

General overview

Corpus	Analytics date	Language
mn_1.jsonl.tsv	3/26/2024	Mongolian (mn)

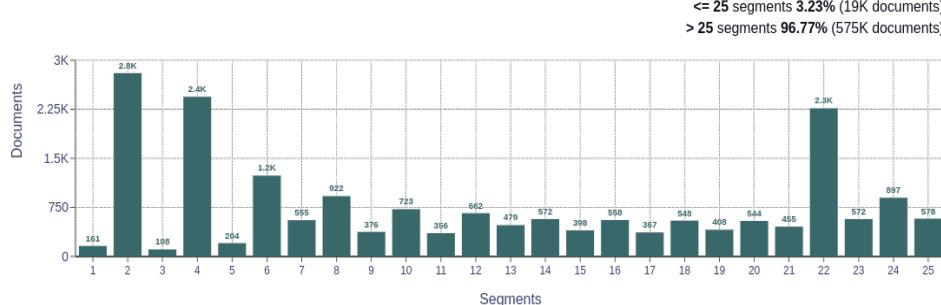
Volumes

Docs	Segments	Unique segments	Tokens	Size
594,905	80,448,202	42,312 (0.05 %)	977M	8.73 GB

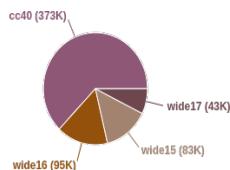
Type-Token Ratio

Mongolian (mn)
0.01

Documents size (in segments)

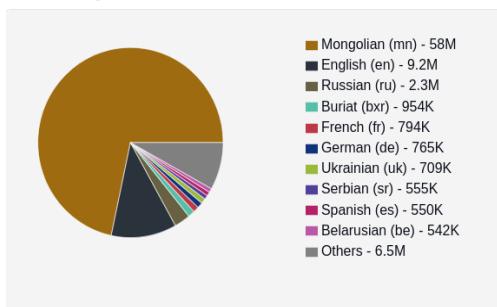


Documents by collection

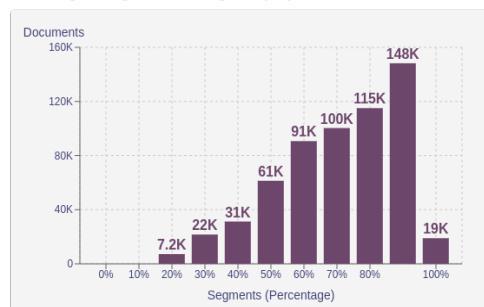


Language Distribution

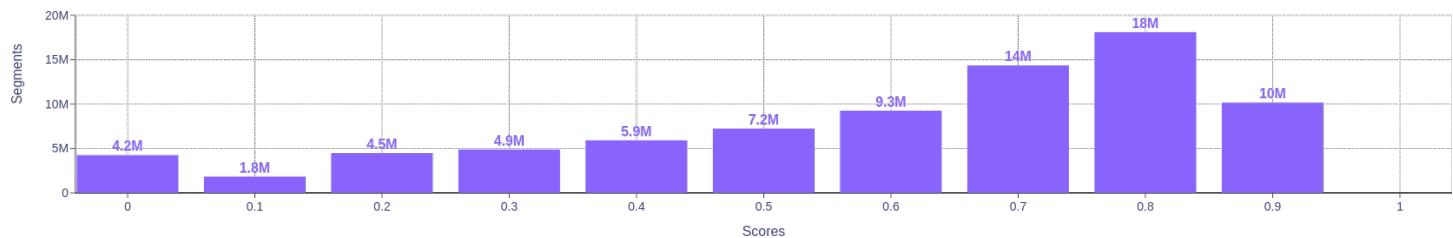
Number of segments



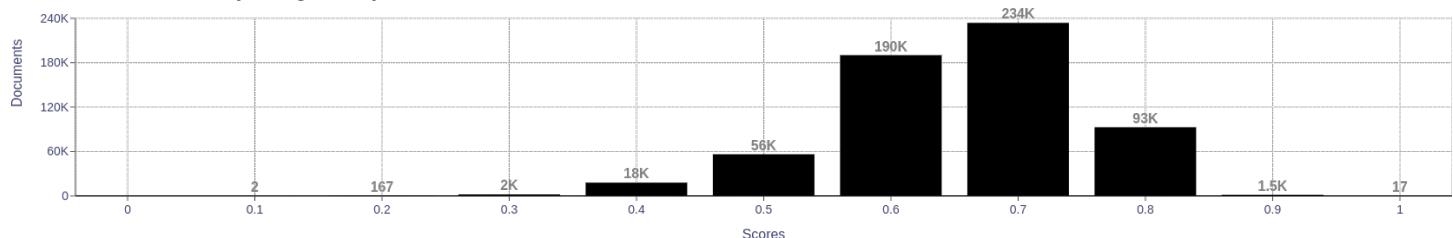
Percentage of segments in Mongolian (mn) inside documents



