

General overview

Corpus	Analytics date	Language
kk_1.jsonl.tsv	3/22/2024	Kazakh (kk)

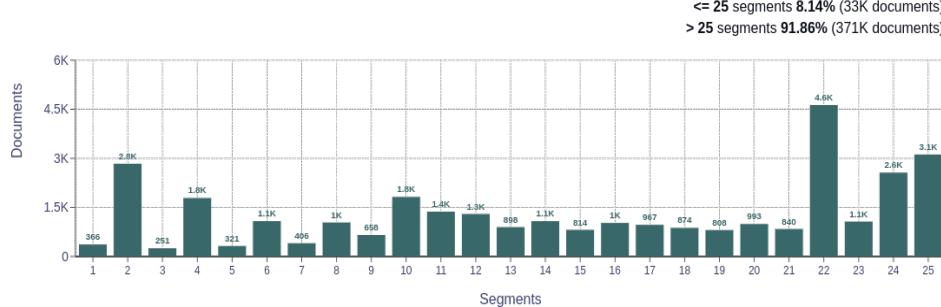
Volumes

Docs	Segments	Unique segments	Tokens	Size
406,351	51,693,572	41,481 (0.08 %)	612M	6.1 GB

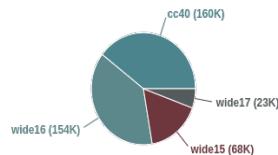
Type-Token Ratio

Kazakh (kk)
0.01

Documents size (in segments)

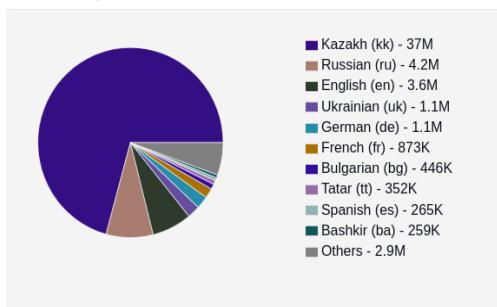


Documents by collection

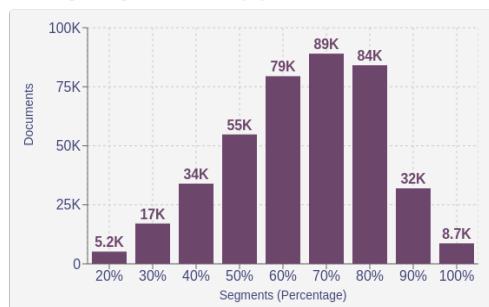


Language Distribution

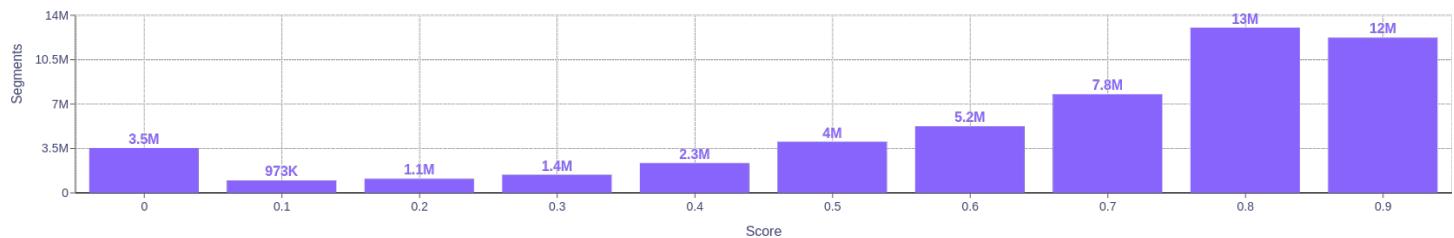
Number of segments



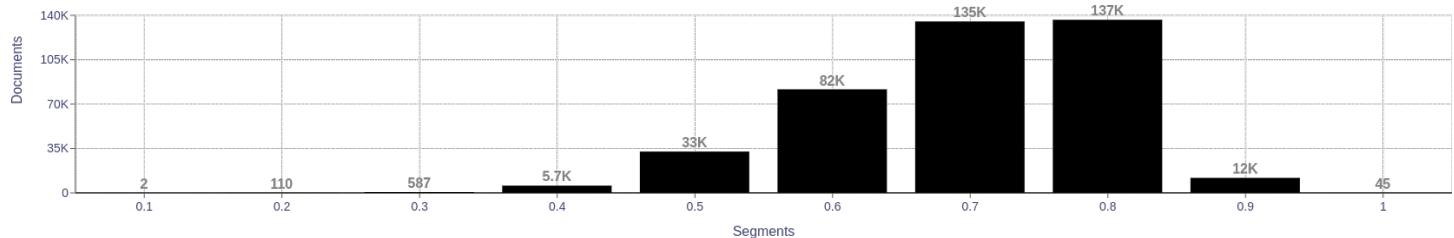
Percentage of segments in Kazakh (kk) inside documents



