# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| ydd_Hebr.jsonl.tsv | 12/5/2024 | Yiddish (ydd) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 128,265 | 2,940,163 | 1,420,745 (48.32 %) | 89M | 455,684,296 | 774.19 MB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| wikipedia.org | 45K | 34.77% |
| kaveshtiebel.com | 15K | 11.49% |
| ivelt.com | 10K | 8.06% |
| yiddish.news | 4.1K | 3.22% |
| soft-free-downl... | 3.7K | 2.87% |
| sgames.org | 2.5K | 1.93% |
| freeplayonlineg... | 1.8K | 1.41% |
| actualidadgadge... | 1.4K | 1.12% |
| itsmygame.org | 1.2K | 0.95% |
| creativosonline... | 1.2K | 0.92% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 58K | 45.11% |
| org | 54K | 42.22% |
| news | 4.3K | 3.34% |
| net | 3.7K | 2.91% |
| ru | 1.1K | 0.84% |
| co.il | 775 | 0.60% |
| zone | 667 | 0.52% |
| gov | 609 | 0.47% |
| de | 549 | 0.43% |
| co | 513 | 0.40% |

## Register labels



- HI - 0.4%
- ID - 17.9%
- IN - 37.0%
- IP - 1.8%
- LY - 0.3%
- MIX - 0.6%
- NA - 5.7%
- OP - 4.7%
- SP - 0.1%
- UNK - 31.6%

🖐 **MT**:30.0% | 38K Documents

Documents

- HI_other - 0.3%
- HI_re - 0.0%
- ID_other - 17.9%
- IN_dtp - 1.7%
- IN_en - 32.0%
- IN_fi - 0.0%
- IN_lt - 0.2%
- IN_other - 3.2%
- IN_ra - 0.0%
- IP_ds - 1.5%
- IP_ed - 0.0%
- IP_other - 0.3%
- LY_other - 0.3%
- MIX - 0.6%
- NA_nb - 0.8%
- NA_ne - 3.9%
- NA_other - 1.1%
- NA_sr - 0.0%
- OP_av - 0.1%
- OP_ob - 0.3%
- OP_other - 1.2%
- OP_rs - 2.8%
- OP_rv - 0.3%
- SP_it - 0.0%
- SP_other - 0.0%
- UNK - 31.6%

## Documents size (in segments)

<= 25 segments **79.98%** (103K documents)
> 25 segments **20.02%** (26K documents)



## Documents by collection

CC = 49.18%
IA = 50.82%



- cc18 (15K)
- wide11 (15K)
- wide16 (13K)
- cc22 (21K)
- 17 Others (64K)

## Language Distribution

### Number of segments in the Yiddish (ydd) corpus



- Yiddish (yi) - 2.4M
- Hebrew (he) - 441K
- English (en) - 64K
- Italian (it) - 38K
- French (fr) - 5.4K
- Somali (so) - 4.6K
- German (de) - 3.9K
- Spanish (es) - 2.9K
- Greek (el) - 2.6K
- Russian (ru) - 2.5K
- 129 Others - 22K

*Yiddish (ydd) identification might be inaccurate because it is not supported by FastSpell

### Percentage of segments in Yiddish (ydd) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (128K documents)

## Segment length distribution by token

≤ **49** tokens = **1.1M** segments | **1.3M** duplicates
> **50** tokens = **579K** segments | **223K** duplicates



Segments

- Unique segments
- Duplicated segments

Number of tokens in the segment

## Segment noise distribution



| | |
|---|---|
| Too long | 0.94 % |
| Too short | 10.66 % |
| URLs | 0.90 % |
| Bad encoding | 0.00 % |
| Contains PII | 0.12 % |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | 1847266 \| פון    1600808 \| איז    1355274 \| צו    853755 \| האט    680390 \| נישט |
| 2 | 206425 \| עס איז    129219 \| איז געווען    127409 \| האט געשריבן    96478 \| ער האט    72789 \| וואס איז |
| 3 | 63357 \| זיך איינגעשריבן אום    15325 \| עס איז נישט    15249 \| עס איז א    15004 \| עס איז געווען    14271 \| אז עס איז |
| 4 | 6270 \| באנוצערס וואס דרייען זיך    6176 \| וואס לייענען דעם פארום    5839 \| וואס דרייען זיך דא    5468 \| איינער פון די מערסט    5110 \| נישטא קיין אנליין באנוצער |
| 5 | 6176 \| באנוצער וואס לייענען דעם פארום    5814 \| באנוצערס וואס דרייען זיך דא    2529 \| טן יאנואר ליטן יוליאנישן קאלענדאר    2527 \| ליטן יוליאנישן קאלענדאר מיט    2527 \| יאנואר ליטן גרעגאריאנישן    2527 \| דאטום דא זענען ליטן |

## 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 | | | | |