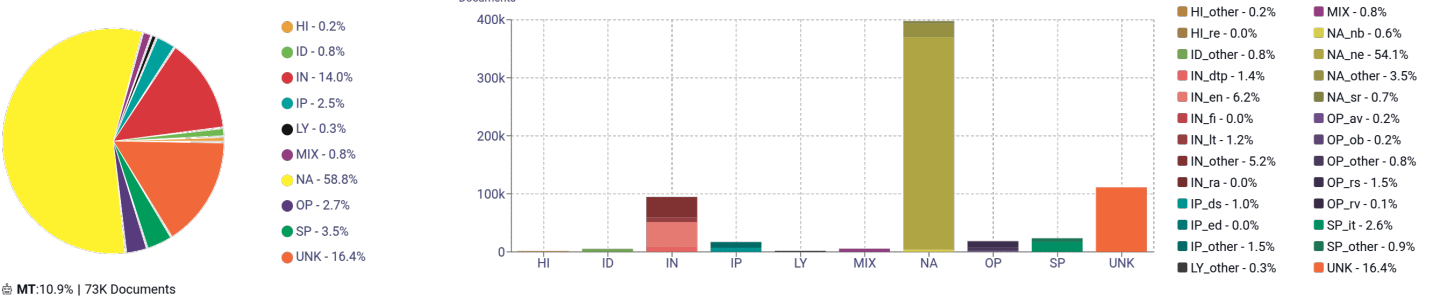| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 676,111 | 10,041,028 | 6,263,608 (62.38 %) | 312M | 1,916,124,485 | 3.27 GB |

### Top 10 domains

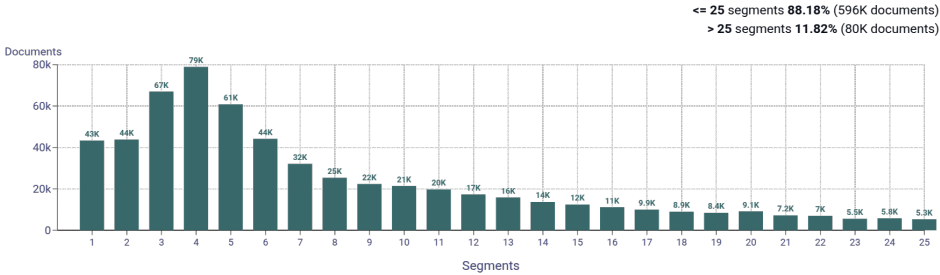| Domain | Docs | % of total |
|---|---|---|
| azattyk.org | 201K | 29.69% |
| wikipedia.org | 54K | 7.92% |
| turmush.kg | 22K | 3.19% |
| sputnik.kg | 12K | 1.71% |
| kyrgyztoday.org | 10K | 1.48% |
| kabar.kg | 9K | 1.34% |
| tyup.net | 7.2K | 1.07% |
| bagyt.kg | 6.6K | 0.98% |
| ykt.ru | 6.6K | 0.97% |
| centralasian.org | 6.2K | 0.91% |

### Top 10 TLDs

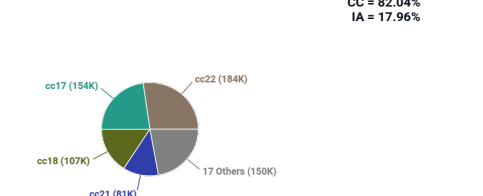| Domain | Docs | % of total |
|---|---|---|
| org | 293K | 43.29% |
| kg | 215K | 31.77% |
| com | 54K | 7.93% |
| ru | 45K | 6.59% |
| net | 9.9K | 1.46% |
| gov.kg | 8.6K | 1.27% |
| media | 6.2K | 0.92% |
| asia | 5.4K | 0.79% |
| info | 5.3K | 0.78% |
| news | 4.1K | 0.61% |

