# HPLT Analytics report

**◉ HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| hplt-v3-ces_Latn | 9/18/2025 | Czech |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 107,802,248 | 2,467,616,974 | 1,166,958,008 (47.29 %) | 66B | 365,482,306,497 | 377.31 GB |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| denik.cz | 2.5M | 2.31% |
| idnes.cz | 1.6M | 1.46% |
| novinky.cz | 1.1M | 1.02% |
| web.app | 1M | 0.97% |
| blog.cz | 972K | 0.90% |
| docplayer.cz | 785K | 0.73% |
| ceskatelevize.cz | 647K | 0.60% |
| blogspot.com | 592K | 0.55% |
| rozhlas.cz | 538K | 0.50% |
| firebaseapp.com | 496K | 0.46% |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| cz | 90M | 83.16% |
| com | 8.3M | 7.74% |
| eu | 2.3M | 2.17% |
| org | 1.3M | 1.23% |
| net | 1.1M | 1.06% |
| app | 1.1M | 0.98% |
| sk | 912K | 0.85% |
| info | 683K | 0.63% |
| ru | 254K | 0.24% |
| tv | 133K | 0.12% |

## Register labels



Pie chart legend:
- HI - 3.0%
- ID - 3.4%
- IN - 11.4%
- IP - 37.1%
- LY - 0.1%
- MIX - 4.1%
- NA - 26.7%
- OP - 5.1%
- SP - 0.7%
- UNK - 8.4%

Bar chart legend:
- HI_other - 1.8%
- HI_re - 1.3%
- ID_other - 3.4%
- IN_dtp - 4.0%
- IN_en - 0.6%
- IN_fi - 0.0%
- IN_lt - 0.7%
- IN_other - 5.9%
- IN_ra - 0.0%
- IP_ds - 34.0%
- IP_other - 3.1%
- LY_other - 0.1%
- MIX - 4.1%
- NA_nb - 4.9%
- NA_ne - 14.7%
- NA_other - 3.5%
- NA_sr - 3.6%
- OP_av - 0.6%
- OP_ob - 0.9%
- OP_other - 0.9%
- OP_rs - 0.4%
- OP_rv - 2.2%
- SP_it - 0.5%
- SP_other - 0.2%
- UNK - 8.4%

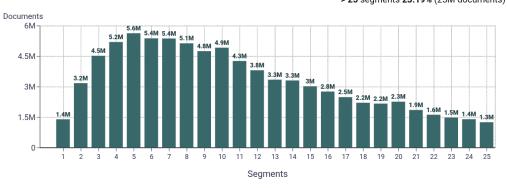🤖 **MT**:4.8% | 5.2M Documents

## Documents size (in segments) ⓘ

≤ 25 segments **76.81%** (83M documents)
> 25 segments **23.19%** (25M documents)



## Document collections

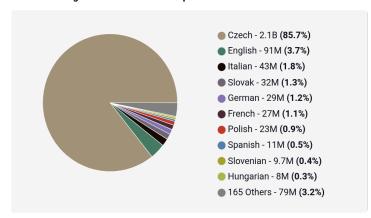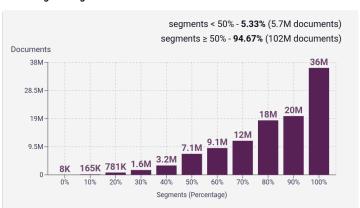CC = 90.67%
IA = 9.33%



67 Others (108M)

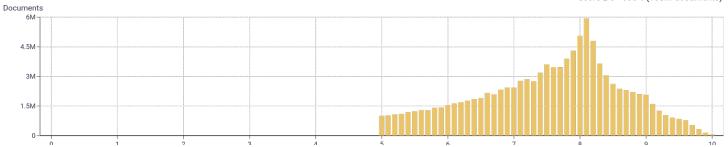## Language Distribution

### Number of segments in the Czech corpus

- Czech - 2.1B **(85.7%)**
- English - 91M **(3.7%)**
- Italian - 43M **(1.8%)**
- Slovak - 32M **(1.3%)**
- German - 29M **(1.2%)**
- French - 27M **(1.1%)**
- Polish - 23M **(0.9%)**
- Spanish - 11M **(0.5%)**
- Slovenian - 9.7M **(0.4%)**
- Hungarian - 8M **(0.3%)**
- 165 Others - 79M **(3.2%)**

### Percentage of segments in Czech inside documents

segments < 50% - **5.33%** (5.7M documents)
segments ≥ 50% - **94.67%** (102M documents)

Documents

| 0% | 10% | 20% | 30% | 40% | 50% | 60% | 70% | 80% | 90% | 100% |
|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|
| 8K | 165K | 781K | 1.6M | 3.2M | 7.1M | 9.1M | 12M | 18M | 20M | 36M |

Segments (Percentage)

### Distribution of documents by document score

score < 5 - **0%** (0 documents)
score ≥ 5 - **100%** (108M documents)



### Segment length distribution by token

≤ **49** tokens = **2.1B** segments | **1.2B** duplicates

> **50** tokens = **381M** segments | **108M** duplicates



### Segment noise distribution

| | % of corpus |
|---|---|
| Too long | 0.84% |
| Too short | 15.87% |
| URLs | 2.49% |
| Bad encoding | 0.01% |
| Contains PII | 0.66% |

■ % of corpus

## Frequent n-grams

| SIZE | N-GRAMS | |
|---|---|---|
| 1 | a \| 1,712,557,567  v \| 1,031,555,076  se \| 936,203,781  s \| 487,047,348  z \| 342,408,726 | ⧉ |
| 2 | v roce \| 29,393,125  že se \| 25,101,426  se v \| 22,491,588  a v \| 20,701,832  v případě \| 18,678,974 | ⧉ |
| 3 | jedná se o \| 7,888,261  v české republice \| 6,198,254  se jedná o \| 4,485,230  v současné době \| 4,324,382  v souladu s \| 3,638,486 | ⧉ |
| 4 | napíše příspěvek k této \| 1,792,064  příspěvek k této položce \| 1,792,019  vltava labe media a.s. \| 1,567,148  že se jedná o \| 1,081,330  zdarma a bez rizika \| 914,590 | ⧉ |
| 5 | napíše příspěvek k této položce \| 1,792,019  obsahu denik.cz je bez písemného \| 795,236  publikování nebo šíření obsahu denik.cz \| 795,235  denik.cz je bez písemného souhlasu \| 795,224  základní škola a mateřská škola \| 760,058 | ⧉ |

# 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 | | | | |