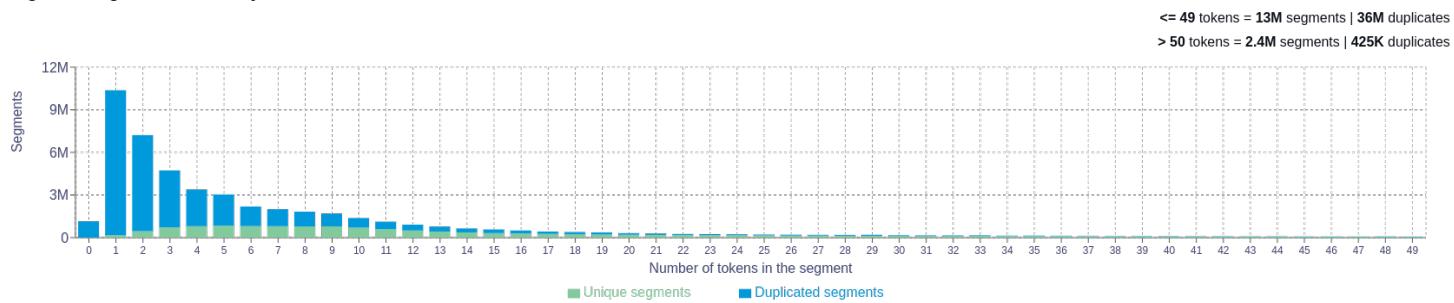
Distribution of segments by fluency score



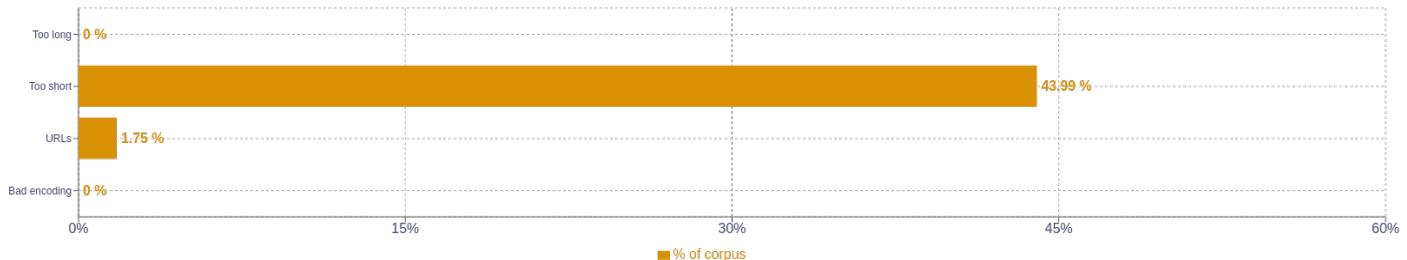
Distribution of documents by average fluency score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	(және 5609164) (қазақстан 1833616) (да 1401038) (бойынша 1400474) (бір 1366271)
2	(қазақстан республикасының 556431) (қазақстан республикасы 385886) (білім беру 268447) (болып табылады 248914) (басқа да 155028)
3	(сыйайлас жемқорлықта қарсы 63128) (өткен соң қолданыска 51531) (және басқа да 51024) (білім және ғылым 50872) (қазақстан республикасы үкіметінің 47140)
4	(күн өткен соң қолданыска 48564) (өткен соң қолданыска енгізіледі 39894) (алғашқы ресми жарияланған күнінен 35181) (күнтізбелік он күн өткен 34256) (тарих және география пәнінен 29632)
5	(он күн өткен соң қолданыска 43772) (күн өткен соң қолданыска енгізіледі 37735) (тарих және география пәнінен үзд 29621) (ресми жарияланған күнінен кейін күнтізбелік 23392) (күнінен кейін күнтізбелік он күн 20478)

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 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>, , , 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>).

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>