

## General overview

Corpus	Analytics date	Language
HPLT-docslite.fi.tsv	6/9/2024	Finnish (fi)

## Volumes

Docs	Segments	Unique segments	Tokens	Size
3,490,538	504,776,624	93,529 (0.02 %)	5.5B	34.89 GB

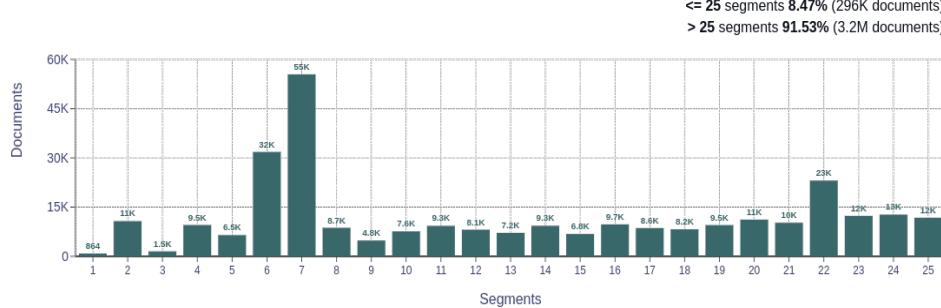
## Top 10 domains

Domain	Docs	% of total
blogspot.com	304K	8.71
docplayer.fi	121K	3.47
diebuchsuche.com	89K	2.54
letsbookhotel.com	59K	1.69
lightinthebox.com	48K	1.37
wordpress.com	33K	0.94
mininthebox.com	27K	0.78
lily.fi	22K	0.63
wikipedia.org	22K	0.62
uusisuomi.fi	20K	0.57

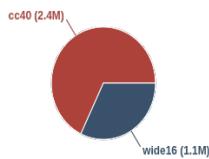
## Top 10 TLDs

Domain	Docs	% of total
fi	1.9M	53.28
com	1.2M	34.94
net	183K	5.24
org	73K	2.08
info	44K	1.26
eu	24K	0.70
blog	8.5K	0.24
pt	4.8K	0.14
ru	4.1K	0.12
ee	3.4K	0.10

## Documents size (in segments)

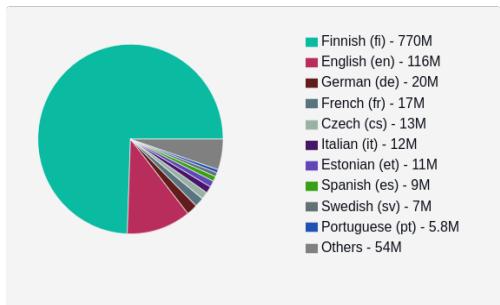


## Documents by collection

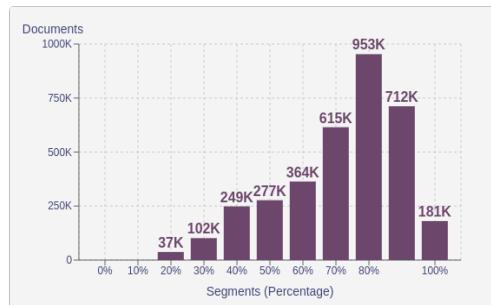


## Language Distribution

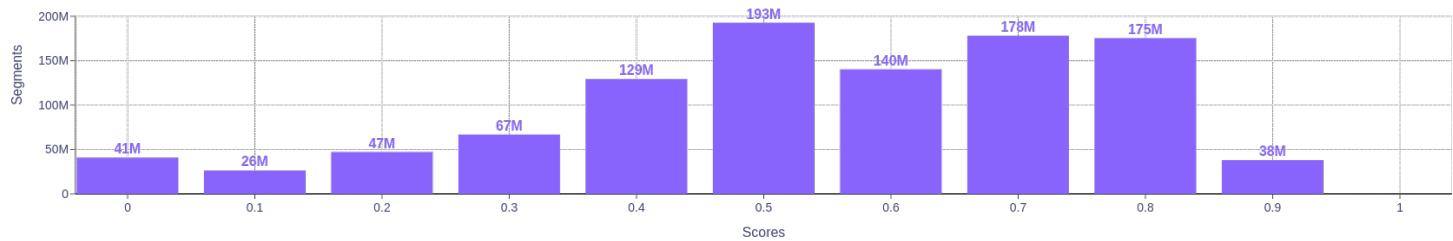
### Number of segments



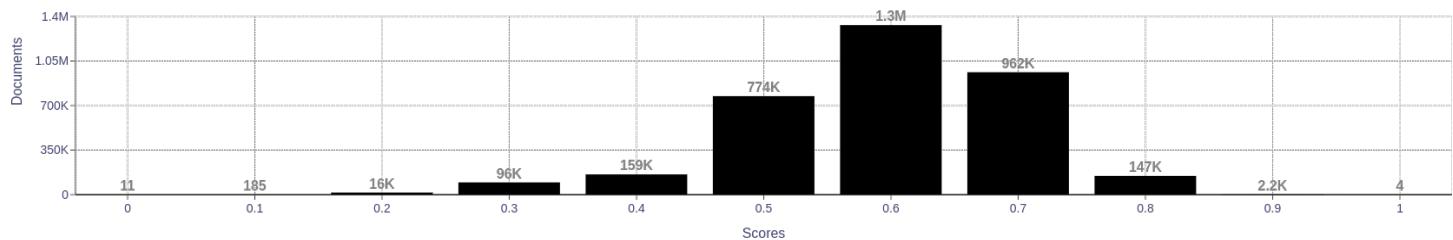
### Percentage of segments in Finnish (fi) inside documents



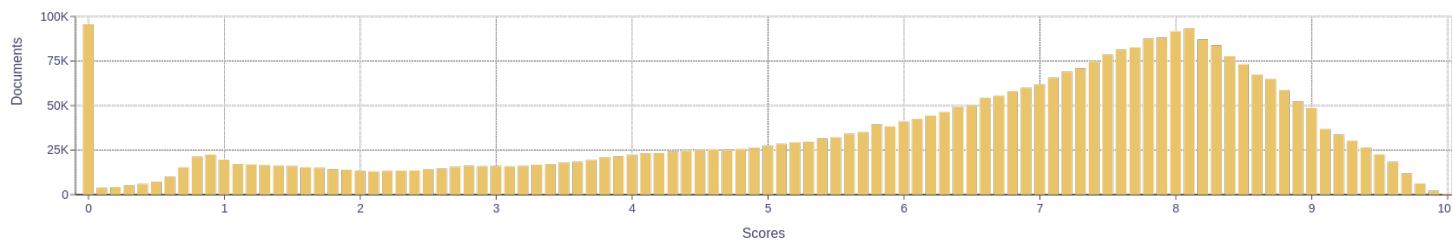
## Distribution of segments by fluency score



## Distribution of documents by average fluency score



## Distribution of documents by document score



## Segment length distribution by token

<= 49 tokens = 103M segments | 382M duplicates

> 50 tokens = 20M segments | 5.1M duplicates



## 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 of 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>