

General overview

Corpus	Analytics date	Language
HPLT-v2-dan_Latn.tsv	9/20/2024	Danish (da)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
33,841,408	873,018,899			126.49 GB	132,541,320,753

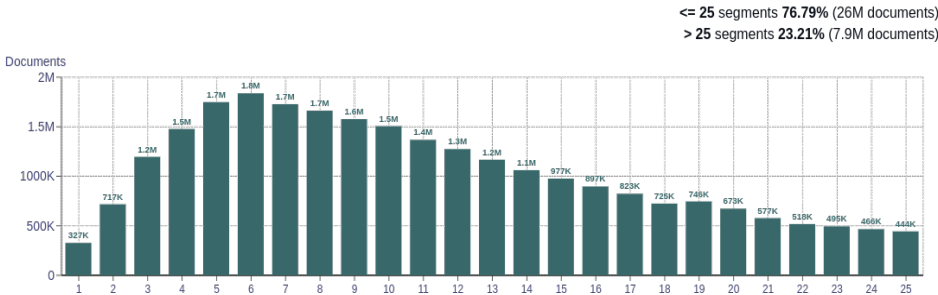
Top 10 domains

Domain	Docs	% of total
docplayer.dk	639K	1.89
blogspot.com	541K	1.60
wikipedia.org	458K	1.35
billedeverden.com	444K	1.31
blogspot.dk	271K	0.80
tripadvisor.dk	239K	0.71
dagens.dk	234K	0.69
avisen.dk	155K	0.46
wordpress.com	147K	0.43
ekstrabladet.dk	144K	0.43

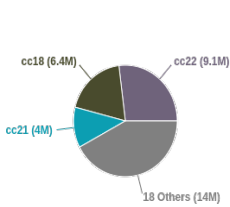
Top 10 TLDs

Domain	Docs	% of total
dk	26M	76.44
com	4.8M	14.11
org	771K	2.28
eu	644K	1.90
net	350K	1.03
nu	210K	0.62
info	179K	0.53
no	105K	0.31
se	92K	0.27
de	71K	0.21

Documents size (in segments)

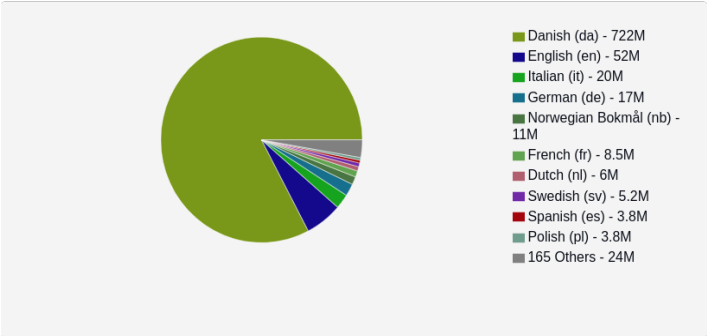


Documents by collection

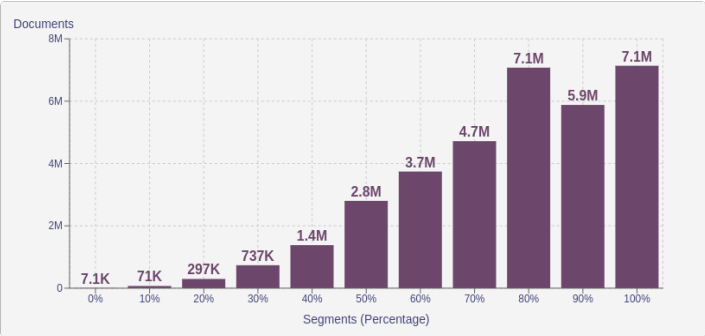


Language Distribution

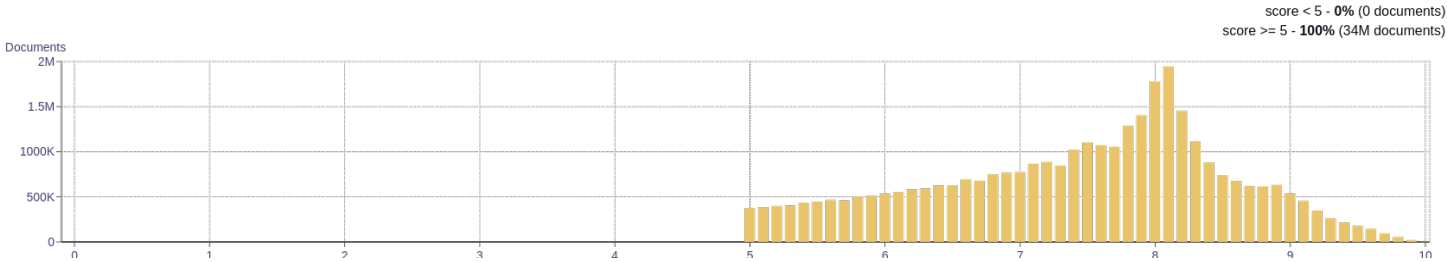
Number of segments



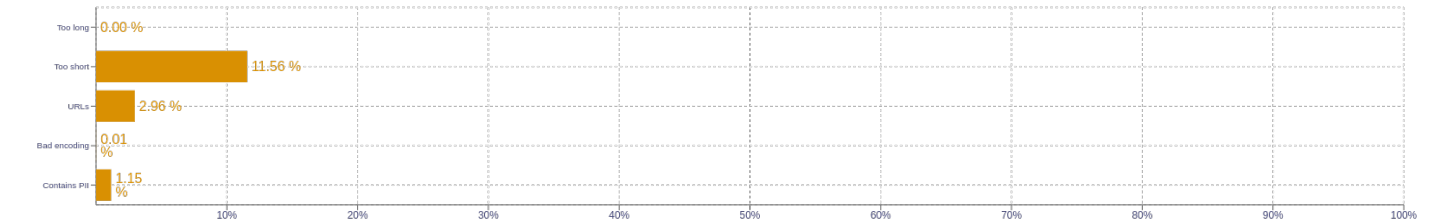
Percentage of segments in Danish (da) inside documents



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>