

General overview

Corpus	Analytics date	Language
HPLT-docslite.nl.tsv	6/27/2024	Dutch (nl)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
31,745,184	3,670,328,667			201.01 GB	

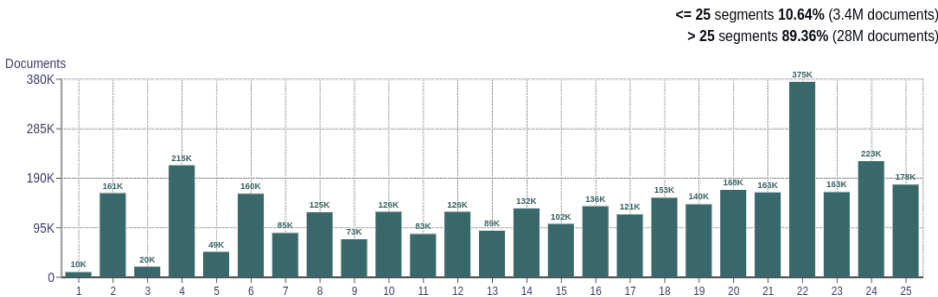
Top 10 domains

Domain	Docs	% of total
blogspot.nl	809K	2.55
alibaba.com	549K	1.73
blogspot.be	336K	1.06
docplayer.nl	233K	0.73
diebuchsuche.com	230K	0.73
blogspot.com	201K	0.63
aliexpress.com	131K	0.41
made-in-china.com	120K	0.38
wordpress.com	114K	0.36
takedrivers.nl	101K	0.32

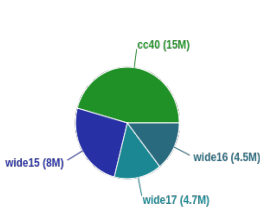
Top 10 TLDs

Domain	Docs	% of total
nl	19M	60.90
com	5.3M	16.67
be	4.2M	13.10
org	521K	1.64
eu	433K	1.36
net	418K	1.32
info	208K	0.65
nu	202K	0.64
de	112K	0.35
me	92K	0.29

Documents size (in segments)

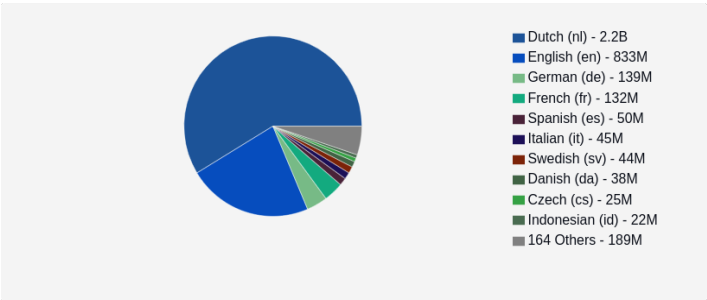


Documents by collection

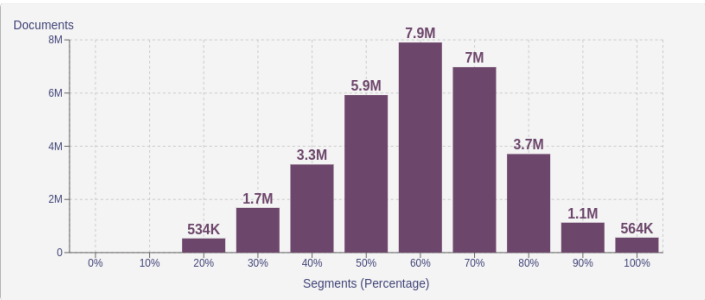


Language Distribution

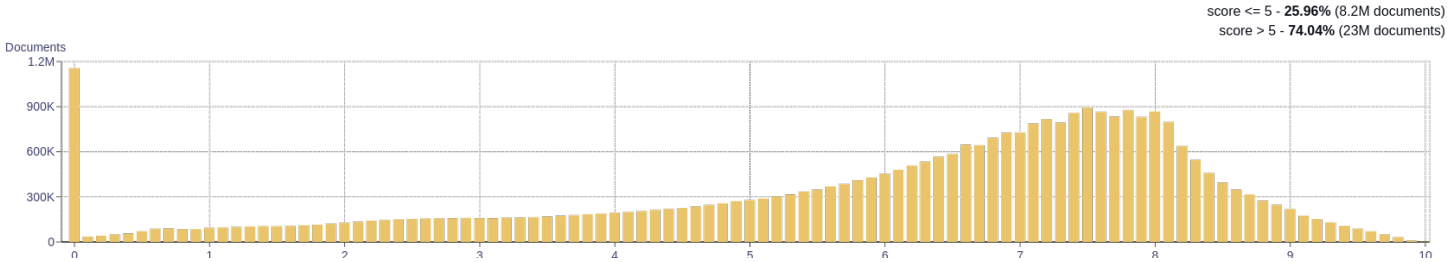
Number of segments



Percentage of segments in Dutch (nl) inside documents



Distribution of documents by document score



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 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>, <ul>, <ol>, etc.) replaced by newlines.

Language distribution

Language identified with FastSpell (<https://github.com/mbanor/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>