

General overview

Corpus	Date	Language
khk_CyrlJsonl.tsv	9/26/2024	Mongolian (khk)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
2,120,983	53,467,285	24,830,052 (46.44 %)	1.6B	9,275,953,976	11.45 GB

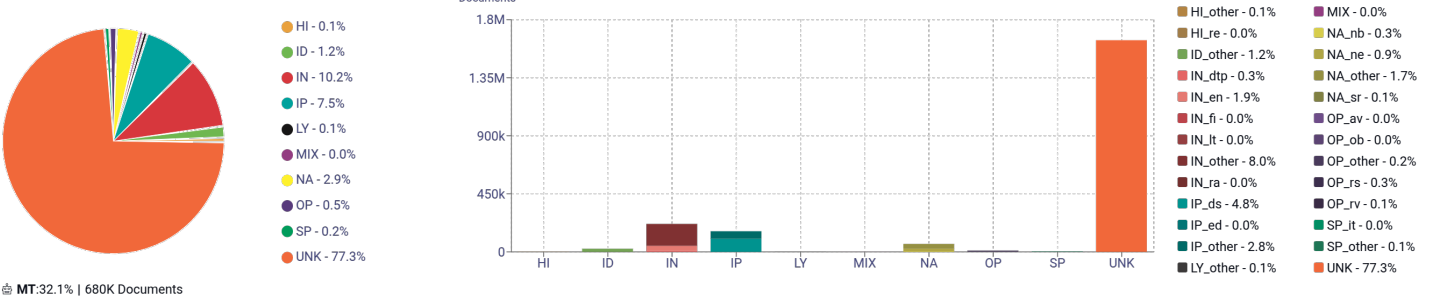
Top 10 domains

Domain	Docs	% of total
wikipedia.org	55K	2.61%
miss.mn	44K	2.06%
fact.mn	37K	1.73%
blogspot.com	36K	1.70%
olloo.mn	31K	1.48%
zindaa.mn	22K	1.02%
vip76.mn	21K	1.00%
shuud.mn	19K	0.92%
montsame.mn	19K	0.88%
ruvr.ru	18K	0.86%

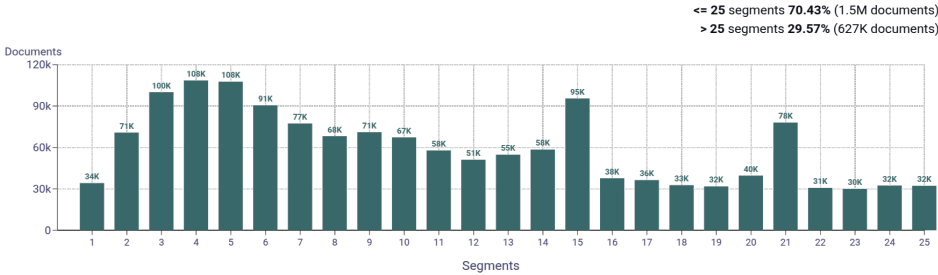
Top 10 TLDs

Domain	Docs	% of total
mn	1.1M	49.73%
com	227K	10.71%
pl	166K	7.81%
nl	110K	5.18%
gov.mn	86K	4.05%
org	83K	3.93%
de	54K	2.55%
be	53K	2.51%
fr	45K	2.11%
es	33K	1.58%

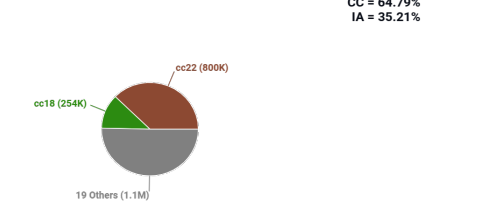
Register labels



Documents size (in segments)

