

General overview

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

Volumes

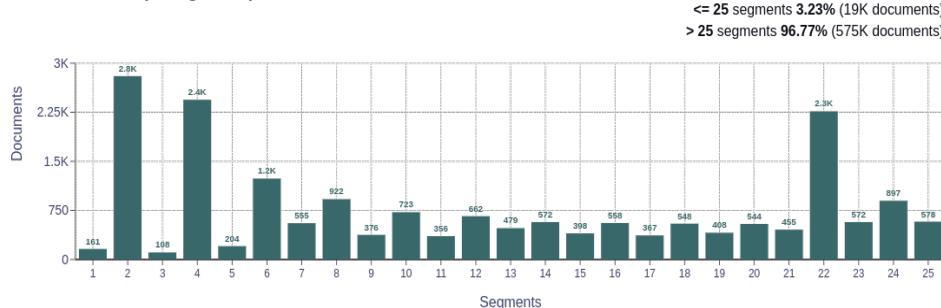
Docs	Segments	Unique segments	Tokens	Size
594,905	80,448,202	42,312 (0.05 %)	977M	8.73 GB

Type-Token Ratio

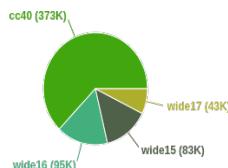
Mongolian (mn)

0.01

Documents size (in segments)

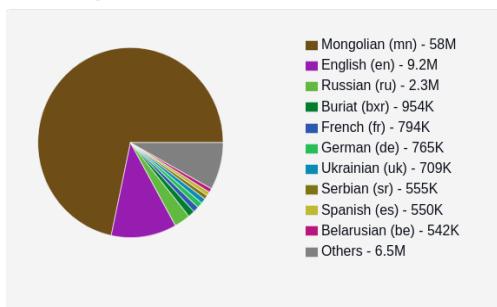


Documents by collection

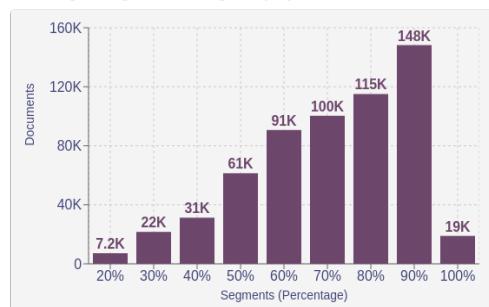


Language Distribution

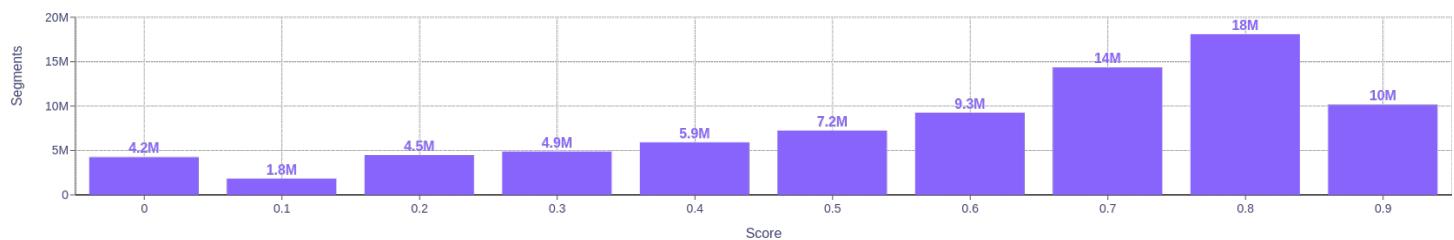
Number of segments



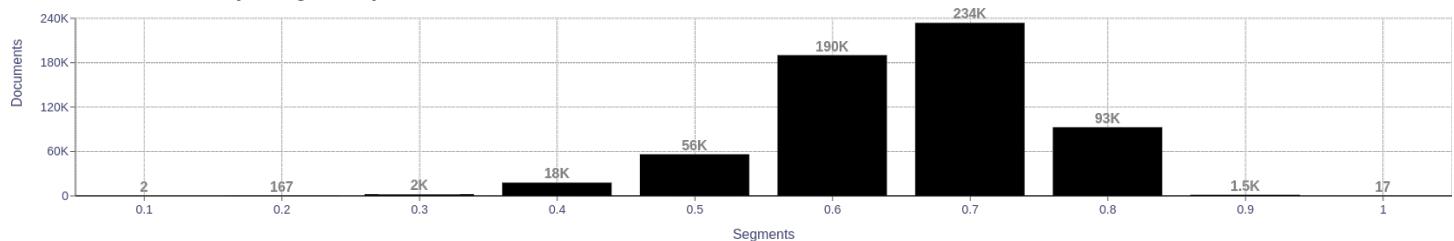
Percentage of segments in Mongolian (mn) inside documents



