

General overview

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

Volumes

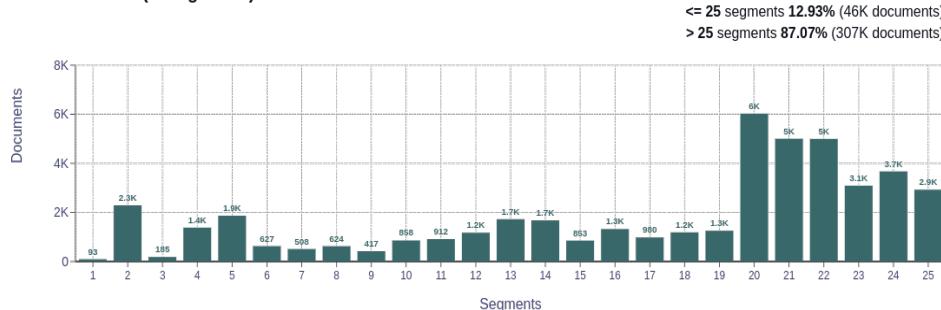
Docs	Segments	Unique segments	Tokens	Size
356,534	38,016,416	33,575 (0.09 %)	517M	4.51 GB

Type-Token Ratio

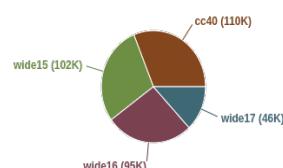
Belarusian (be)

0.01

Documents size (in segments)

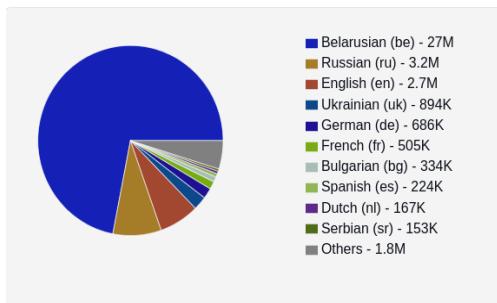


Documents by collection

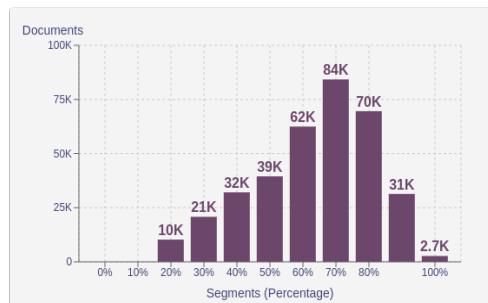


Language Distribution

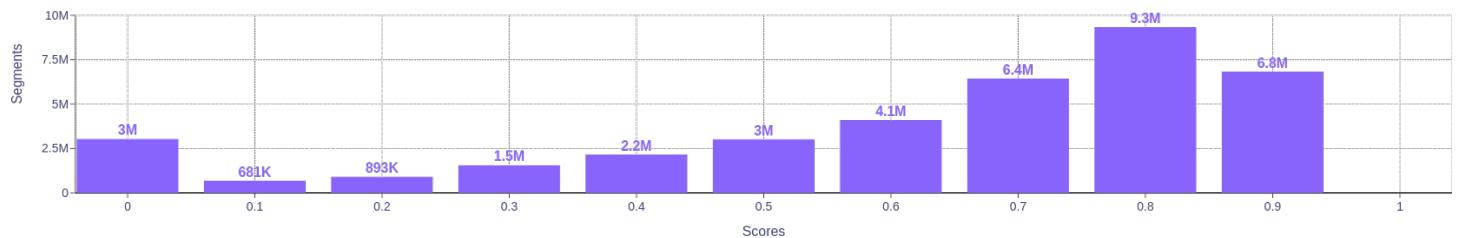
Number of segments



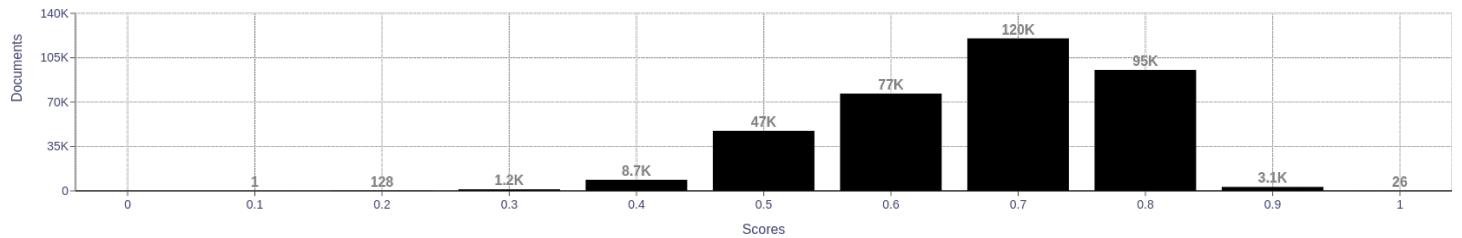
Percentage of segments in Belarusian (be) inside documents



