# HPLT Analytics report

## General overview

| Corpus | Date | Language |
|---|---|---|
| hplt-v3-bjn_Arab | 9/17/2025 | Banjar |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 1,306 | 30,400 | 24,867 (81.80 %) | 1.1M | 4,582,094 | 7.88 MB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| utusanmelayu.co... | 662 | 50.69% |
| blogspot.com | 153 | 11.72% |
| wordpress.com | 147 | 11.26% |
| ahmadalikarim.com | 54 | 4.13% |
| utusantv.com | 33 | 2.53% |
| wikimedia.org | 19 | 1.45% |
| harakahdaily.net | 12 | 0.92% |
| ulamasedunia.org | 11 | 0.84% |
| co.cc | 11 | 0.84% |
| urusniaga.my | 9 | 0.69% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com.my | 667 | 51.07% |
| com | 475 | 36.37% |
| org | 59 | 4.52% |
| my | 30 | 2.30% |
| net | 25 | 1.91% |
| cc | 11 | 0.84% |
| moe | 7 | 0.54% |
| org.my | 6 | 0.46% |
| eu | 4 | 0.31% |
| edu.my | 4 | 0.31% |

## Documents size (in segments) ⓘ

≤ 25 segments **84.69%** (1.1K documents)
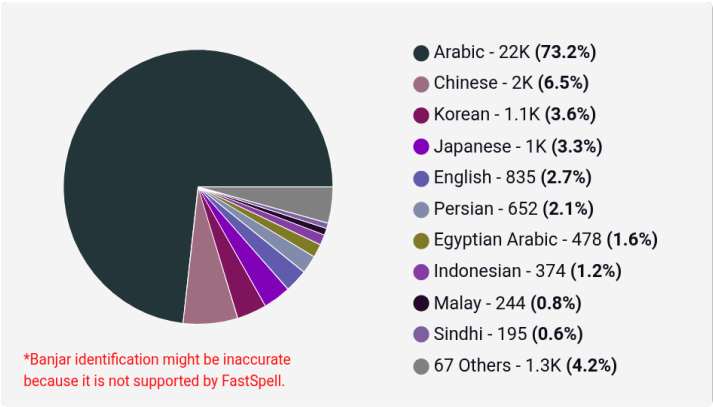> 25 segments **15.31%** (200 documents)
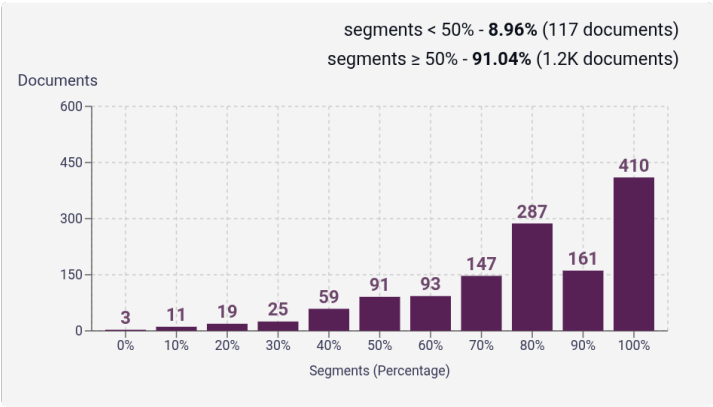


## Document collections

CC = **91.58%**
IA = **8.42%**



CC-MAIN-2018-39 (159)
CC-MAIN-2018-2...
61 Others (795)

## Language Distribution

### Number of segments in the Banjar corpus



- Arabic - 22K **(73.2%)**
- Chinese - 2K **(6.5%)**
- Korean - 1.1K **(3.6%)**
- Japanese - 1K **(3.3%)**
- English - 835 **(2.7%)**
- Persian - 652 **(2.1%)**
- Egyptian Arabic - 478 **(1.6%)**
- Indonesian - 374 **(1.2%)**
- Malay - 244 **(0.8%)**
- Sindhi - 195 **(0.6%)**
- 67 Others - 1.3K **(4.2%)**

