

## General overview

Corpus	Analytics date	Language
HPLT-docsite.ro.tsv	6/13/2024	Romanian (ro)

## Volumes

Docs	Segments	Unique segments	Tokens	Size
14,468,448	2,001,534,126	204,352 (0.01 %)	23B	118.88 GB

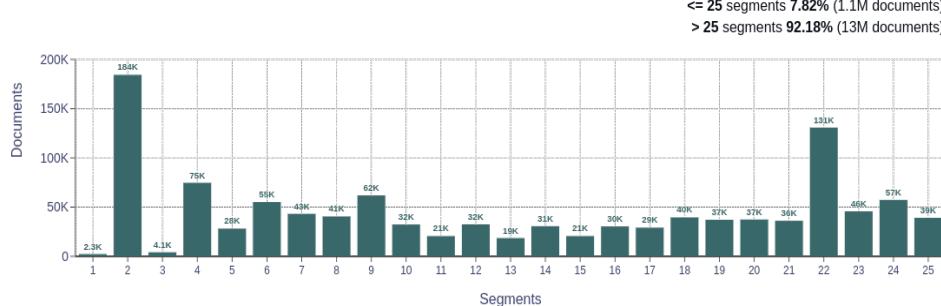
## Top 10 domains

Domain	Docs	% of total
blogspot.ro	949K	6.56
clubafaceri.ro	248K	1.71
reverso.net	213K	1.48
blogspot.com	197K	1.36
business24.ro	186K	1.29
wordpress.com	161K	1.12
9am.ro	111K	0.77
citatepedia.ro	104K	0.72
diebuchsueche.com	102K	0.70
dex.ro	81K	0.56

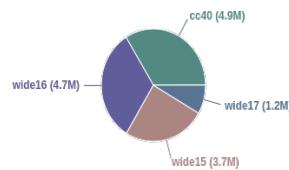
## Top 10 TLDs

Domain	Docs	% of total
ro	10M	70.33
com	2M	13.77
net	578K	3.99
md	348K	2.40
org	226K	1.56
eu	199K	1.37
info	165K	1.14
me	93K	0.64
it	54K	0.37
be	38K	0.26

## Documents size (in segments)

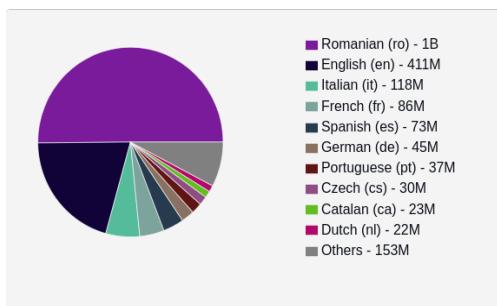


## Documents by collection

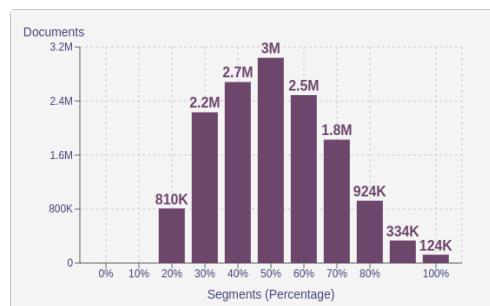


## Language Distribution

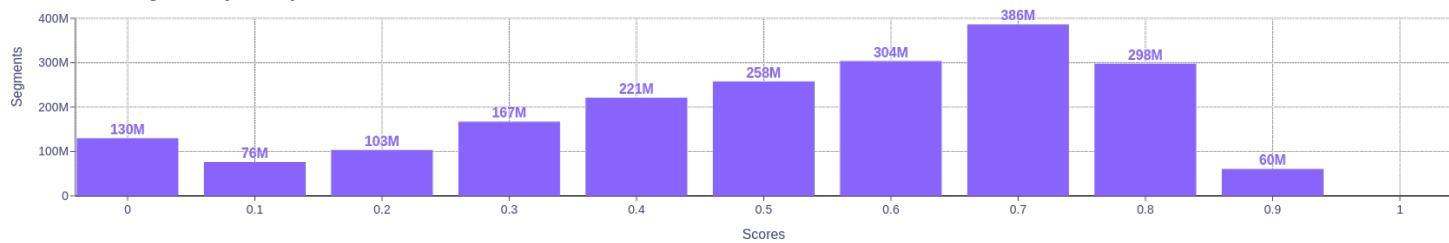
### Number of segments



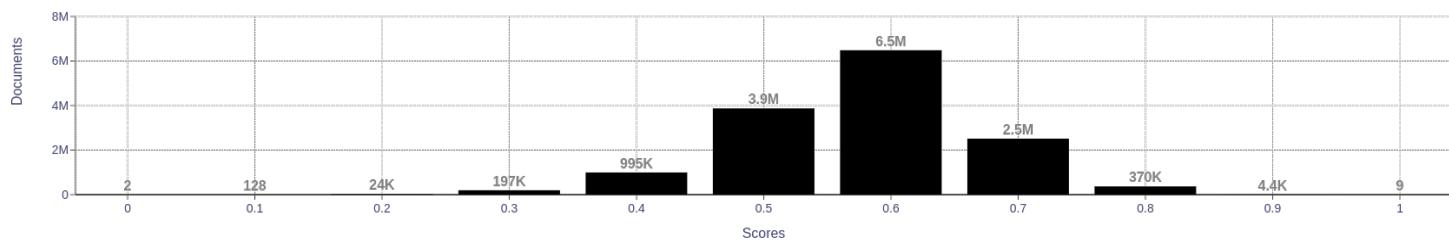
### Percentage of segments in Romanian (ro) inside documents



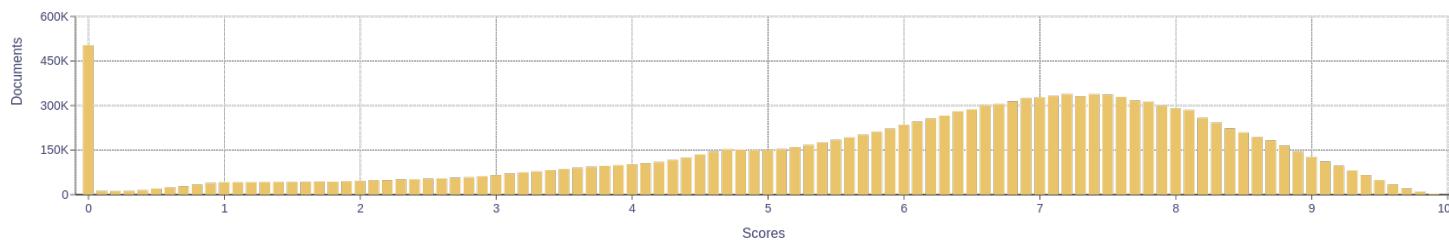
## Distribution of segments by fluency score



## Distribution of documents by average fluency score



## Distribution of documents by document score



## Segment length distribution by token

<= 49 tokens = 305M segments | 1.6B duplicates

> 50 tokens = 91M segments | 32M 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/mabanon/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>