## Register labels



- HI - 0.2%
- ID - 0.8%
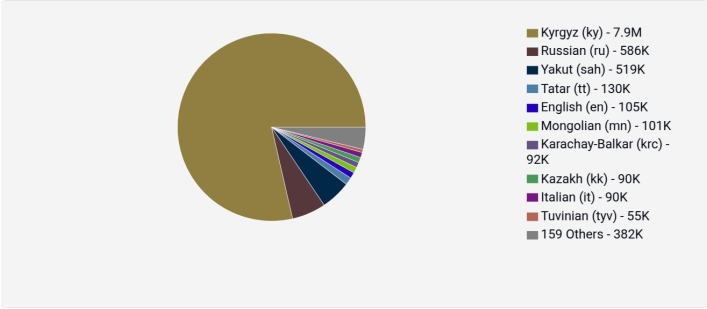- IN - 14.0%
- IP - 2.5%
- LY - 0.3%
- MIX - 0.8%
- NA - 58.8%
- OP - 2.7%
- SP - 3.5%
- UNK - 16.4%

🤖 **MT**:10.9% | 73K Documents

- HI_other - 0.2%
- HI_re - 0.0%
- ID_other - 0.8%
- IN_dtp - 1.4%
- IN_en - 6.2%
- IN_fi - 0.0%
- IN_lt - 1.2%
- IN_other - 5.2%
- IN_ra - 0.0%
- IP_ds - 1.0%
- IP_ed - 0.0%
- IP_other - 1.5%
- LY_other - 0.3%
- MIX - 0.8%
- NA_nb - 0.6%
- NA_ne - 54.1%
- NA_other - 3.5%
- NA_sr - 0.7%
- OP_av - 0.2%
- OP_ob - 0.2%
- OP_other - 0.8%
- OP_rs - 1.5%
- OP_rv - 0.1%
- SP_it - 2.6%
- SP_other - 0.9%
- UNK - 16.4%

## Documents size (in segments)

**<= 25** segments **88.18%** (596K documents)
**> 25** segments **11.82%** (80K documents)



## Documents by collection

**CC = 82.04%**
**IA = 17.96%**



- cc17 (154K)
- cc22 (184K)
- cc18 (107K)
- cc21 (81K)
- 17 Others (150K)

## Language Distribution

### Number of segments in the Kyrgyz (ky) corpus



- Kyrgyz (ky) - 7.9M
- Russian (ru) - 586K
- Yakut (sah) - 519K
- Tatar (tt) - 130K
- English (en) - 105K
- Mongolian (mn) - 101K
- Karachay-Balkar (krc) - 92K
- Kazakh (kk) - 90K
- Italian (it) - 90K
- Tuvinian (tyv) - 55K
- 159 Others - 382K

### Percentage of segments in Kyrgyz (ky) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (676K documents)

## Segment length distribution by token

≤ 49 tokens = **4.9M** segments | **3.3M** duplicates
> 50 tokens = **1.8M** segments | **496K** duplicates



■ Unique segments   ■ Duplicated segments

## Segment noise distribution



| | |
|---|---|
| Too long | 1.27 % |
| Too short | 10.00 % |
| URLs | 0.80 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.11 % |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | жылы \| 356235    иш \| 344897    а \| 340362    республикасынын \| 333084    жөнүндө \| 329113 |
| 2 | билим берүү \| 90359    орун басары \| 63181    булагын оңдоо \| 53054    жылдан бери \| 52625    ички иштер \| 48316 |
| 3 | жергиликтүү өз алдынча \| 23652    тил жана энциклопедия \| 19418    республикасынын жогорку кеңешинин \| 17184    тышкы иштер министри \| 15110    президент алмазбек атамбаев \| 12351 |
| 4 | тил жана энциклопедия борбору \| 19395    жергиликтүү өз алдынча башкаруу \| 19337    билим берүү жана илим \| 11521    тик так тик так \| 8035    так тик так тик \| 8035 |
| 5 | тик так тик так тик \| 8030    так тик так тик так \| 8020    жергиликтүү өз алдынча башкаруу органдарынын \| 5659    жергиликтүү өз алдынча башкаруу органдары \| 4119    билим берүү жана илим министрлигинин \| 3473 |

## 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 | | | | |