# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|--------|------|----------|
| ace_Arab.jsonl.tsv | 12/3/2024 | Acehnese (ace) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|------|----------|-----------------|--------|------------|------|
| 16 | 117 | 93 (79.49 %) | 12K | 49,626 | 84.34 KB |

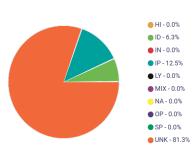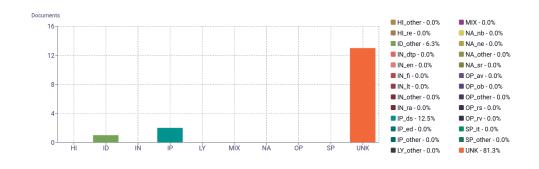## Top 10 domains

| Domain | Docs | % of total |
|--------|------|-----------|
| iran-eng.ir | 4 | 25.00% |
| blogspot.com | 2 | 12.50% |
| bayardservicewe... | 1 | 6.25% |
| utusanmelayu.co... | 1 | 6.25% |
| office-converte... | 1 | 6.25% |
| salemalanzi.sa | 1 | 6.25% |
| onlineyoutube.com | 1 | 6.25% |
| all.biz | 1 | 6.25% |
| blogspot.it | 1 | 6.25% |
| blogspot.my | 1 | 6.25% |

## Top 10 TLDs

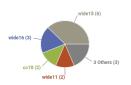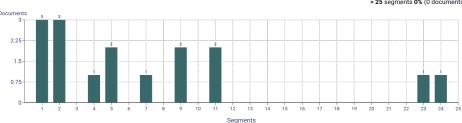| Domain | Docs | % of total |
|--------|------|-----------|
| com | 6 | 37.50% |
| ir | 4 | 25.00% |
| com.my | 1 | 6.25% |
| sa | 1 | 6.25% |
| biz | 1 | 6.25% |
| it | 1 | 6.25% |
| my | 1 | 6.25% |
| net | 1 | 6.25% |

## Register labels



- HI - 0.0%
- ID - 6.3%
- IN - 0.0%
- IP - 12.5%
- LY - 0.0%
- MIX - 0.0%
- NA - 0.0%
- OP - 0.0%
- SP - 0.0%
- UNK - 81.3%

🤖 **MT**:56.3% | 9 Documents

- HI_other - 0.0%
- HI_re - 0.0%
- ID_other - 6.3%
- IN_dtp - 0.0%
- IN_en - 0.0%
- IN_fi - 0.0%
- IN_lt - 0.0%
- IN_other - 0.0%
- IN_ra - 0.0%
- IP_ds - 12.5%
- IP_ed - 0.0%
- IP_other - 0.0%
- LY_other - 0.0%
- MIX - 0.0%
- NA_nb - 0.0%
- NA_ne - 0.0%
- NA_other - 0.0%
- NA_sr - 0.0%
- OP_av - 0.0%
- OP_ob - 0.0%
- OP_other - 0.0%
- OP_rs - 0.0%
- OP_rv - 0.0%
- SP_it - 0.0%
- SP_other - 0.0%
- UNK - 81.3%

## Documents size (in segments)

<= **25** segments **100%** (16 documents)
> **25** segments **0%** (0 documents)



## Documents by collection

CC = 18.75%
IA = 81.25%



## Language Distribution

### Number of segments in the Acehnese (ace) corpus



- Arabic (ar) - 76
- Persian (fa) - 16
- French (fr) - 7
- English (en) - 6
- Malay (ms) - 4
- German (de) - 3
- Egyptian Arabic (arz) - 3
- Indonesian (id) - 2

*Acehnese (ace) identification might be inaccurate because it is not supported by FastSpell

### Percentage of segments in Acehnese (ace) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (16 documents)

## Segment length distribution by token

Segments

Number of tokens in the segment

■ Unique segments  ■ Duplicated segments

## Segment noise distribution



| | |
|---|---|
| Too long | 7.69 % |
| Too short | 9.40 % |
| URLs | 0.00 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.85 % |

## Frequent n-grams

| Size | n-grams |
|------|---------|
| 1 | کن \| 1683    مدون \| 1680    کحاس \| 1530    خو نئون \| 1530    اوک \| 756 |
| 2 | مدون کن \| 1676    مدون کن \| 1673    خو نئون کحاس \| 1530    اوک اوک \| 755    ل ppm \| 336 |
| 3 | کن مدون کن \| 1670    مدون کن مدون \| 1666    مدون کن مدون \| 754    اوک اوک اوک \| 56    هیبوندای جینیسیس کوبیه \| 24    چت ازئتا ازئتا چت \| 24 |
| 4 | مدون کن مدون کن \| 1665    کن مدون کن مدون \| 1660    کن مدون کن مدون \| 753    اوک اوک اوک اوک \| 13    سعر هیبوندای جینیسیس کوبیه \| 8    فورتل این اکن دافت \| 8 |
| 5 | کن مدون کن مدون کن \| 1659    مدون کن مدون کن مدون \| 1655    مدون کن مدون کن مدون \| 752    اوک اوک اوک اوک اوک \| 6    مدونکن مدون کن مدون کن \| 6 |

# 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 | | | | |