

**General overview**

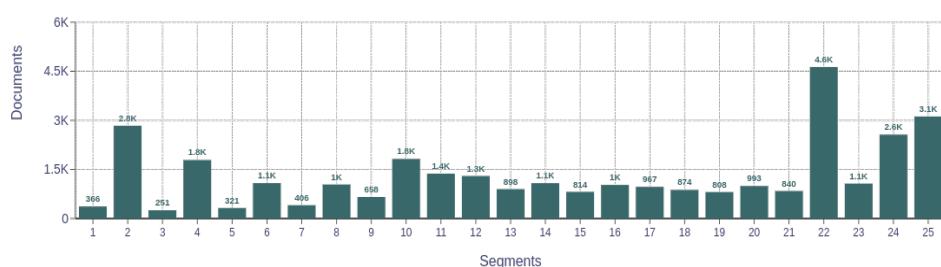
Corpus	Analytics date	Language
kk_1.jsonl.tsv	3/22/2024	Kazakh (kk)

**Volumes**

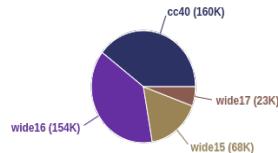
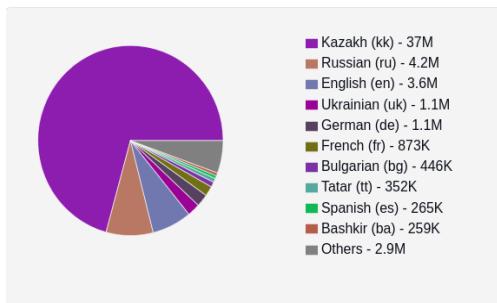
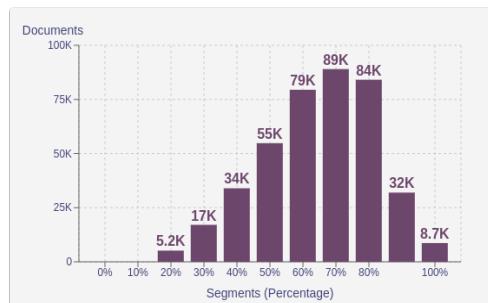
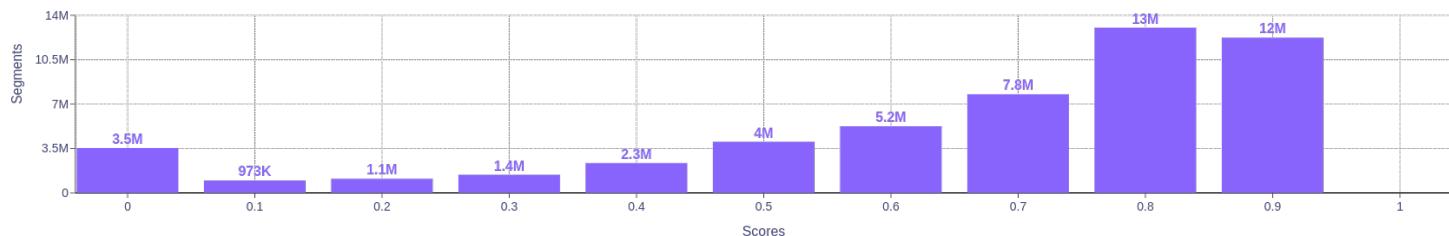
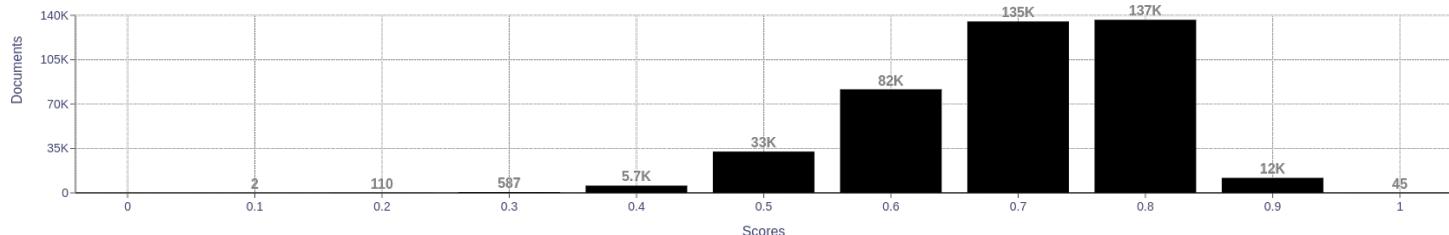
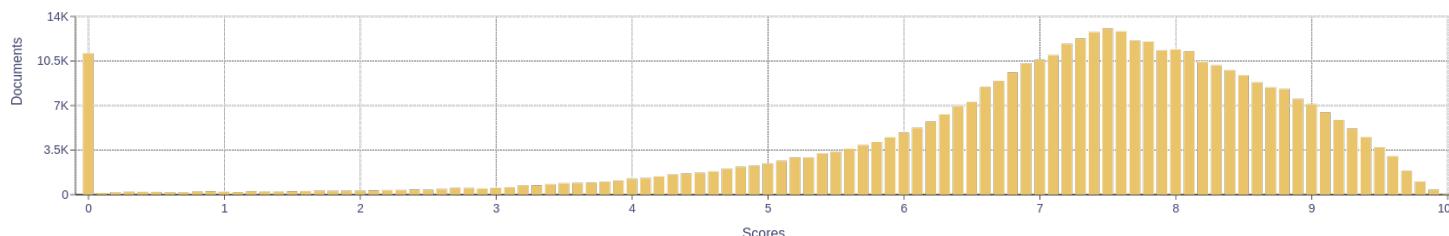
Docs	Segments	Unique segments	Tokens	Size
406,351	51,693,572	41,481 (0.08 %)	612M	6.1 GB

