# HPLT Analytics report

## General overview

| Corpus | Date | Language |
|---|---|---|
| dan_Latn.jsonl.tsv | 6/16/2025 | Danish (da) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 33,841,355 | 872,886,142 | 340,944,296 (39.06 %) | 24B | 132,541,320,753 | 126.49 GB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| docplayer.dk | 639K | 1.89% |
| blogspot.com | 541K | 1.60% |
| wikipedia.org | 458K | 1.35% |
| billedeverden.com | 444K | 1.31% |
| blogspot.dk | 271K | 0.80% |
| tripadvisor.dk | 239K | 0.71% |
| dagens.dk | 234K | 0.69% |
| avisen.dk | 155K | 0.46% |
| wordpress.com | 147K | 0.43% |
| ekstrabladet.dk | 144K | 0.43% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| dk | 26M | 76.44% |
| com | 4.8M | 14.11% |
| org | 771K | 2.28% |
| eu | 644K | 1.90% |
| net | 350K | 1.03% |
| nu | 210K | 0.62% |
| info | 179K | 0.53% |
| no | 105K | 0.31% |
| se | 92K | 0.27% |
| de | 71K | 0.21% |

## Register labels



- HI - 3.8%
- ID - 1.5%
- IN - 14.0%
- IP - 33.5%
- LY - 0.0%
- MIX - 6.8%
- NA - 25.7%
- OP - 6.9%
- SP - 0.2%
- UNK - 7.5%

MT:6.2% | 2.1M Documents



- HI_other - 2.3%
- HI_re - 1.5%
- ID_other - 1.5%
- IN_dtp - 5.2%
- IN_en - 1.7%
- IN_fi - 0.1%
- IN_lt - 1.2%
- IN_other - 5.8%
- IN_ra - 0.1%
- IP_ds - 30.6%
- IP_ed - 0.0%
- IP_other - 2.9%
- LY_other - 0.0%
- MIX - 6.8%
- NA_nb - 8.9%
- NA_ne - 11.7%
- NA_other - 3.1%
- NA_sr - 2.1%
- OP_av - 1.3%
- OP_ob - 1.4%
- OP_other - 1.0%
- OP_rs - 0.3%
- OP_rv - 2.8%
- SP_it - 0.2%
- SP_other - 0.1%
- UNK - 7.5%

## Documents size (in segments)

<= 25 segments **76.79%** (26M documents)
> 25 segments **23.21%** (7.9M documents)



## Documents by collection

CC = 69.09%
IA = 30.91%



- cc18 (6.4M)
- cc22 (9.1M)
- cc21 (4M)
- 18 Others (14M)

## Language Distribution

### Number of segments in the Danish (da) corpus



- Danish (da) - 722M
- English (en) - 52M
- Italian (it) - 20M
- German (de) - 17M
- Norwegian Bokmål (nb) - 11M
- French (fr) - 8.5M
- Dutch (nl) - 6M
- Swedish (sv) - 5.2M
- Spanish (es) - 3.8M
- Polish (pl) - 3.8M
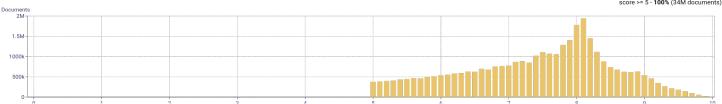- 165 Others - 24M

### Percentage of segments in Danish (da) inside documents



## Distribution of documents by document score

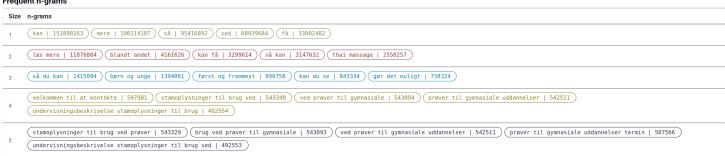score < 5 - **0%** (0 documents)
score >= 5 - **100%** (34M documents)

## Segment length distribution by token

Segments

Number of tokens in the segment

■ Unique segments  ■ Duplicated segments

## Segment noise distribution



| | |
|---|---|
| Too long | 1.03 % |
| Too short | 14.36 % |
| URLs | 2.99 % |
| Bad encoding | 0.01 % |
| Contains PII | 1.15 % |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | kan \| 151800263   mere \| 100114187   så \| 95416892   ved \| 68939684   få \| 33002482 |
| 2 | læs mere \| 11876884   blandt andet \| 4161626   kan få \| 3299614   så kan \| 3147632   thai massage \| 2550257 |
| 3 | så du kan \| 1415094   børn og unge \| 1394061   først og fremmest \| 999758   kan du se \| 843334   gør det muligt \| 738324 |
| 4 | velkommen til at kontakte \| 597981   stamoplysninger til brug ved \| 543348   ved prøver til gymnasiale \| 543094   prøver til gymnasiale uddannelser \| 542511   undervisningsbeskrivelse stamoplysninger til brug \| 492554 |
| 5 | stamoplysninger til brug ved prøver \| 543329   brug ved prøver til gymnasiale \| 543093   ved prøver til gymnasiale uddannelser \| 542511   prøver til gymnasiale uddannelser termin \| 507566   undervisningsbeskrivelse stamoplysninger til brug ved \| 492553 |

## 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 | | | | |