

General overview

Corpus	Analytics date	Language
ky_1.jsonl.tsv	3/17/2024	Kyrgyz (ky)

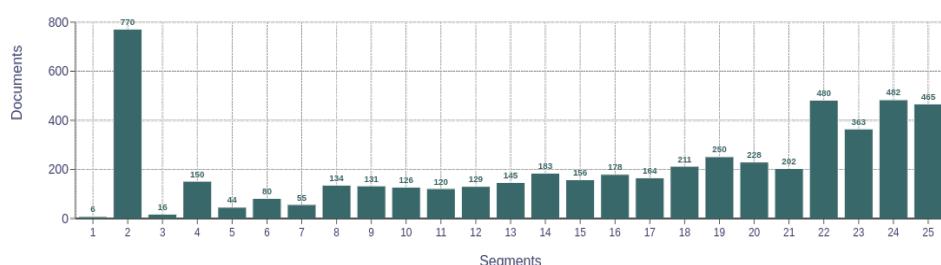
Volumes

Docs	Segments	Unique segments	Tokens	Size
88,322	11,748,340	14,112 (0.12 %)	132M	1.29 GB

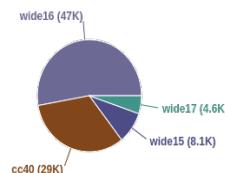
Type-Token Ratio

Kyrgyz (ky)
0.02

Documents size (in segments)

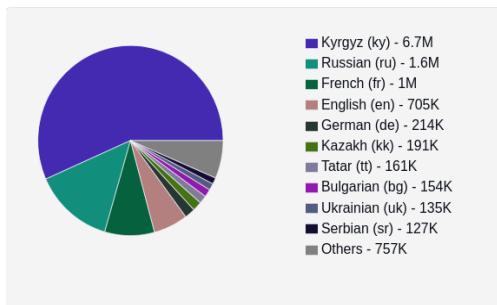


Documents by collection

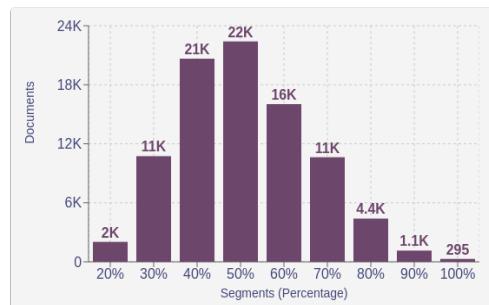


Language Distribution

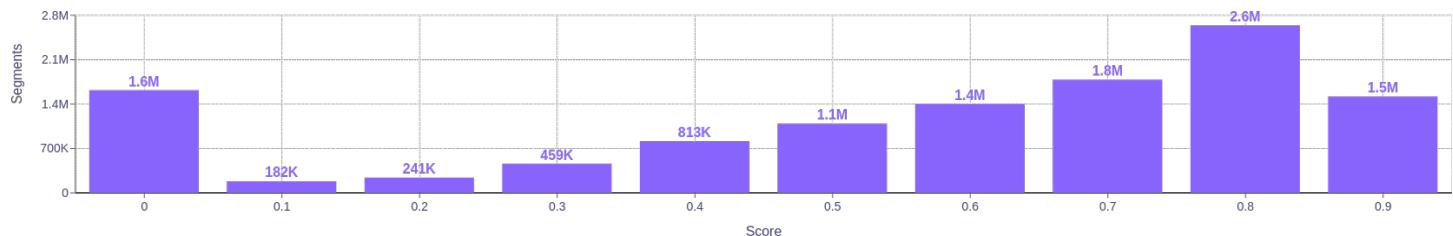
Number of segments



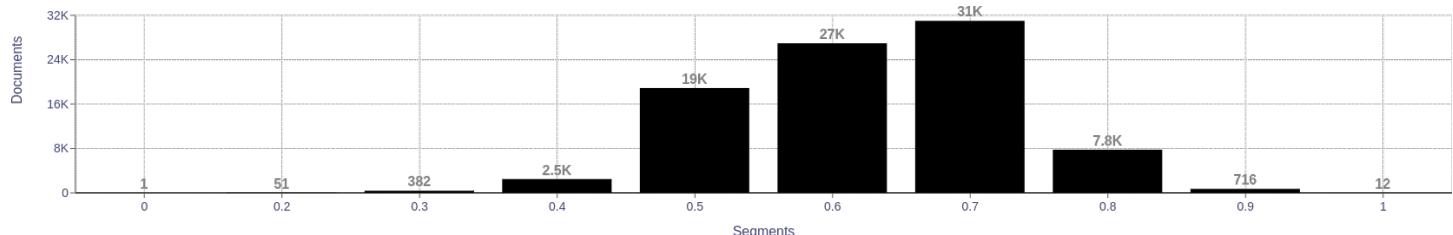
Percentage of segments in Kyrgyz (ky) inside documents