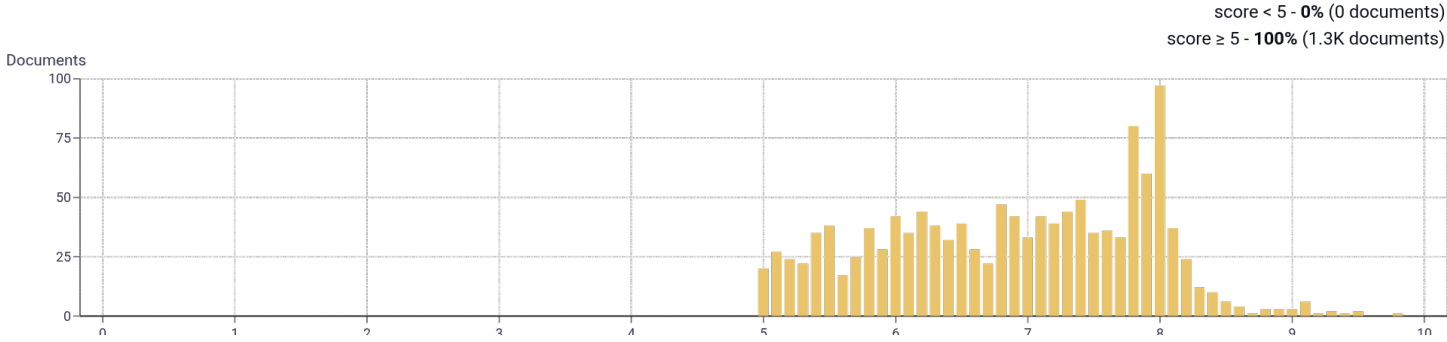*Banjar identification might be inaccurate because it is not supported by FastSpell.
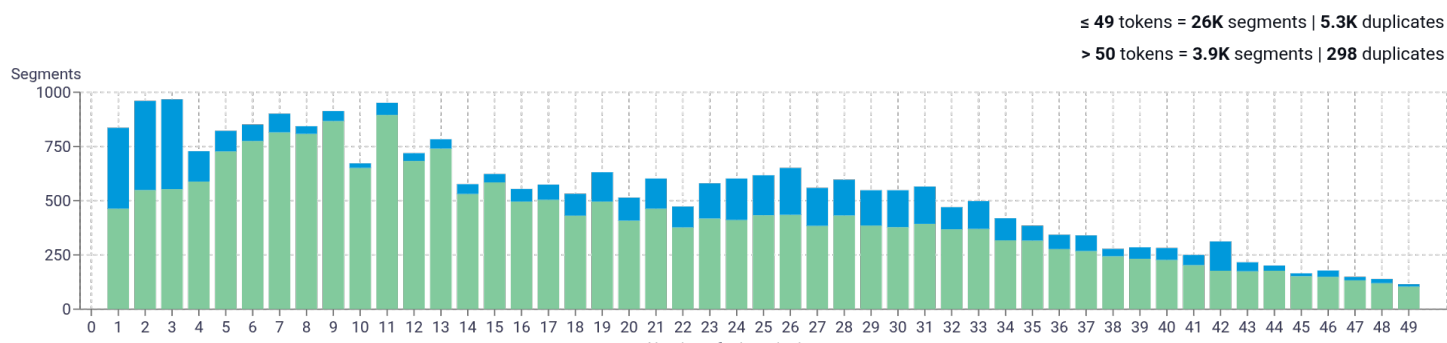
### Percentage of segments in Banjar inside documents

segments < 50% - **8.96%** (117 documents)
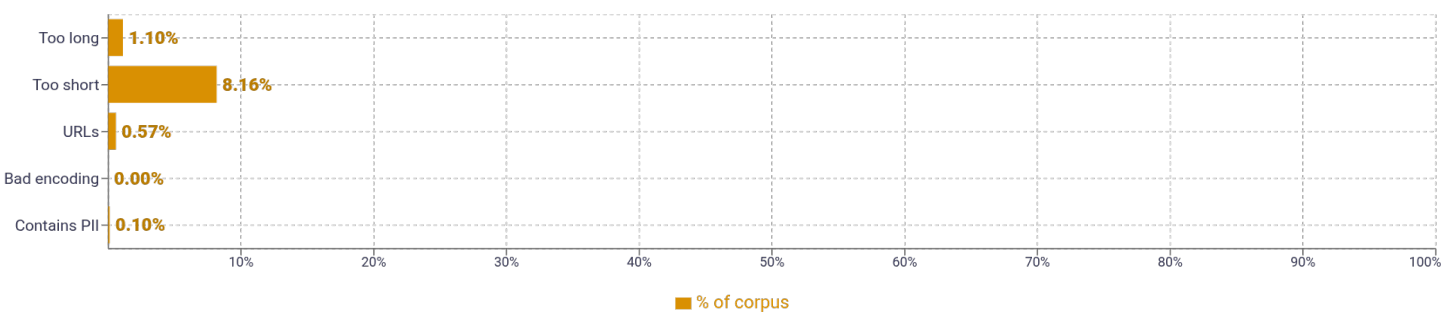segments ≥ 50% - **91.04%** (1.2K documents)

## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score ≥ 5 - **100%** (1.3K documents)

Documents



## Segment length distribution by token

**≤ 49** tokens = **26K** segments | **5.3K** duplicates
**> 50** tokens = **3.9K** segments | **298** duplicates

Segments



Number of tokens in the segment

## Segment noise distribution



- Too long — **1.10%**
- Too short — **8.16%**
- URLs — **0.57%**
- Bad encoding — **0.00%**
- Contains PII — **0.10%**

■ % of corpus

## Frequent n-grams

| SIZE | N-GRAMS |
|------|---------|
| 1 | يغ \| 16,988  ‏م \| 8,255  ‏ن \| 7,869  ‏و \| 7,212  ‏الله \| 6,597 |
| 2 | arabic script \| 4,102  yue chinese \| 998  ‏رسول الله \| 644  ‏هاري ابن \| 590  ‏اورڠ يغ \| 564 |
| 3 | ‏صلى الله عليه \| 464  ‏الله عليه وسلم \| 463  ‏الله سبحانه وتعالى \| 399  ‏الله صلى الله \| 183  ‏رسول الله صلى \| 182 |
| 4 | ‏صلى الله عليه وسلم \| 448  ‏رسول الله صلى الله \| 179  ‏الله صلى الله عليه \| 179  ‏فنديديقن اسلم تيغكاتن ساتو \| 123  ‏سوكاتن بارو کريکولوم برسفادو \| 123 |
| 5 | ‏رسول الله صلى الله عليه \| 175  ‏الله صلى الله عليه وسلم \| 170  ‏فنديديقن اسلم تيغكاتن \| 123  ‏رنخغن تاهونن فنديديقن اسلم تيغكاتن \| 123  ‏تاهونن فنديديقن اسلم تيغکاتن ساتو \| 123  ‏محمد صلى الله عليه وسلم \| 112 |

# 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 | | | | |