# HPLT Analytics report

## General overview

| Corpus | Date | Language |
|---|---|---|
| hplt-v3-taq_Latn | 9/18/2025 | Tamasheq |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 827 | 48,115 | 43,515 (90.44 %) | 2.3M | 8,927,633 | 9.65 MB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| bible.is | 519 | 62.76% |
| biblehub.com | 141 | 17.05% |
| newchristianbib... | 79 | 9.55% |
| ebible.org | 39 | 4.72% |
| baebol.org | 6 | 0.73% |
| stepbible.org | 5 | 0.60% |
| gospelgo.com | 4 | 0.48% |
| bible.com | 4 | 0.48% |
| stalk.info | 3 | 0.36% |
| omniglot.com | 3 | 0.36% |

### Top 10 TLDs

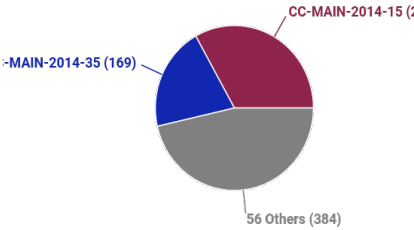| Domain | Docs | % of total |
|---|---|---|
| is | 519 | 62.76% |
| com | 162 | 19.59% |
| org | 135 | 16.32% |
| net | 4 | 0.48% |
| info | 3 | 0.36% |
| gq | 1 | 0.12% |
| es | 1 | 0.12% |
| edu.vn | 1 | 0.12% |
| ca | 1 | 0.12% |

## Documents size (in segments) ⓘ

≤ **25** segments **73.88%** (611 documents)
> **25** segments **26.12%** (216 documents)
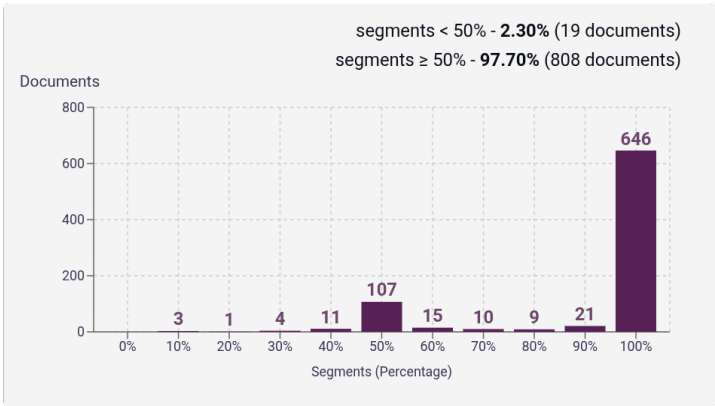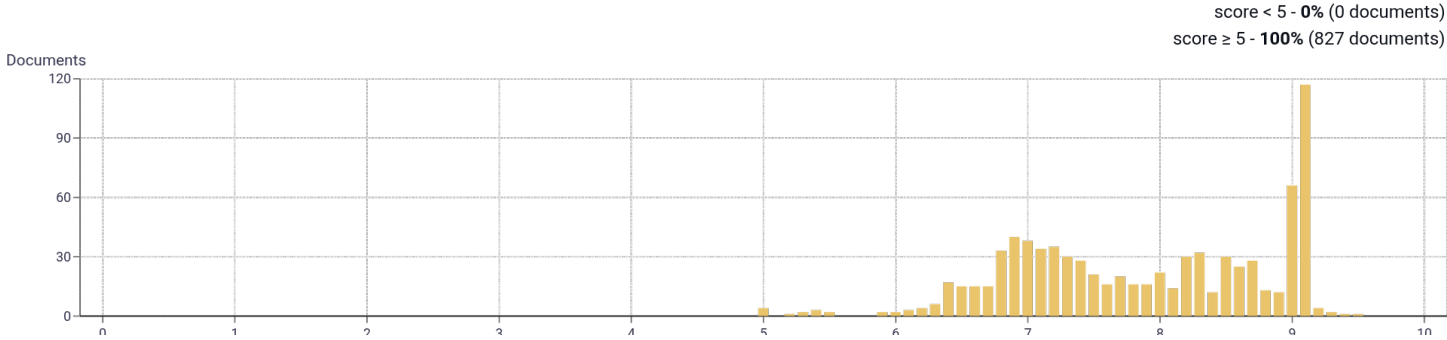


## Document collections

CC = 98.43%
IA = 1.57%



CC-MAIN-2014-15 (?
-MAIN-2014-35 (169)
56 Others (384)

## Language Distribution

### Number of segments in the Tamasheq corpus



- Azerbaijani - 12K **(25.6%)**
- English - 9.1K **(19.0%)**
- Romanian - 7.2K **(15.0%)**
- Swahili - 2K **(4.1%)**
- French - 1.6K **(3.3%)**
- German - 1.5K **(3.1%)**
- Upper Sorbian - 1.4K **(3.0%)**
- Indonesian - 1.3K **(2.6%)**
- Norwegian Bokmål - 1.2K **(2.6%)**
- Luxembourgish - 987 **(2.1%)**
- 127 Others - 9.5K **(19.7%)**

*Tamasheq identification might be inaccurate because it is not supported by FastSpell.

### Percentage of segments in Tamasheq inside documents

segments < 50% - **2.30%** (19 documents)
segments ≥ 50% - **97.70%** (808 documents)

## Distribution of documents by document score

Documents



## Segment length distribution by token

≤ 49 tokens = **40K** segments | **4.6K** duplicates
> 50 tokens = **8.6K** segments | **8** duplicates

Segments



## Segment noise distribution



| | |
|---|---|
| Too long | **1.58%** |
| Too short | **8.13%** |
| URLs | **0.05%** |
| Bad encoding | **0.07%** |
| Contains PII | **0.00%** |

■ % of corpus

## Frequent n-grams

| SIZE | N-GRAMS | | | | |
|---|---|---|---|---|---|
| 1 | n \| 38,547 | s \| 28,992 | dăɣ \| 26,710 | net \| 18,871 | ahay \| 18,049 |
| 2 | ata awan \| 3,617 | ɗo ahay \| 3,181 | ata nà \| 2,685 | ɗo sə \| 2,573 | anà ɗo \| 2,422 |
| 3 | mer su way \| 1,719 | anà ɗo ahay \| 1,026 | wulen su doh \| 657 | dăɣ a făl \| 654 | à man ata \| 627 |
| 4 | ga mer su way \| 432 | ɗo a yesu ahay \| 402 | à wulen su doh \| 378 | gəɗan dungo anà way \| 368 | sə gəɗan dungo anà \| 366 |
| 5 | sə gəɗan dungo anà way \| 364 | doh sə mazlab a mbərom \| 291 | ɗo si mer su way \| 279 | sa ga mer su way \| 239 | gəɗan dungo anà way ahay \| 239 |

# 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 | | | | |