

General overview

Corpus	Analytics date	Language
HPLT-v2-heb_Hebr.tsv	9/17/2024	Hebrew (he)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
17,116,813	466,627,647			92.68 GB	56,374,343,947

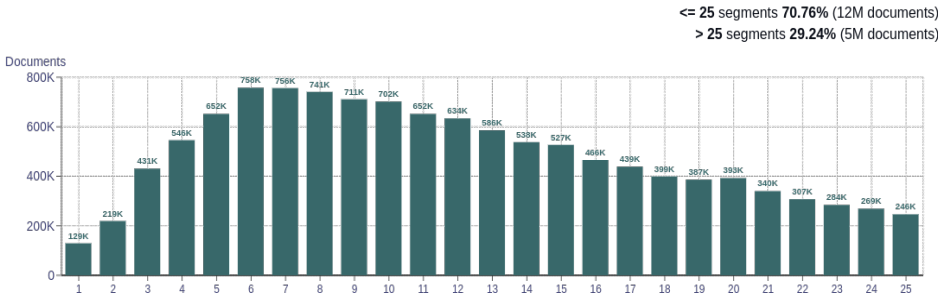
Top 10 domains

Domain	Docs	% of total
wikipedia.org	862K	5.03
psika.net	318K	1.85
blogspot.com	244K	1.42
blogspot.co.il	237K	1.39
articles.co.il	234K	1.37
walla.co.il	195K	1.14
wordpress.com	185K	1.08
mako.co.il	165K	0.96
nana10.co.il	146K	0.85
ynet.co.il	142K	0.83

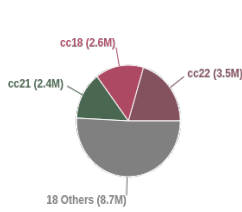
Top 10 TLDs

Domain	Docs	% of total
co.il	9M	52.69
com	3.8M	22.01
org	1.5M	9.01
org.il	1.1M	6.57
net	721K	4.21
ac.il	237K	1.39
info	92K	0.54
gov.il	56K	0.33
muni.il	51K	0.30
me	36K	0.21

Documents size (in segments)

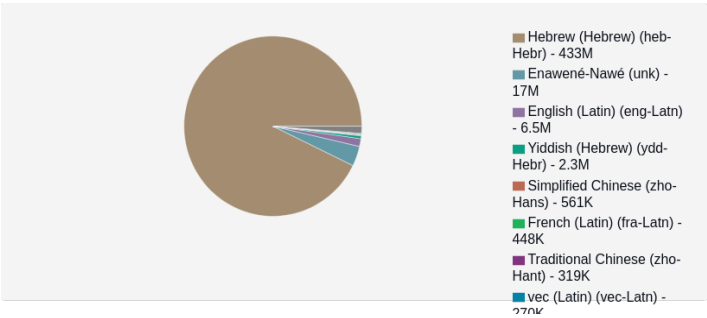


Documents by collection

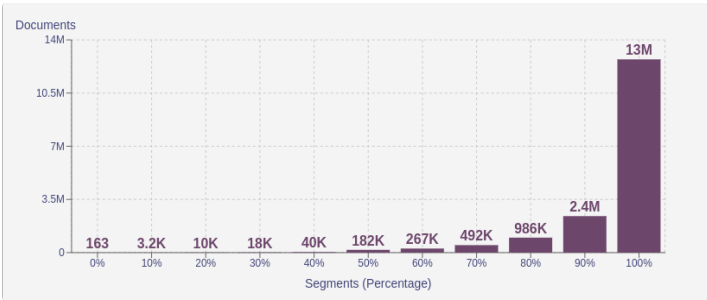


Language Distribution

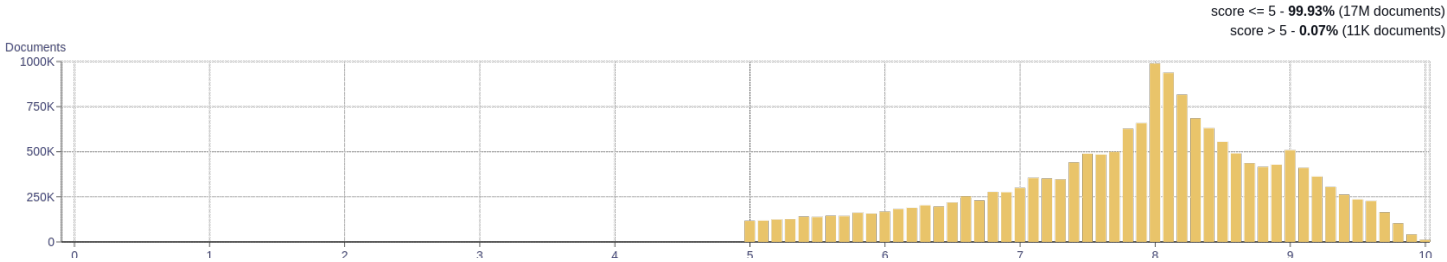
Number of segments



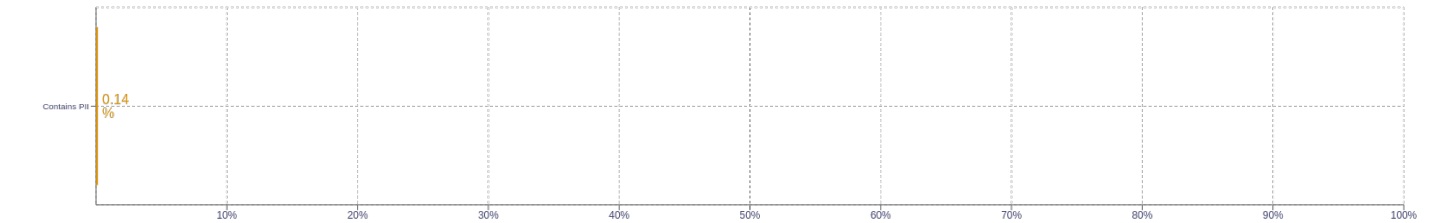
Percentage of segments in Hebrew (he) inside documents



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>, <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>