# HPLT Analytics report

## General overview

| Corpus | Analytics date | Language |
|---|---|---|
| HPLT-v2-nld_Latn.tsv | 9/28/2024 | Dutch (nl) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Size | Characters |
|---|---|---|---|---|---|
| 138,651,084 | 3,074,576,227 | | | 419.09 GB | 448,140,483,636 |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| blogspot.com | 1.9M | 1.37 |
| wikipedia.org | 1.7M | 1.26 |
| blogspot.nl | 1.1M | 0.80 |
| knack.be | 982K | 0.71 |
| wordpress.com | 895K | 0.65 |
| docplayer.nl | 697K | 0.50 |
| nrc.nl | 535K | 0.39 |
| blogspot.be | 488K | 0.35 |
| tripadvisor.nl | 369K | 0.27 |
| viva.nl | 344K | 0.25 |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| nl | 91M | 65.64 |
| com | 19M | 13.37 |
| be | 17M | 12.12 |
| org | 3.4M | 2.48 |
| net | 1.9M | 1.38 |
| eu | 1.9M | 1.35 |
| nu | 1.1M | 0.80 |
| info | 781K | 0.56 |
| de | 298K | 0.21 |
| tv | 172K | 0.12 |

## Documents size (in segments)

<= 25 segments **78.98%** (110M documents)
> 25 segments **21.02%** (29M documents)



## Documents by collection



cc22 (43M), cc18 (24M), cc21 (18M), 18 Others (54M)

## Language Distribution

### Number of segments



- Dutch (nl) - 2.5B
- English (en) - 229M
- Italian (it) - 86M
- German (de) - 57M
- French (fr) - 55M
- Spanish (es) - 17M
- Swedish (sv) - 12M
- Norwegian Bokmål (nb) - 7.1M
- Polish (pl) - 7M
- Afrikaans (af) - 5.7M
- 165 Others - 72M

### Percentage of segments in Dutch (nl) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (139M documents)



## Segment noise distribution



- Too long: 0.00 %
- Too short: 7.85 %
- URLs: 2.26 %
- Bad encoding: 0.01 %
- Contains PII: 1.06 %

## 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