

General overview

Corpus	Analytics date	Language
khm_khmr.jsonl.tsv	9/26/2024	Khmer (km)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
700,992	9,864,172	5,800,059 (58.80 %)	1.5B	5.48 GB	2,113,232,593

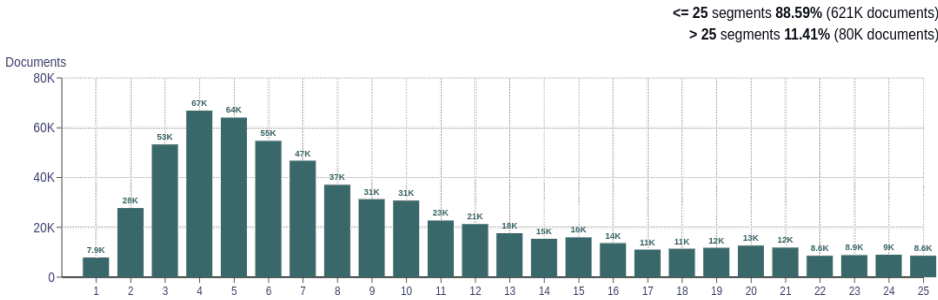
Top 10 domains

Domain	Docs	% of total
khtoem.com	22K	3.16
voanews.com	21K	3.07
wordpress.com	18K	2.56
monoroom.info	16K	2.30
khmread.com	14K	1.98
nokorwatnews.com	12K	1.78
freshnewsasia.com	11K	1.57
sabay.com.kh	10K	1.44
rasmeinews.com	9.2K	1.31
postkhmer.com	8.6K	1.23

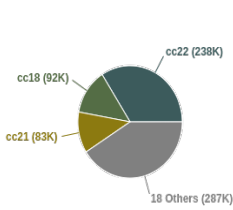
Top 10 TLDs

Domain	Docs	% of total
com	419K	59.75
com.kh	58K	8.27
icu	38K	5.44
org	37K	5.31
gov.kh	34K	4.86
info	26K	3.67
net	17K	2.39
vn	13K	1.90
news	11K	1.61
org.kh	11K	1.52

Documents size (in segments)