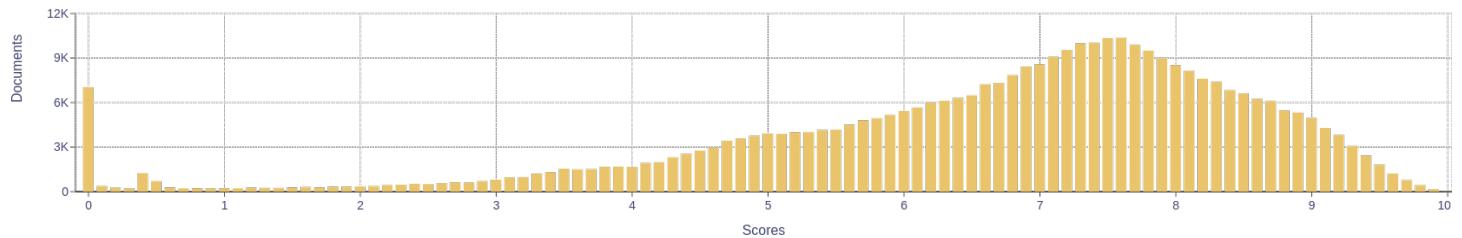
Distribution of segments by fluency score



Distribution of documents by average fluency score

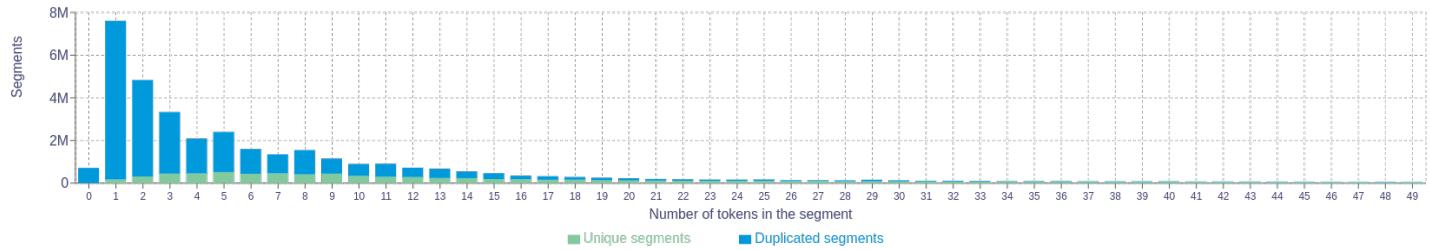


Distribution of documents by document score

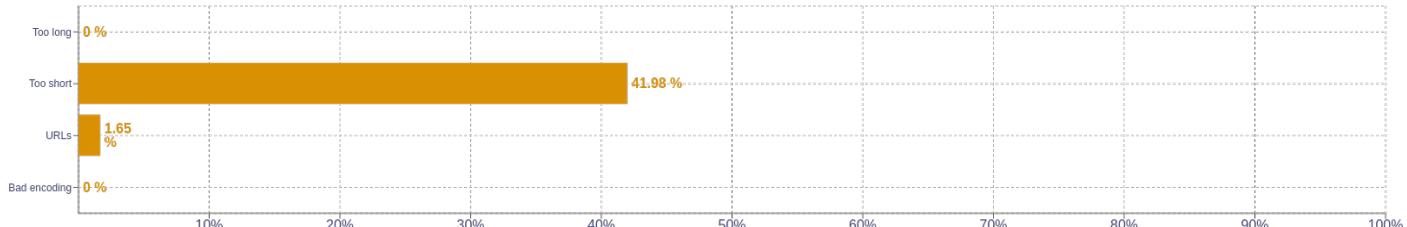


Segment length distribution by token

<= 49 tokens = 8.3M segments | 27M duplicates
 > 50 tokens = 2.4M segments | 436K duplicates



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 беларуси 122596 кропка расы 115818 судносіны тэмпературы 105229
3	тэмпературы і вільготнасці 105329 ападкаў не чакаеца 70865 օ օ օ 64652 լ մ ն 63522 կ լ մ 63423
4	судносіны тэмпературы і вільготнасці 105227 կ լ մ ն 63205 պ ր ս տ 63107 խ չ շ 62529 ֆ խ չ 62292
5	ֆ խ չ շ 62215 լ մ ն օ պ 61155 կ լ մ ն օ 61112 մ ն օ պ ր 61105 հ ո ն օ ր ս 61009

About HPLT Analytics

Volumes - Segments

Segments correspond to paragraph and list boundaries as defined by HTML elements (<p>, , , 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.slinfo.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>, , , 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>