# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|--------|------|----------|
| amh_Ethi.jsonl.tsv | 9/6/2024 | Amharic (am) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|------|----------|-----------------|--------|------------|------|
| 295,542 | 7,005,835 | 3,859,756 (55.09 %) | 226M | 1,024,632,898 | 2.39 GB |

## Top 10 domains

| Domain | Docs | % of total |
|--------|------|------------|
| ethiopianreport... | 15K | 4.92% |
| wordpress.com | 14K | 4.88% |
| voanews.com | 12K | 4.15% |
| addisadmassnews... | 11K | 3.84% |
| ethioreference.com | 8.8K | 2.99% |
| blogspot.com | 7.1K | 2.40% |
| blogspot.no | 7K | 2.35% |
| goolgule.com | 6.2K | 2.10% |
| wikipedia.org | 5.9K | 1.98% |
| dw.com | 5.6K | 1.89% |

## Top 10 TLDs

| Domain | Docs | % of total |
|--------|------|------------|
| com | 203K | 68.65% |
| org | 33K | 11.11% |
| net | 8.5K | 2.86% |
| no | 7.4K | 2.49% |
| gov.et | 5.3K | 1.81% |
| et | 5.3K | 1.79% |
| news | 2.6K | 0.87% |
| info | 2.4K | 0.80% |
| ch | 2.2K | 0.73% |
| us | 1.9K | 0.65% |

## Register labels



- HI - 0.3%
- ID - 0.8%
- IN - 8.2%
- IP - 5.5%
- LY - 0.1%
- MIX - 2.2%
- NA - 35.6%
- OP - 18.7%
- SP - 1.4%
- UNK - 27.1%

**MT**:17.9% | 53K Documents



- HI_other - 0.3%
- HI_re - 0.0%
- ID_other - 0.8%
- IN_dtp - 1.8%
- IN_en - 1.5%
- IN_fi - 0.0%
- IN_lt - 1.2%
- IN_other - 3.7%
- IN_ra - 0.0%
- IP_ds - 0.4%
- IP_ed - 0.0%
- IP_other - 5.1%
- LY_other - 0.1%
- MIX - 2.2%
- NA_nb - 0.7%
- NA_ne - 29.7%
- NA_other - 3.1%
- NA_sr - 2.1%
- OP_av - 0.3%
- OP_ob - 3.5%
- OP_other - 5.0%
- OP_rs - 9.6%
- OP_rv - 0.2%
- SP_it - 0.8%
- SP_other - 0.6%
- UNK - 27.1%

## Documents size (in segments)

<= 25 segments **76.14%** (225K documents)
> 25 segments **23.86%** (71K documents)



## Documents by collection

CC = 70.50%
IA = 29.50%



- cc22 (101K)
- cc18 (49K)
- cc21 (38K)
- 18 Others (108K)

## Language Distribution

### Number of segments in the Amharic (am) corpus



- Amharic (am) - 6.1M
- English (en) - 551K
- Italian (it) - 91K
- Korean (ko) - 31K
- French (fr) - 23K
- Russian (ru) - 15K
- German (de) - 14K
- Chinese (zh) - 11K
- Greek (el) - 11K
- Spanish (es) - 9.8K
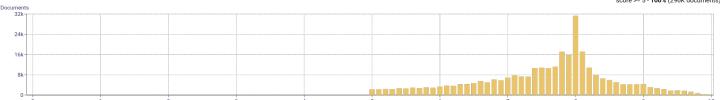- 155 Others - 139K

### Percentage of segments in Amharic (am) inside documents



## Distribution of documents by document score

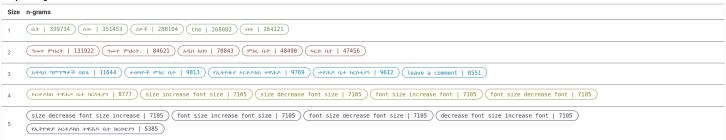score < 5 - **0%** (0 documents)
score >= 5 - **100%** (296K documents)

## Segment length distribution by token

■ Unique segments  ■ Duplicated segments

## Segment noise distribution



| Noise type | Percentage |
|---|---|
| Too long | 0.61% |
| Too short | 10.07% |
| URLs | 0.78% |
| Bad encoding | 0.01% |
| Contains PII | 0.15% |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | ቤት \| 399734   ሰው \| 351453   ሰዎች \| 288104   the \| 268082   ብዙ \| 264121 |
| 2 | ዓመተ ምህረት \| 131922   ዓመተ ምህረት. \| 84621   አዲስ አበባ \| 70843   ምክር ቤት \| 48490   ፍርድ ቤት \| 47456 |
| 3 | አዲስ ግምግማዎች በይፋ \| 11644   ተወካዮች ምክር ቤት \| 9813   የኢትዮጵያ ኦርቶዶክስ ተዋሕዶ \| 9769   ተዋሕዶ ቤተ ክርስቲያን \| 9612   leave a comment \| 8551 |
| 4 | ኦርቶዶክስ ተዋሕዶ ቤተ ክርስቲያን \| 8777   size increase font size \| 7105   size decrease font size \| 7105   font size increase font \| 7105   font size decrease font \| 7105 |
| 5 | size decrease font size increase \| 7105   font size increase font size \| 7105   font size decrease font size \| 7105   decrease font size increase font \| 7105   የኢትዮጵያ ኦርቶዶክስ ተዋሕዶ ቤተ ክርስቲያን \| 5385 |

## 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 | | | | |