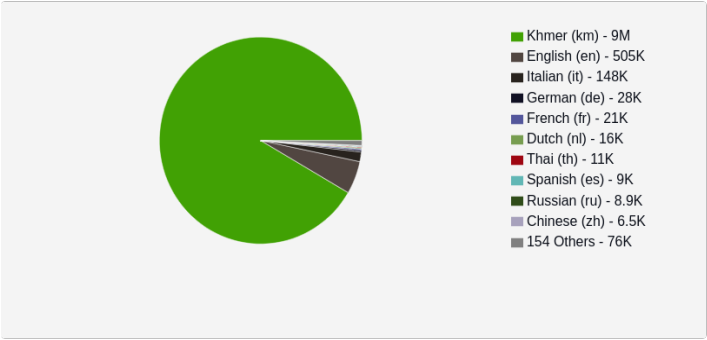


Documents by collection

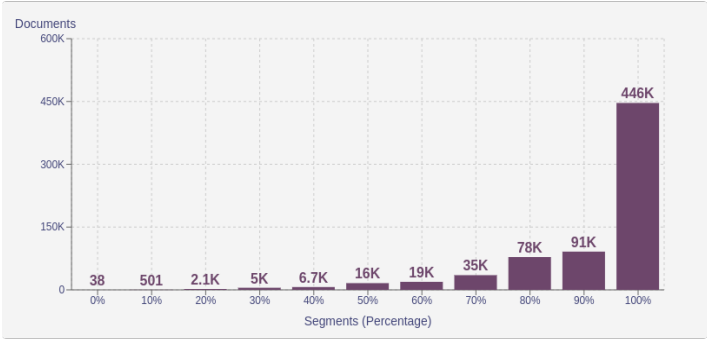


Language Distribution

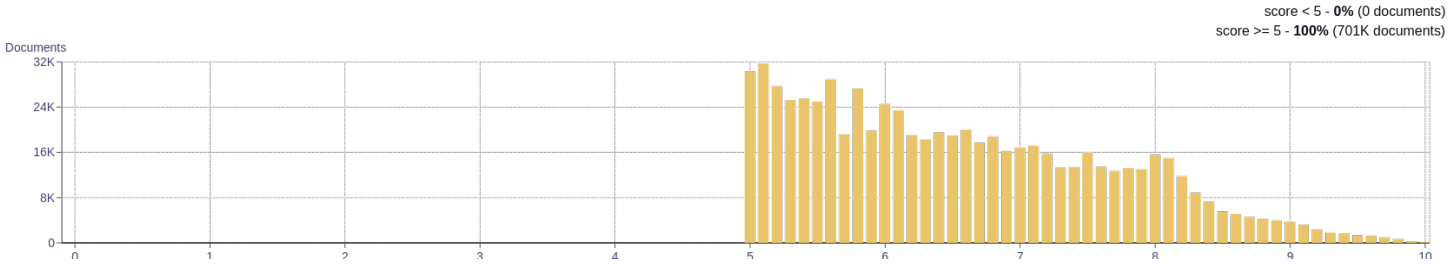
Number of segments



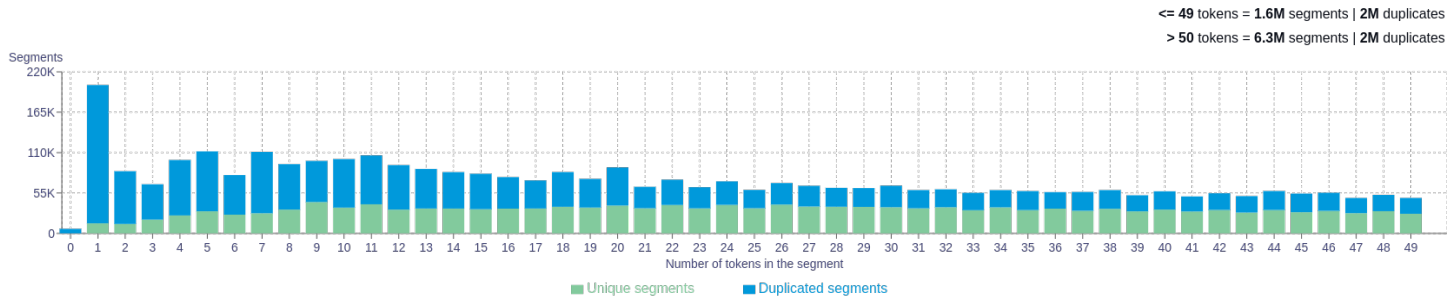
Percentage of segments in Khmer (km) inside documents



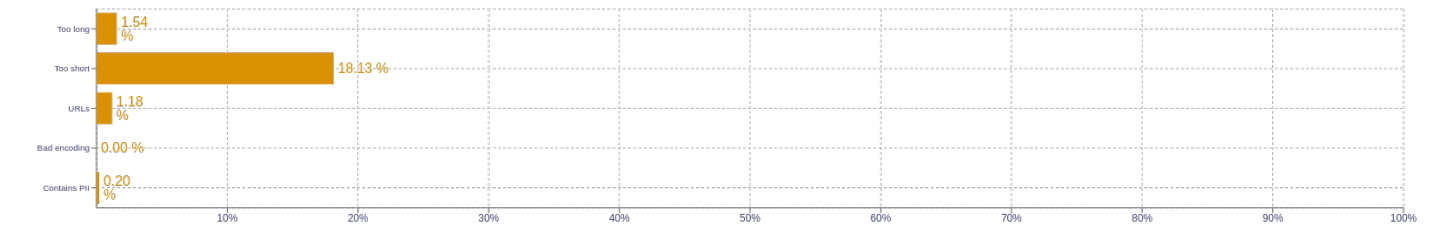
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div><div>៖   47850531</div><div>ិ   40593396</div><div>ី   35971657</div><div>ី   31393488</div><div>ុ   31022979</div></div>
2	<div><div>៖ ៖   503541</div><div>world cup   502430</div><div>ិ ៖   414159</div><div>៖ ័   407861</div><div>ំ ៖   405408</div></div>
3	<div><div>fifa world cup   45367</div><div>ំ world cup   45159</div><div>qatar world cup   35895</div><div>ី ័ ័   32137</div><div>world cup qatar   31281</div></div>
4	<div><div>fifa world cup qatar   15154</div><div>pro evolution soccer ័   7634</div><div>opens in new window   5251</div><div>click to share on   5250</div><div>័័័័ world cup ៖   4833</div></div>
5	<div><div>to cnrp radio at national   2547</div><div>read more khmer news and   2547</div><div>radio at national rescue party   2547</div><div>please read more khmer news   2547</div><div>news and listen to cnrp   2547</div></div>

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>