# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| kaz_Cyrl.jsonl.tsv | 9/21/2024 | Kazakh (kk) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 2,637,363 | 81,006,479 | 42,567,471 (52.55 %) | 1.8B | 11,053,418,553 | 18.78 GB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| wikipedia.org | 154K | 5.84% |
| azattyq.org | 143K | 5.43% |
| strategy2050.kz | 116K | 4.41% |
| kodeksy-kz.com | 62K | 2.37% |
| tengrinews.kz | 58K | 2.22% |
| nur.kz | 57K | 2.17% |
| inform.kz | 50K | 1.91% |
| stud.kz | 47K | 1.77% |
| massaget.kz | 37K | 1.40% |
| baq.kz | 35K | 1.32% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| kz | 1.8M | 67.51% |
| org | 350K | 13.25% |
| com | 216K | 8.21% |
| ru | 77K | 2.93% |
| gov.kz | 64K | 2.41% |
| net | 29K | 1.10% |
| info | 29K | 1.08% |
| edu.kz | 23K | 0.86% |
| mobi | 6.4K | 0.24% |
| uz | 6K | 0.23% |

## Register labels



- HI - 0.6%
- ID - 0.5%
- IN - 28.3%
- IP - 3.6%
- LY - 0.3%
- MIX - 1.3%
- NA - 47.6%
- OP - 3.4%
- SP - 1.5%
- UNK - 12.9%

🤖 **MT**:7.4% | 194K Documents

- HI_other - 0.5%
- HI_re - 0.1%
- ID_other - 0.5%
- IN_dtp - 3.6%
- IN_en - 5.4%
- IN_fi - 0.0%
- IN_lt - 4.6%
- IN_other - 14.6%
- IN_ra - 0.1%
- IP_ds - 1.9%
- IP_ed - 0.0%
- IP_other - 1.7%
- LY_other - 0.3%
- MIX - 1.3%
- NA_nb - 1.1%
- NA_ne - 40.3%
- NA_other - 4.4%
- NA_sr - 1.8%
- OP_av - 0.4%
- OP_ob - 0.4%
- OP_other - 1.1%
- OP_rs - 1.4%
- OP_rv - 0.1%
- SP_it - 1.1%
- SP_other - 0.4%
- UNK - 12.9%

## Documents size (in segments)

**<= 25** segments **78.51%** (2.1M documents)
**> 25** segments **21.49%** (567K documents)



## Documents by collection

**CC = 68.75%**
**IA = 31.25%**



- cc18 (433K)
- cc22 (696K)
- cc21 (320K)
- 18 Others (1.2M)

## Language Distribution

### Number of segments in the Kazakh (kk) corpus



- Kazakh (kk) - 71M
- Russian (ru) - 3.9M
- Italian (it) - 1.2M
- Ukrainian (uk) - 1.1M
- English (en) - 988K
- Tatar (tt) - 425K
- German (de) - 382K
- French (fr) - 281K
- Serbian (sr) - 216K
- Kyrgyz (ky) - 196K
- 165 Others - 1.4M

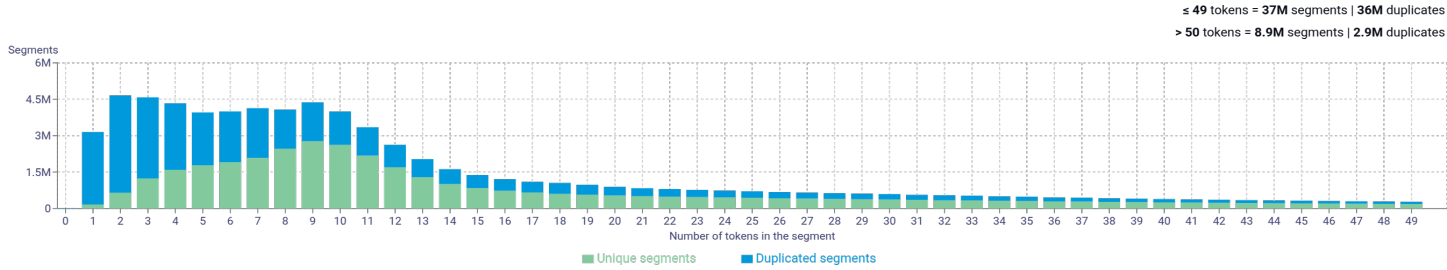### Percentage of segments in Kazakh (kk) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (2.6M documents)

## Segment length distribution by token

Segments

6M
4.5M
3M
1.5M
0

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

Number of tokens in the segment

■ Unique segments  ■ Duplicated segments

## Segment noise distribution



| Category | Value |
|---|---|
| Too long | 0.88 % |
| Too short | 14.14 % |
| URLs | 1.14 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.15 % |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | және \| 17900838    да \| 5290698    бір \| 5108395    бойынша \| 4820816    қазақстан \| 4770376 |
| 2 | қазақстан республикасының \| 1407502    болып табылады \| 1009956    қазақстан республикасы \| 946957    білім беру \| 703405    басқа да \| 560917 |
| 3 | өткен соң қолданысқа \| 193307    және басқа да \| 191300    он күн өткен \| 166733    күнтізбелік он күн \| 146253    ресми жарияланған күнінен \| 144430 |
| 4 | күн өткен соң қолданысқа \| 182442    өткен соң қолданысқа енгізіледі \| 158912    күнтізбелік он күн өткен \| 137602    алғашқы ресми жарияланған күнінен \| 134652    жарияланған күнінен кейін күнтізбелік \| 93685 |
| 5 | он күн өткен соң қолданысқа \| 166096    күн өткен соң қолданысқа енгізіледі \| 150197    ресми жарияланған күнінен кейін күнтізбелік \| 93542    күнінен кейін күнтізбелік он күн \| 85030    жарияланған күнінен кейін күнтізбелік он \| 84740 |

## 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 | | | | |