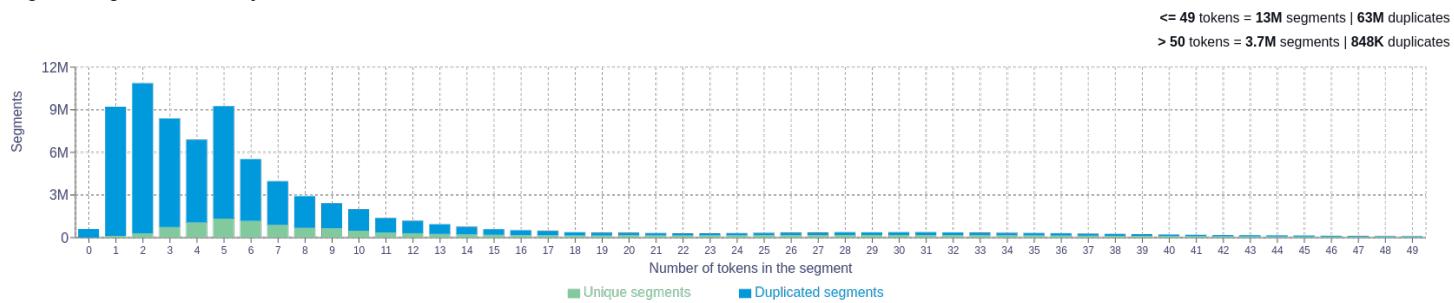
Distribution of segments by fluency score



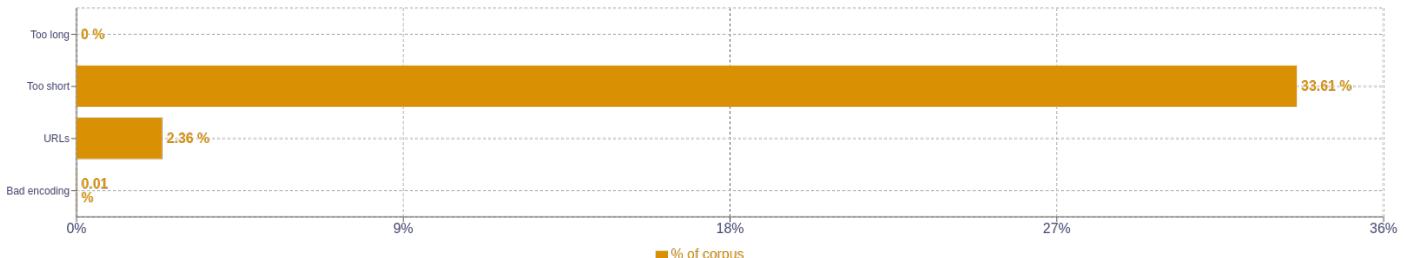
Distribution of documents by average fluency score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	(бутлуур 12542558) (чулуу 6178063) (бутлуурын 4837091) (машин 4018551) (үнэ 4051676)
2	(чулуу бутлуур 2249495) (тоног төхөөрөмж 2081635) (үнэ авах 1931351) (хацарт бутлуур 1782974) (уул уурхайн 1601339)
3	(бидэнтэй холбоо барина 428512) (хоёр дахь гар 398470) (уул уурхайн тоног 381526) (урхайн тоног төхөөрөмж 326418) (чулуу бутлах машин 286444)
4	(уул уурхайн тоног төхөөрөмж 268677) (яг одоо бидэнтэй нэгдээрэй 218968) (худалдах хоёр дахь гар 166978) (144398) (كسارة الحجر كسارة الحجر 124304)
5	(كسارة الحجر كسارة الحجر 109413) (الحجر كسارة الحجر 109455) (бутлах 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 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/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>).

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>