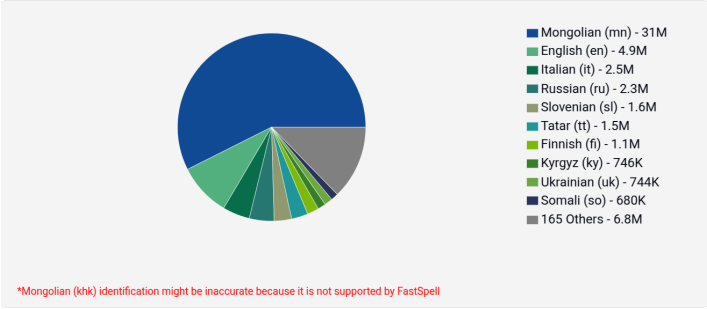


Documents by collection

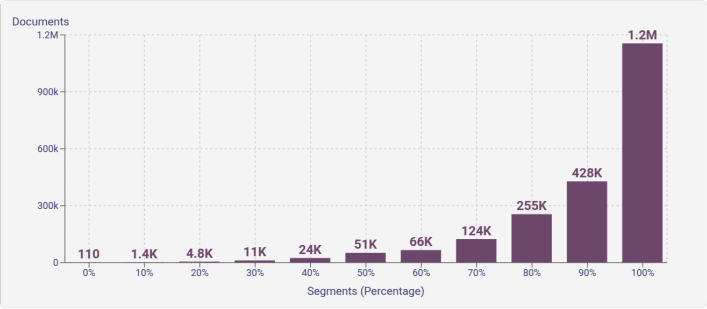


Language Distribution

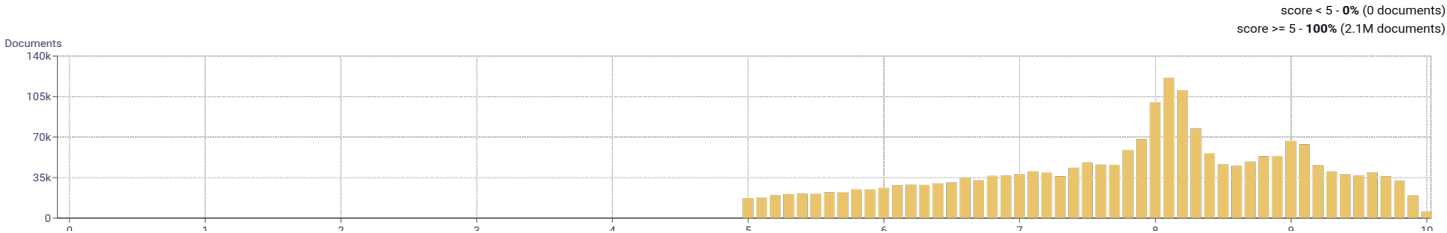
Number of segments in the Mongolian (khk) corpus



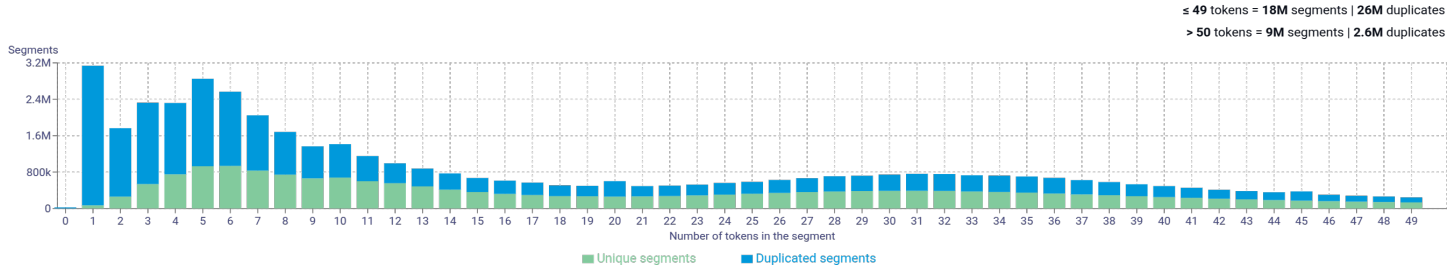
Percentage of segments in Mongolian (khk) inside documents



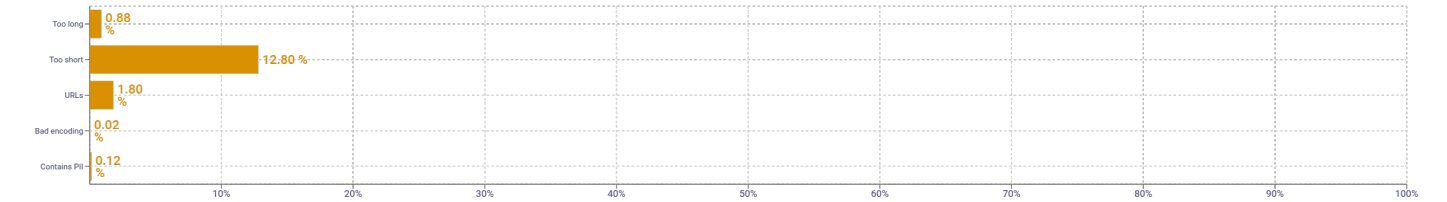
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	hb   23447199   6yтnypp   17520513   чyлny   9467368   6айha   7938528   maиh   7280864
2	чулуу 6yтnypp   3344143   тоhор тeхeepэмж   3061392   хацарт 6yтnypp   2484266   yул yypхайh   2371590   6yтnypp hb   1877872
3	xoёр дахb гар   664007   yул yypхайh тоhор   578192   tph чyлу 6yтnypp   505749   yypхайh тоhор тeхeepэмж   487415   чулуу 6yтлax maиh   465725
4	yул yypхайh тоhор тeхeepэмж   406797   худалдах xoёр дахb гар   299706   xoёр дахb гар hb   288490   6yтnypp hb шoxойh чулуу   236101   6yтлax maиh хийх элс   203175
5	худалдах xoёр дахb гар hb   277487   6yтлax maиh хийх элс maиh   195862   чулуу 6yтлax maиh хийх элс   190785   хуй хуй хуй хуй хуй   167870   хацарт 6yтnypp pe цyвpал хацарт   163694

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/pablopt6n/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>

Register labels

Name	Abbr.	Name	Abbr.	Name	Abbr.
Machine-translated	MT	How-to or instructions	HI	Description of a thing or person	dtp
Lyrical	LY	Recipe	re	FAQ	fi
Spoken	SP	Informational persuasion	IP	Legal terms & conditions	lt
Interview	it	Description with intent to sell	ds	Opinion	OP
Interactive discussion	ID	News & opinion blog or editorial	ed	Review	rv
Narrative	NA	Informational description	IN	Opinion blog	ob
News report	ne	Enciclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				