# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| ukr_Cyrl.jsonl.tsv | 6/26/2025 | Ukrainian (uk) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 47,395,779 | 1,168,793,792 | 537,002,359 (45.95 %) | 32B | 181,746,006,077 | 306.5 GB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| wikipedia.org | 858K | 1.81% |
| unian.ua | 787K | 1.66% |
| uapatents.com | 448K | 0.95% |
| korrespondent.net | 425K | 0.90% |
| radiosvoboda.org | 367K | 0.77% |
| co.ua | 350K | 0.74% |
| 24tv.ua | 334K | 0.70% |
| pp.ua | 325K | 0.69% |
| rada.gov.ua | 301K | 0.64% |
| referatcentral.... | 298K | 0.63% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 8.4M | 17.75% |
| ua | 8.1M | 17.04% |
| com.ua | 7.6M | 16.04% |
| org.ua | 2.8M | 5.88% |
| org | 2.8M | 5.88% |
| gov.ua | 2.7M | 5.76% |
| net | 1.9M | 4.03% |
| in.ua | 1.7M | 3.61% |
| info | 1.5M | 3.14% |
| ru | 1.2M | 2.57% |

## Register labels



Documents

- HI - 3.5%
- ID - 0.6%
- IN - 22.3%
- IP - 15.3%
- LY - 0.2%
- MIX - 3.2%
- NA - 45.2%
- OP - 4.3%
- SP - 0.7%
- UNK - 4.7%

- HI_other - 2.6%
- HI_re - 1.0%
- ID_other - 0.6%
- IN_dtp - 4.4%
- IN_en - 2.4%
- IN_fi - 0.0%
- IN_lt - 2.9%
- IN_other - 12.4%
- IN_ra - 0.3%
- IP_ds - 13.4%
- IP_ed - 0.0%
- IP_other - 1.9%
- LY_other - 0.2%
- MIX - 3.2%
- NA_nb - 1.9%
- NA_ne - 37.1%
- NA_other - 4.4%
- NA_sr - 1.7%
- OP_av - 1.0%
- OP_ob - 0.9%
- OP_other - 1.2%
- OP_rs - 0.7%
- OP_rv - 0.4%
- SP_it - 0.5%
- SP_other - 0.2%
- UNK - 4.7%
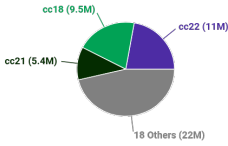
**MT**:1.2% | 573K Documents

## Documents size (in segments)

<= 25 segments **77.1%** (37M documents)
> 25 segments **22.9%** (11M documents)



## Documents by collection
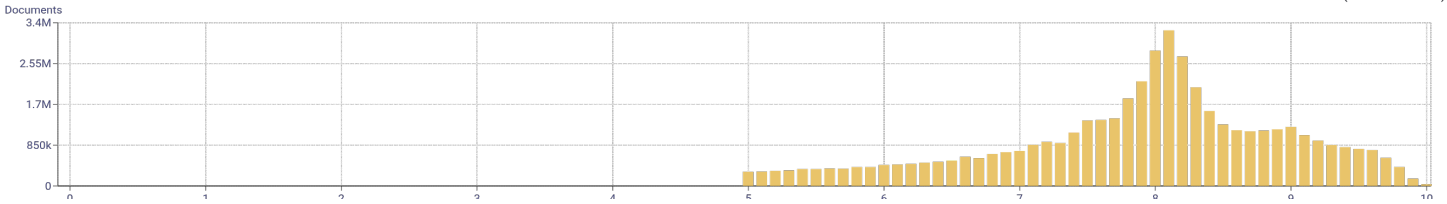
CC = 65.40%
IA = 34.60%



- cc18 (9.5M)
- cc22 (11M)
- cc21 (5.4M)
- 18 Others (22M)

## Language Distribution

### Number of segments in the Ukrainian (uk) corpus



- Ukrainian (uk) - 1.1B
- Russian (ru) - 36M
- Italian (it) - 17M
- English (en) - 17M
- German (de) - 5M
- French (fr) - 4.7M
- Serbian (sr) - 3M
- Belarusian (be) - 2.3M
- Macedonian (mk) - 1.4M
- Bulgarian (bg) - 1.2M
- 165 Others - 12M

### Percentage of segments in Ukrainian (uk) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (47M documents)
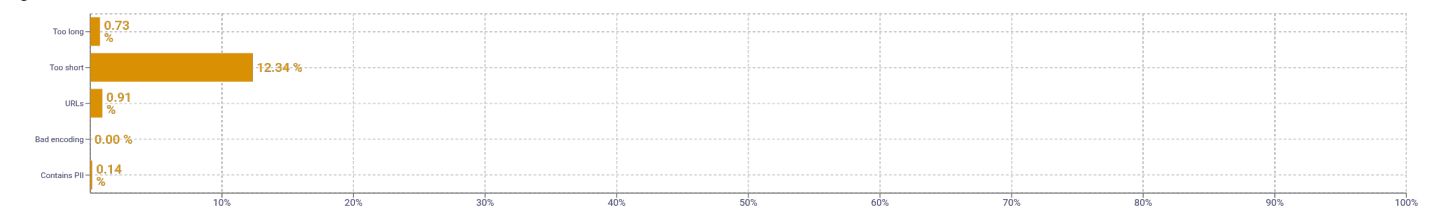
## Segment length distribution by token

≤ 49 tokens = **990M** segments | **556M** duplicates
> 50 tokens = **179M** segments | **78M** duplicates



Legend: ■ Unique segments ■ Duplicated segments
X-axis: Number of tokens in the segment
Y-axis: Segments

## Segment noise distribution



| Category | Value |
|---|---|
| Too long | 0.73 % |
| Too short | 12.34 % |
| URLs | 0.91 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.14 % |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | україни \| 73192772   його \| 60796919   при \| 50653118   із \| 47680333   він \| 43800126 |
| 2 | під час \| 16668196   може бути \| 8655692   при цьому \| 8407996   крім того \| 5494569   можуть бути \| 4196695 |
| 3 | верховної ради україни \| 1532593   відповідно до закону \| 657530   включає в себе \| 626415   другої світової війни \| 595833   автономної республіки крим \| 497519 |
| 4 | внесення змін до деяких \| 396244   відповідно до закону україни \| 332998   внесення змін до закону \| 213279   вітаю тебе з днем \| 208832   внесеними згідно із законом \| 180985 |
| 5 | внесення змін до деяких законодавчих \| 250845   деяких законодавчих актів україни щодо \| 191895   внесення змін до закону україни \| 187395   вітаю тебе з днем народження \| 174670   верховної ради україни з питань \| 149134 |

## 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 | | | | |