# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| hplt-v3-heb_Hebr | 9/18/2025 | Hebrew |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 26,082,588 | 647,272,554 | 410,742,381 (63.46 %) | 16B | 78,486,012,143 | 129.13 GB |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| walla.co.il | 931K | 3.57% |
| ynet.co.il | 387K | 1.48% |
| wikipedia.org | 347K | 1.33% |
| haaretz.co.il | 256K | 0.98% |
| maariv.co.il | 250K | 0.96% |
| blogspot.com | 228K | 0.88% |
| mako.co.il | 219K | 0.84% |
| wordpress.com | 218K | 0.84% |
| psakdin.co.il | 194K | 0.74% |
| calcalist.co.il | 180K | 0.69% |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| co.il | 15M | 58.30% |
| com | 5.4M | 20.84% |
| org | 1.6M | 6.10% |
| org.il | 1.4M | 5.50% |
| net | 806K | 3.09% |
| ac.il | 346K | 1.32% |
| tv | 145K | 0.56% |
| info | 118K | 0.45% |
| muni.il | 63K | 0.24% |
| gov.il | 52K | 0.20% |

## Register labels



Legend (pie):
- HI - 3.4%
- ID - 3.3%
- IN - 14.4%
- IP - 26.2%
- LY - 0.1%
- MIX - 4.6%
- NA - 26.7%
- OP - 11.3%
- SP - 0.5%
- UNK - 9.4%

Legend (bar chart):
- HI_other - 1.8%
- HI_re - 1.5%
- ID_other - 3.3%
- IN_dtp - 4.9%
- IN_en - 1.8%
- IN_fi - 0.0%
- IN_lt - 1.3%
- IN_other - 6.2%
- IN_ra - 0.1%
- IP_ds - 22.9%
- IP_other - 3.3%
- LY_other - 0.1%
- MIX - 4.6%
- NA_nb - 4.1%
- NA_ne - 16.9%
- NA_other - 3.1%
- NA_sr - 2.6%
- OP_av - 2.0%
- OP_ob - 1.9%
- OP_other - 2.0%
- OP_rs - 2.8%
- OP_rv - 2.6%
- SP_it - 0.4%
- SP_other - 0.2%
- UNK - 9.4%

**MT**:6.5% | 1.7M Documents

## Documents size (in segments) ⓘ

≤ 25 segments **75.43%** (20M documents)
> 25 segments **24.57%** (6.4M documents)



## Document collections

CC = 88.07%
IA = 11.93%

67 Others (26M)

## Language Distribution

### Number of segments in the Hebrew corpus

- Hebrew - 611M **(94.5%)**
- English - 14M **(2.2%)**
- Italian - 9.6M **(1.5%)**
- French - 2.5M **(0.4%)**
- Yiddish - 1.6M **(0.2%)**
- German - 1.1M **(0.2%)**
- Spanish - 974K **(0.2%)**
- Russian - 499K **(0.1%)**
- Portuguese - 477K **(0.1%)**
- Greek - 433K **(0.1%)**
- 165 Others - 4.6M **(0.7%)**

### Percentage of segments in Hebrew inside documents

segments < 50% - **0.46%** (119K documents)
segments ≥ 50% - **99.54%** (26M documents)

Documents

20M

15M

10M

5M

0

246  8K  21K  30K  60K  259K  358K  642K  1.3M  3.4M  20M

0%  10%  20%  30%  40%  50%  60%  70%  80%  90%  100%

Segments (Percentage)

## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score ≥ 5 - **100%** (26M documents)

Documents

1.8M

1.35M

900k

450k

0

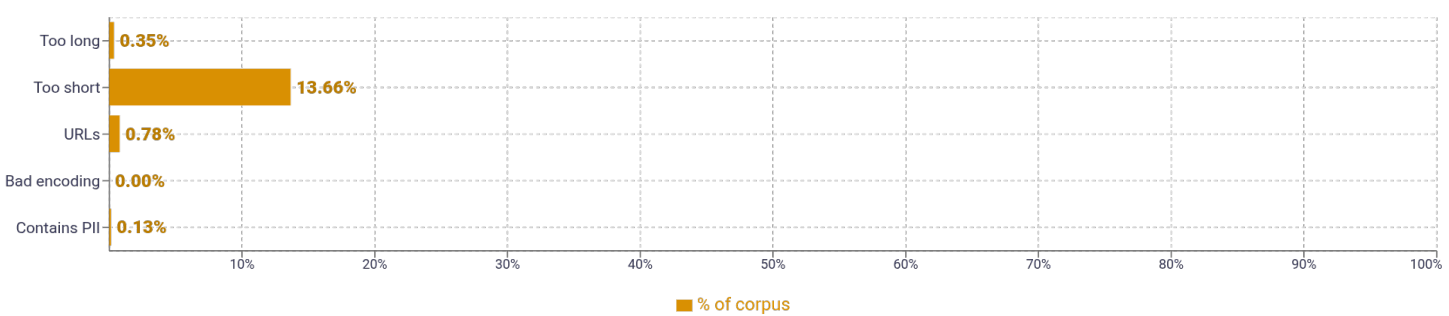0    1    2    3    4    5    6    7    8    9    10

## Segment length distribution by token

≤ 49 tokens = **551M** segments | **223M** duplicates
> 50 tokens = **97M** segments | **14M** duplicates

Segments

40M

30M

20M

10M

0

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49

## Segment noise distribution

| | % of corpus |
|---|---|
| Too long | 0.35% |
| Too short | 13.66% |
| URLs | 0.78% |
| Bad encoding | 0.00% |
| Contains PII | 0.13% |

10%  20%  30%  40%  50%  60%  70%  80%  90%  100%

■ % of corpus

## Frequent n-grams

| SIZE | N-GRAMS |
|---|---|
| 1 | כדי \| 26,307,297   ב \| 25,200,621   אחד \| 22,614,309   ה \| 20,435,743   בכל \| 18,650,872 |
| 2 | תל אביב \| 2,946,106   בית המשפט \| 2,509,445   לאחר מכן \| 2,056,397   בתל אביב \| 1,679,738   חוות דעת \| 1,645,229 |
| 3 | בסופו של דבר \| 1,433,400   עריכת קוד מקור \| 1,249,605   בתי מלון ב \| 399,233   בית המשפט העליון \| 372,262   תגובתך נקלטה בהצלחה \| 316,527 |
| 4 | ותפורסם על פי מדיניות \| 316,278   ראש הממשלה בנימין נתניהו \| 162,066   נסה שנית במועד מאוחר \| 160,738   אנא נסה שנית במועד \| 160,738   התראה בדוא"ל כאשר תגובתך \| 159,206 |
| 5 | ותפורסם על פי מדיניות המערכת \| 316,278   אנא נסה שנית במועד מאוחר \| 160,737   לקבל התראה בדוא"ל כאשר תגובתך \| 159,206   התראה בדוא"ל כאשר תגובתך תאושר \| 159,206   בדוא"ל כאשר תגובתך תאושר ותפורסם \| 159,206 |

# 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 | | | | |