**Type-Token Ratio**

Kazakh (kk)
0.01

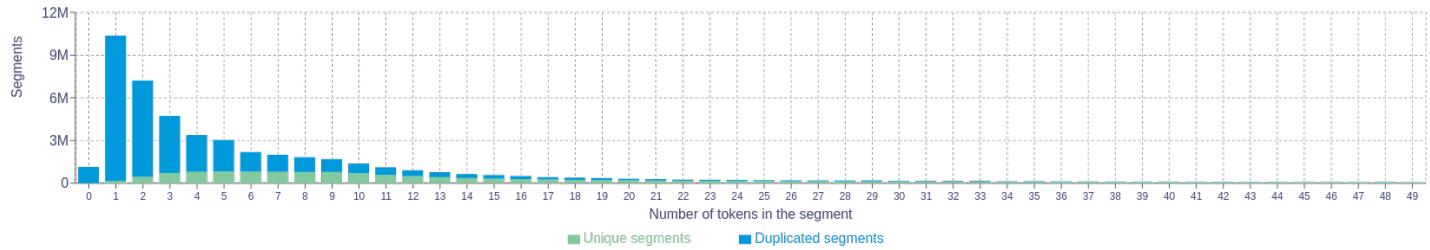
**Documents size (in segments)**

<= 25 segments 8.14% (33K documents)  
> 25 segments 91.86% (371K documents)

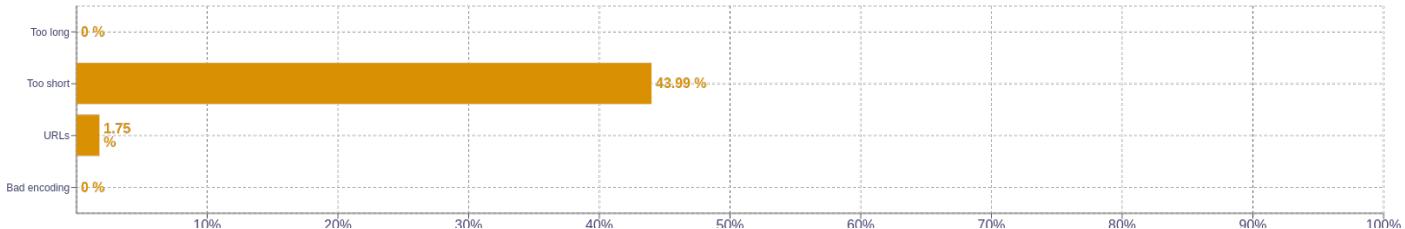
**Documents by collection****Language Distribution****Number of segments****Percentage of segments in Kazakh (kk) inside documents****Distribution of segments by fluency score****Distribution of documents by average fluency score****Distribution of documents by document score**

## Segment length distribution by token

<= 49 tokens = 13M segments | 36M duplicates  
> 50 tokens = 2.4M segments | 425K duplicates



## Segment noise distribution



## Frequent n-grams

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

## 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>