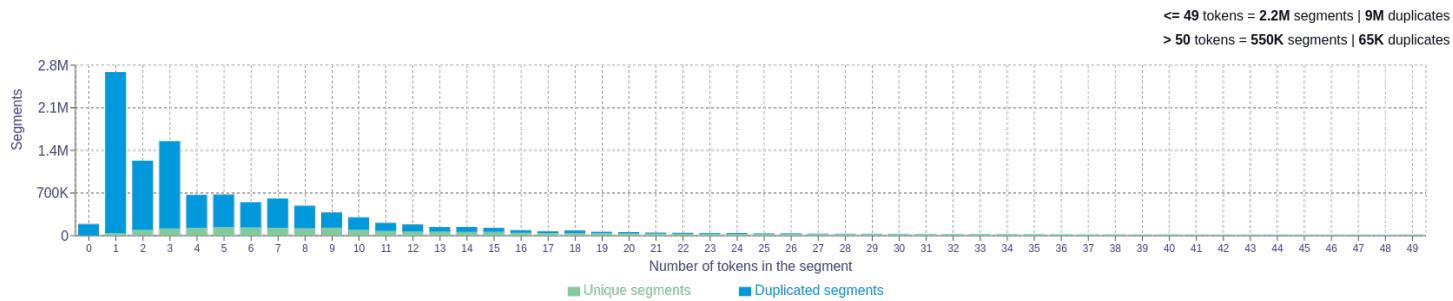
Distribution of segments by fluency score



Distribution of documents by average fluency score

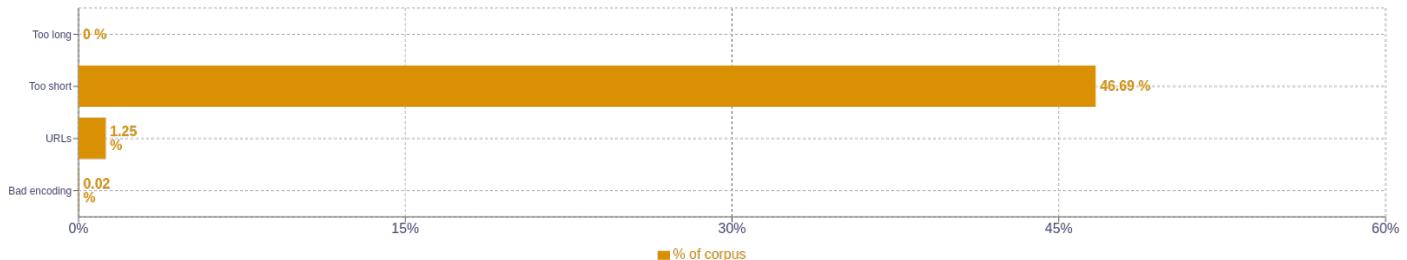


Segment length distribution by token



<= 49 tokens = 2.2M segments | 9M duplicates
> 50 tokens = 550K segments | 65K duplicates

Segment noise distribution



Frequent n-grams

Size	n-grams
1	(бет 382383) (kg 234760) (киргыстан 182462) (жөнүндө 168895) (республикасынын 150480)
2	(билим берүү 33791) (кат келиптири 25909) (курманбек бакиев 25701) (аскар акаев 24623) (министрликтин жообу 23172)
3	(министрликтин жаш кадрларынын 23162) (жаш кадрларынын жоруктары 23162) (уну тарыхый мурас 13710) (эркинтоо биримдик пресс 13710) (энесай новости иссык 13710)
4	(министрликтин жаш кадрларынын жоруктары 23162) (өкмөтү нур эл де 13710) (уну тарыхый мурас обон 13710) (эркинтоо биримдик пресс kg 13710) (экспресс саратан diesel айыл 13710)
5	(өкмөтү нур эл де факт 13710) (уну тарыхый мурас обон ош 13710) (эркинтоо биримдик пресс kg нур 13710) (экспресс саратан diesel айыл өкмөтү 13710) (шамы ош жаңырығы кыргыз руху 13710)

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>