

## General overview

Corpus	Analytics date	Language
be_1.jsonl.tsv	3/21/2024	Belarusian (be)

## Volumes

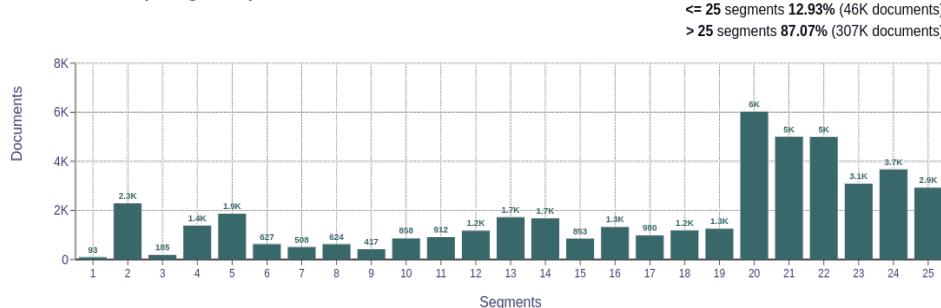
Docs	Segments	Unique segments	Tokens	Size
356,534	38,016,416	33,575 (0.09 %)	517M	4.51 GB

## Type-Token Ratio

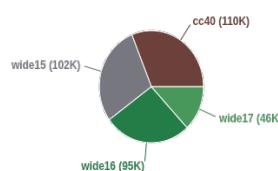
Belarusian (be)

0.01

## Documents size (in segments)

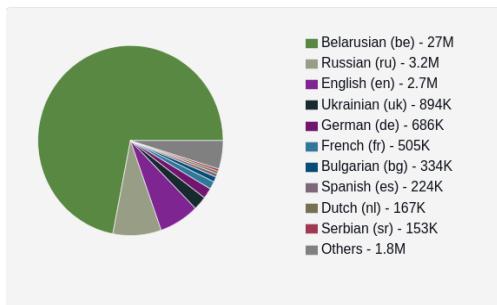


## Documents by collection

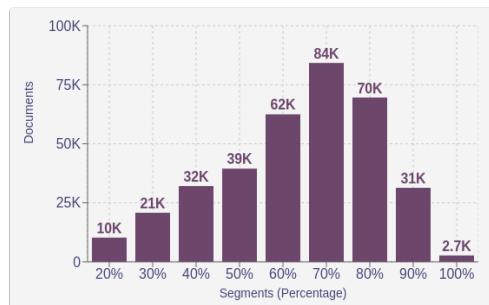


## Language Distribution

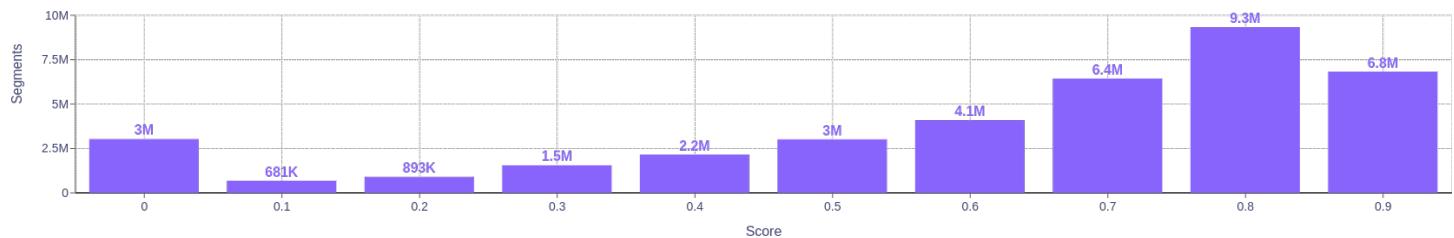
## Number of segments



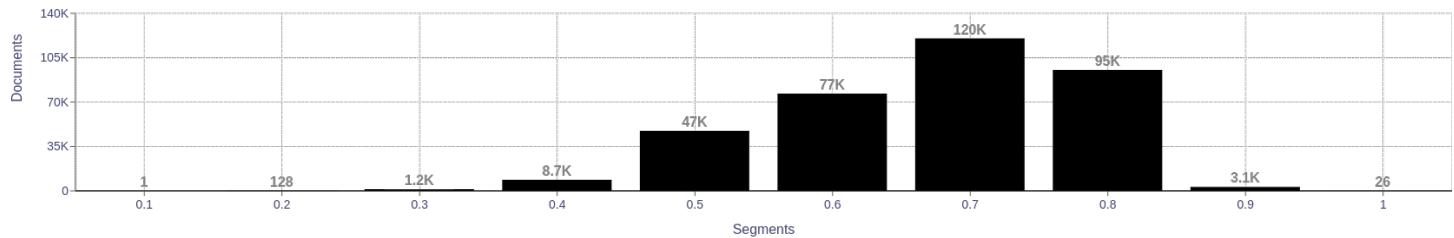
## Percentage of segments in Belarusian (be) inside documents



## Distribution of segments by fluency score



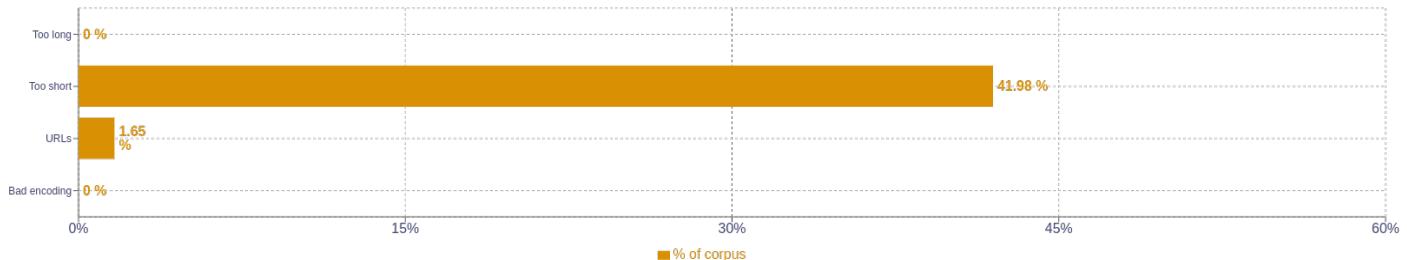
## Distribution of documents by average fluency score



## Segment length distribution by token



## Segment noise distribution



## Frequent n-grams

Size	n-grams
1	( ў   6176577 ) ( да   2343371 ) ( як   1613889 ) ( ад   1160746 ) ( пра   1117709 )
2	( е ў   150570 ) ( рэспублікі беларусь   124645 ) ( ў беларусі   122596 ) ( кропка расы   115818 ) ( судносіны тэмпературы   105229 )
3	( тэмпературы і вільготнасці   105329 ) ( ападкаў не чакаеца   70865 ) ( б в г   64652 ) ( л м н   63522 ) ( к л м   63423 )
4	( судносіны тэмпературы і вільготнасці   105227 ) ( к л м н   63205 ) ( п р с т   63107 ) ( х ц ч ш   62529 ) ( ф х ц ч   62292 )
5	( ф х ц ч ш   62215 ) ( л м н о п   61155 ) ( к л м н о   61112 ) ( м н о п р   61105 ) ( н о п р с   61009 )

## 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>).

### 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.

### 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>