

General overview

Corpus	Analytics date	Language
HPLT-v2-tha_Thai.tsv	9/23/2024	Thai (th)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
17,703,323	339,050,942			152.12 GB	59,655,264,386

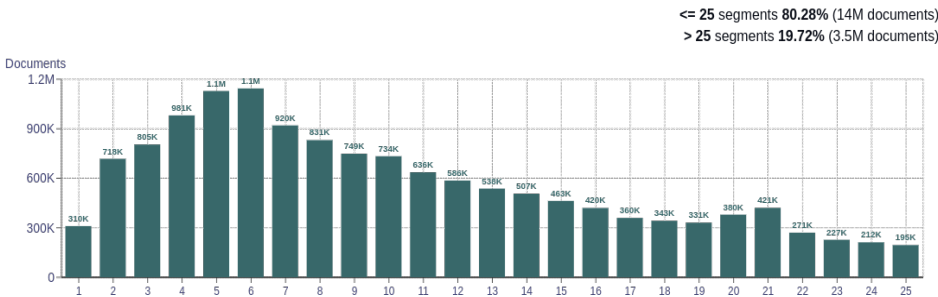
Top 10 domains

Domain	Docs	% of total
blogspot.com	330K	1.86
sanook.com	255K	1.44
wikipedia.org	229K	1.29
mthai.com	221K	1.25
thairath.co.th	194K	1.10
tripadvisor.com	172K	0.97
wordpress.com	150K	0.85
plazathai.com	133K	0.75
ryt9.com	129K	0.73
newsmit.com	119K	0.67

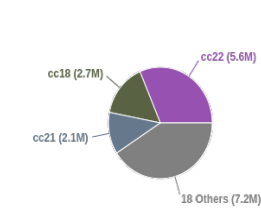
Top 10 TLDs

Domain	Docs	% of total
com	12M	69.22
net	997K	5.63
org	995K	5.62
co.th	766K	4.33
in.th	320K	1.81
ac.th	220K	1.24
co	217K	1.23
go.th	198K	1.12
or.th	191K	1.08
tk	171K	0.97

Documents size (in segments)

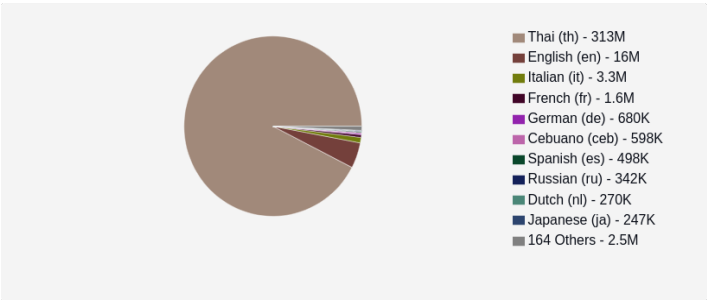


Documents by collection

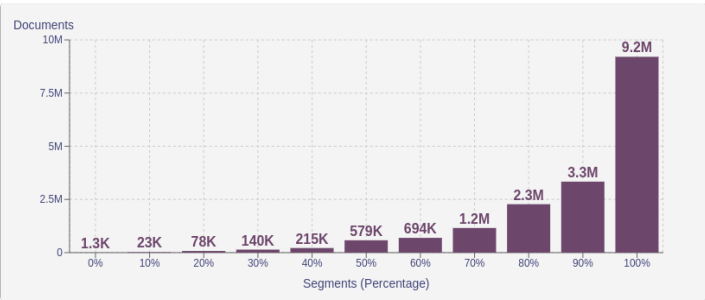


Language Distribution

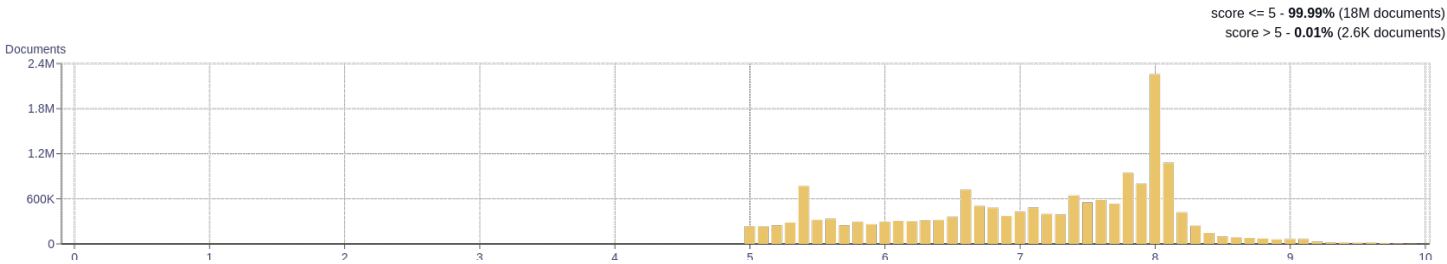
Number of segments



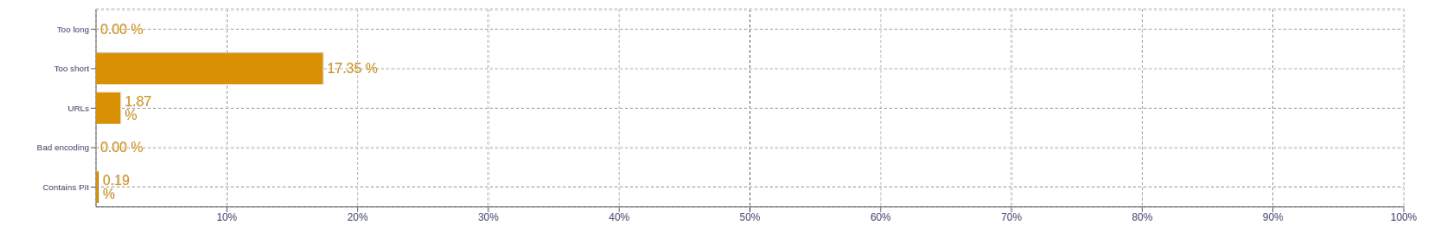
Percentage of segments in Thai (th) inside documents



Distribution of documents by document score



Segment noise distribution



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>).

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>