# HPLT Analytics report

**HPLT** Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| hplt-v3-tam_Taml | 9/18/2025 | Tamil |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 11,266,205 | 205,078,222 | 130,153,460 (63.47 %) | 4.4B | 32,535,155,530 | 80.85 GB |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| blogspot.com | 546K | 4.85% |
| dinakaran.com | 350K | 3.11% |
| dinamani.com | 319K | 2.83% |
| dinamalar.com | 309K | 2.74% |
| vikatan.com | 239K | 2.12% |
| hindutamil.in | 197K | 1.75% |
| maalaimalar.com | 159K | 1.41% |
| dailythanthi.com | 152K | 1.35% |
| asianetnews.com | 143K | 1.27% |
| news18.com | 138K | 1.22% |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 8.4M | 74.88% |
| in | 1M | 9.07% |
| org | 371K | 3.29% |
| net | 351K | 3.12% |
| lk | 332K | 2.94% |
| ca | 61K | 0.55% |
| news | 51K | 0.45% |
| com.my | 47K | 0.42% |
| tv | 46K | 0.41% |
| ch | 41K | 0.37% |

## Register labels



Pie chart:
- HI - 1.5%
- ID - 0.7%
- IN - 8.6%
- IP - 1.6%
- LY - 0.2%
- MIX - 1.0%
- NA - 65.8%
- OP - 9.0%
- SP - 0.5%
- UNK - 11.2%

Documents bar chart legend:
- HI_other - 0.7%
- HI_re - 0.8%
- ID_other - 0.7%
- IN_dtp - 2.7%
- IN_en - 0.9%
- IN_fi - 0.0%
- IN_lt - 0.0%
- IN_other - 4.9%
- IN_ra - 0.0%
- IP_ds - 0.9%
- IP_other - 0.7%
- LY_other - 0.2%
- MIX - 1.0%
- NA_nb - 1.6%
- NA_ne - 57.5%
- NA_other - 4.0%
- NA_sr - 2.8%
- OP_av - 0.4%
- OP_ob - 1.9%
- OP_other - 2.5%
- OP_rs - 2.8%
- OP_rv - 1.4%
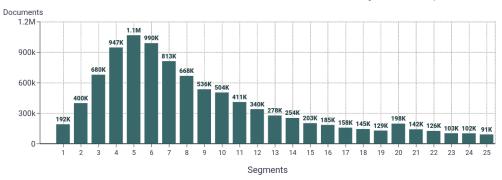- SP_it - 0.2%
- SP_other - 0.3%
- UNK - 11.2%

**MT**:3.8% | 425K Documents

## Documents size (in segments) ⓘ

≤ 25 segments **85.77%** (9.7M documents)
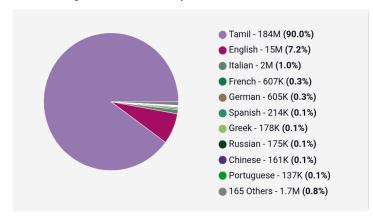> 25 segments **14.23%** (1.6M documents)



Documents by segments: 1: 192K, 2: 400K, 3: 680K, 4: 947K, 5: 1.1M, 6: 990K, 7: 813K, 8: 668K, 9: 536K, 10: 504K, 11: 411K, 12: 340K, 13: 278K, 14: 254K, 15: 203K, 16: 185K, 17: 158K, 18: 145K, 19: 129K, 20: 198K, 21: 142K, 22: 126K, 23: 103K, 24: 102K, 25: 91K

## Document collections

**CC = 91.26%**
**IA = 8.74%**



67 Others (11M)

## Language Distribution

### Number of segments in the Tamil corpus



- Tamil - 184M **(90.0%)**
- English - 15M **(7.2%)**
- Italian - 2M **(1.0%)**
- French - 607K **(0.3%)**
- German - 605K **(0.3%)**
- Spanish - 214K **(0.1%)**
- Greek - 178K **(0.1%)**
- Russian - 175K **(0.1%)**
- Chinese - 161K **(0.1%)**
- Portuguese - 137K **(0.1%)**
- 165 Others - 1.7M **(0.8%)**

### Percentage of segments in Tamil inside documents

segments < 50% - **1.80%** (203K documents)
segments ≥ 50% - **98.20%** (11M documents)



### Distribution of documents by document score

score < 5 - **0%** (0 documents)
score ≥ 5 - **100%** (11M documents)



### Segment length distribution by token

**≤ 49** tokens = **185M** segments | **71M** duplicates
**> 50** tokens = **20M** segments | **4.1M** duplicates



### Segment noise distribution



- Too long — 0.56%
- Too short — 13.26%
- URLs — 0.85%
- Bad encoding — 0.00%
- Contains PII — 0.10%

■ % of corpus

## Frequent n-grams

| SIZE | N-GRAMS | |
|---|---|---|
| 1 | ஆனால் \| 7,592,065   செய்து \| 5,283,014   இல்லை \| 5,203,824   உங்கள் \| 5,044,951   அரசு \| 4,872,464 | ⧉ |
| 2 | read more \| 804,105   a comment \| 799,956   post a \| 683,261   மத்திய அரசு \| 678,336   ஆம் ஆண்டு \| 581,194 | ⧉ |
| 3 | post a comment \| 681,856   by ayyasamy ram \| 190,947   in your inbox \| 163,009   to receive our \| 162,815   up to receive \| 162,751 | ⧉ |
| 4 | in your inbox every \| 162,738   up to receive our \| 162,736   your inbox every day \| 162,735   to receive our newsletter \| 162,735   receive our newsletter in \| 162,735 | ⧉ |
| 5 | up to receive our newsletter \| 162,735   to receive our newsletter in \| 162,735   receive our newsletter in your \| 162,735   our newsletter in your inbox \| 162,735   newsletter in your inbox every \| 162,735 | ⧉ |

# 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 | | | | |