# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Analytics date | Language |
|---|---|---|
| bjn_Arab.jsonl.tsv | 12/5/2024 | Banjar (bjn) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Size | Characters |
|---|---|---|---|---|---|
| 1,112 | 19,529 | 11,025 (56.45 %) | 810K | 5.63 MB | 3,298,215 |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| utusanmelayu.com.my | 514 | 46.22 |
| wordpress.com | 283 | 25.45 |
| blogspot.com | 131 | 11.78 |
| blogspot.my | 18 | 1.62 |
| utusantv.com | 14 | 1.26 |
| wikisource.org | 13 | 1.17 |
| urusniaga.my | 9 | 0.81 |
| unicode.org | 8 | 0.72 |
| wikimedia.org | 8 | 0.72 |
| tronik.org | 7 | 0.63 |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com.my | 516 | 46.40 |
| com | 478 | 42.99 |
| org | 43 | 3.87 |
| my | 27 | 2.43 |
| net | 5 | 0.45 |
| sg | 5 | 0.45 |
| eu | 4 | 0.36 |
| web.id | 4 | 0.36 |
| org.my | 3 | 0.27 |
| moe | 3 | 0.27 |

## Documents size (in segments)

<= 25 segments **84.53%** (940 documents)
> 25 segments **15.47%** (172 documents)



## Documents by collection



cc18 (533)
cc22 (204)
cc21 (117)
16 Others (258)

## Language Distribution

### Number of segments



- Arabic (ar) - 16K
- English (en) - 700
- Indonesian (id) - 639
- Persian (fa) - 381
- Malay (ms) - 193
- Egyptian Arabic (arz) - 169
- Italian (it) - 143
- Nepali (ne) - 138
- South Azerbaijani (azb) - 124
- Mazanderani (mzn) - 113
- 41 Others - 503

*Banjar (bjn) identification might be inaccurate because language is not supported by Fasttext

### Percentage of segments in Banjar (bjn) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (1.1K documents)



## Segment length distribution by token

<= 49 tokens = **9.1K** segments | **7.2K** duplicates
> 50 tokens = **3.3K** segments | **1.3K** duplicates



- Unique segments
- Duplicated segments

## Segment noise distribution



| | |
|---|---|
| Too long | 1.15 % |
| Too short | 7.31 % |
| URLs | 0.73 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.03 % |

**Frequent n-grams**

| Size | n-grams |
|------|---------|
| 1 | يغ \| 16456    م \| 9389    و \| 8403    ن \| 8149    ى \| 7257 |
| 2 | فردان منتري \| 454    يغ تله \| 461    اق يغ \| 478    هاري اين \| 515    اورغ يغ \| 602 |
| 3 | الله سبحانه وتعالى \| 239    صلى الله عليه \| 289    الله عليه وسلم \| 290    يغ دفرتوان اكوغ \| 364    دركن دركن دركن \| 436 |
| 4 | ستياق اورغ عداله برحق \| 136    چت ماسيس چت ماسيس \| 154    ماسيس چت ماسيس چت \| 168    صلى الله عليه وسلم \| 283    دركن دركن دركن دركن \| 434 |
| 5 | رسول الله صلى الله عليه \| 98    ستياق اورغ عداله برحق كڤد \| 128    ماسيس چت ماسيس چت ماسيس \| 154    چت ماسيس چت ماسيس چت \| 154    دركن دركن دركن دركن دركن \| 432 |