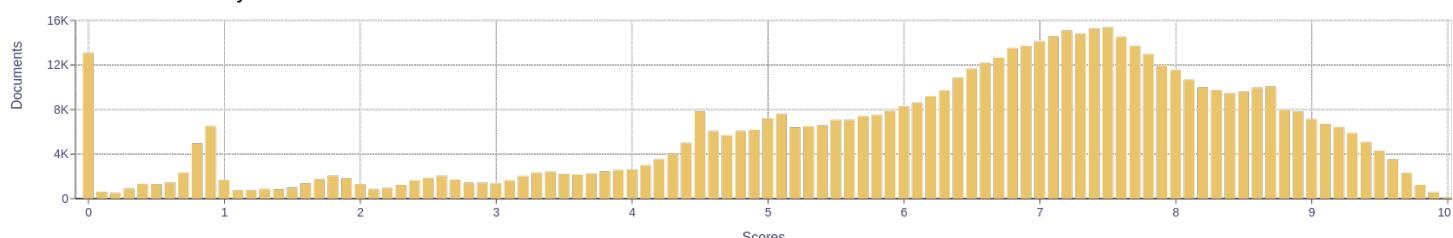
Distribution of segments by fluency score



Distribution of documents by average fluency score

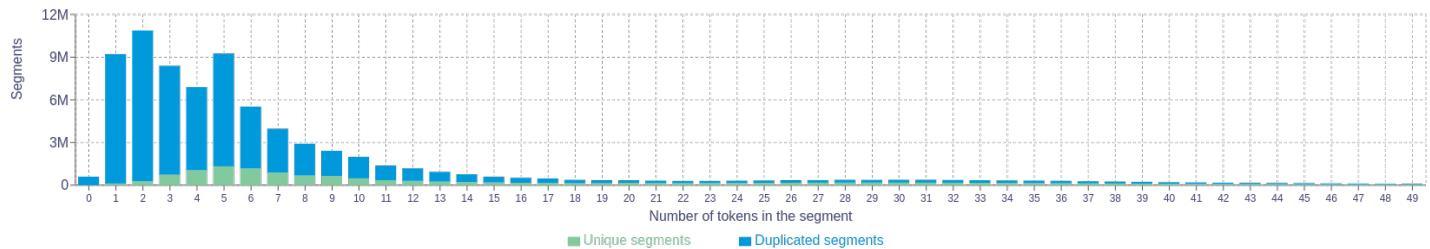


Distribution of documents by document score

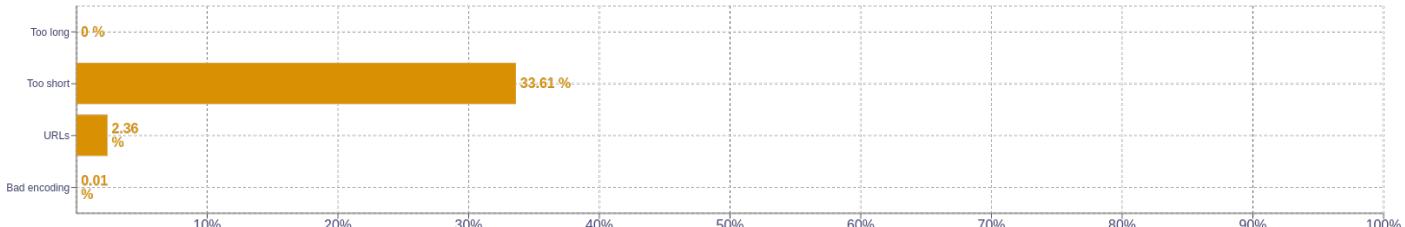


Segment length distribution by token

<= 49 tokens = 13M segments | 63M duplicates
> 50 tokens = 3.7M segments | 848K duplicates



Segment noise distribution



Frequent n-grams

Size	n-grams
1	бутлуур 12542558 чулуу 6178063 бутлуурын 4837091 машин 4818551 үнэ 4051676
2	чулуу бутлуур 2249495 тоног төхөөрөмж 2081635 үнэ авах 1931351 хасарт бутлуур 1782974 уул уурхайн 1601339
3	бидэнтэй холбоо барина 428512 хоёр дахь гар 398470 уул уурхайн тоног 381526 уурхайн тоног төхөөрөмж 326418 чулуу бутлах машин 286444
4	уул уурхайн тоног төхөөрөмж 268677 яг одоо бидэнтэй нэгдээрэй 218968 худалдах хоёр дахь гар 166970 144398 كمساره الحجر كساره الحجر 124304
5	الحجر كساره الحجر كساره الحجر 109413 بутلاх машин хийх элс машин 105919 чулуу бутлах машин хийх элс 104051 машин хийх элс машин чулуу 83600

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.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/mbanov/fastspell>)

Distribution of segments by fluency score

Obtained with Monocleaner (<https://github.com/biteuter/monocleaner>)

Distribution of the *gut* microbiome

Distribution of documents by average fluency score

Obtained with Monocleiner (<https://github.com/bioinfopedia/Monocleiner>)

Distribution of documents by document score

Obtained with Web Docs Scorer (<https://>

Segment length distribution by token

Tokenized with http://

Segment noise distribution

Obtained with Bicle

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->