# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|--------|------|----------|
| azb_Arab.jsonl.tsv | 9/27/2024 | South Azerbaijani (azb) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|------|----------|-----------------|--------|------------|------|
| 66,112 | 2,389,200 | 1,055,552 (44.18 %) | 49M | 257,866,139 | 446.12 MB |

### Top 10 domains

| Domain | Docs | % of total |
|--------|------|------------|
| blogfa.com | 14K | 21.68% |
| wikipedia.org | 11K | 16.86% |
| axar.az | 4.7K | 7.09% |
| trt.net.tr | 4.5K | 6.80% |
| arzublog.com | 2.9K | 4.45% |
| ishiq.net | 2.9K | 4.35% |
| baybak.com | 2.1K | 3.12% |
| bilimsesi.com | 1.5K | 2.23% |
| blogsky.com | 1.4K | 2.12% |
| mihanblog.com | 1.2K | 1.88% |

### Top 10 TLDs

| Domain | Docs | % of total |
|--------|------|------------|
| com | 33K | 50.20% |
| org | 13K | 19.92% |
| ir | 5K | 7.60% |
| az | 4.7K | 7.09% |
| net.tr | 4.5K | 6.80% |
| net | 3.3K | 4.96% |
| info | 380 | 0.57% |
| biz | 271 | 0.41% |
| se | 253 | 0.38% |
| ca | 209 | 0.32% |

## Register labels



- HI - 0.0%
- ID - 0.9%
- IN - 17.3%
- IP - 0.4%
- LY - 6.4%
- MIX - 0.0%
- NA - 0.7%
- OP - 0.8%
- SP - 0.1%
- UNK - 73.3%

**MT**:43.3% | 29K Documents

- HI_other - 0.0%
- HI_re - 0.0%
- ID_other - 0.9%
- IN_dtp - 0.5%
- IN_en - 8.2%
- IN_fi - 0.0%
- IN_lt - 0.0%
- IN_other - 8.6%
- IN_ra - 0.0%
- IP_ds - 0.2%
- IP_ed - 0.0%
- IP_other - 0.2%
- LY_other - 6.4%
- MIX - 0.0%
- NA_nb - 0.1%
- NA_ne - 0.2%
- NA_other - 0.2%
- NA_sr - 0.2%
- OP_av - 0.0%
- OP_ob - 0.0%
- OP_other - 0.2%
- OP_rs - 0.6%
- OP_rv - 0.0%
- SP_it - 0.0%
- SP_other - 0.1%
- UNK - 73.3%

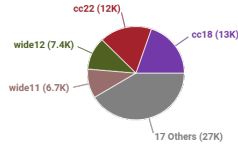## Documents size (in segments)

<= **25** segments **69.41%** (46K documents)
> **25** segments **30.59%** (20K documents)



## Documents by collection

CC = 48.77%
IA = 51.23%



- cc22 (12K)
- cc18 (13K)
- wide12 (7.4K)
- wide11 (6.7K)
- 17 Others (27K)

## Language Distribution

### Number of segments in the South Azerbaijani (azb) corpus



- South Azerbaijani (azb) - 1.4M
- Persian (fa) - 679K
- Arabic (ar) - 102K
- English (en) - 29K
- Azerbaijani (az) - 28K
- Urdu (ur) - 20K
- Western Panjabi (pnb) - 19K
- Turkish (tr) - 14K
- Egyptian Arabic (arz) - 12K
- Central Kurdish (ckb) - 10K
- 140 Others - 71K

### Percentage of segments in South Azerbaijani (azb) inside documents
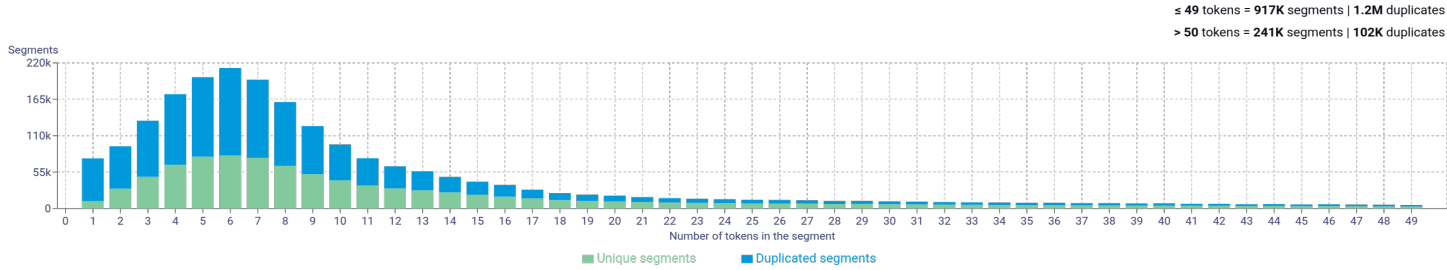


## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (66K documents)

## Segment length distribution by token

≤ 49 tokens = **917K** segments | **1.2M** duplicates

\> 50 tokens = **241K** segments | **102K** duplicates

Segments

220k

165k

110k

55k

0

Number of tokens in the segment

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

■ Unique segments   ■ Duplicated segments

## Segment noise distribution

| | |
|---|---|
| Too long | 0.88 % |
| Too short | 11.34 % |
| URLs | 0.50 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.03 % |

10%  20%  30%  40%  50%  60%  70%  80%  90%  100%

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | ابله \| 172533    دا \| 156791    کیمی \| 109314    جوخ \| 98960    دیر \| 92750 |
| 2 | جی ابله \| 14647    read more \| 9442    گۆره \| 6535    بونا گۆره \| 6317    داها جوخ \| 5500    داها آرتیق \| 5500 |
| 3 | ویکیپدیاسینین ایشلدنلری طرفیندن \| 2718    ایشلد نلری طرفیندن یارا نمیش \| 2740    آرتیق بیلگیلر تا با بیلرسینیز \| 3078    داها آرتیق بیلگیلر \| 3079    گۆره داها آرتیق \| 3082 |
| 4 | ایران ممالیکی محروسه سینده \| 1537    گۆره داها آرتیق بیلگیلر \| 3078    داها آرتیق بیلگیلر تا با بیلرسینیز \| 3078    داها آرتیق بیلگیلر تا با بیلرسینیز \| 2718    ویکیپدیاسینین ایشلدنلری طرفیندن \| 2099    اینگیلیسجه ویکیپدیاسینین ایشلدنلری طرفیندن \| 2099 |
| 5 | ایلام ایلام ایلام ایلام ایلام \| 414    ی تیمی ترکیبیننده چیخیش ائدیب \| 3078    گۆره داها آرتیق بیلگیلر تا با بیلرسینیز \| 2099    اینگیلیسجه ویکیپدیاسینین ایشلدنلری طرفیندن یارا نمیش \| 561    د ایله ابله و بی موبایل یئشیک \| 496    یئشیک موبایل و بی ابله د \| 496 |

## 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 | | | | |