

General overview

Corpus	Analytics date	Language
mk_1.jsonl.tsv	3/22/2024	Macedonian (mk)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
734,687	86,189,730	14,489,281 (16.81 %)	845M	7.59 GB	

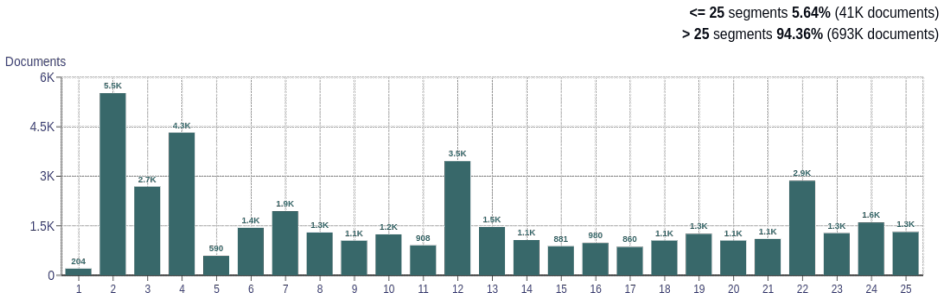
Top 10 domains

Domain	Docs	% of total
sport.com.mk	23K	3.19
daily.mk	20K	2.78
meta.mk	15K	2.10
kurir.mk	15K	2.07
wikipedia.org	15K	1.98
rbth.com	13K	1.81
voanews.com	12K	1.58
mia.mk	9.9K	1.35
nmd.mk	7.6K	1.03
press24.mk	7.1K	0.97

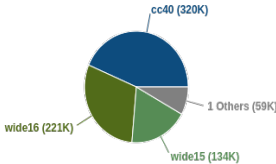
Top 10 TLDs

Domain	Docs	% of total
mk	428K	58.24
com	111K	15.14
com.mk	93K	12.62
org	26K	3.58
org.mk	21K	2.87
gov.mk	16K	2.12
net	7.7K	1.05
edu.mk	7.6K	1.03
news	3.3K	0.45
info	2.4K	0.32

Documents size (in segments)

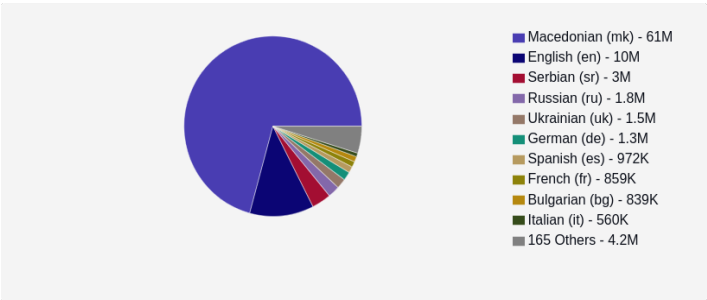


Documents by collection

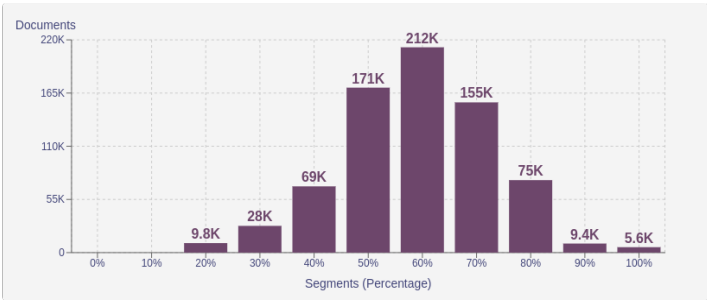


Language Distribution

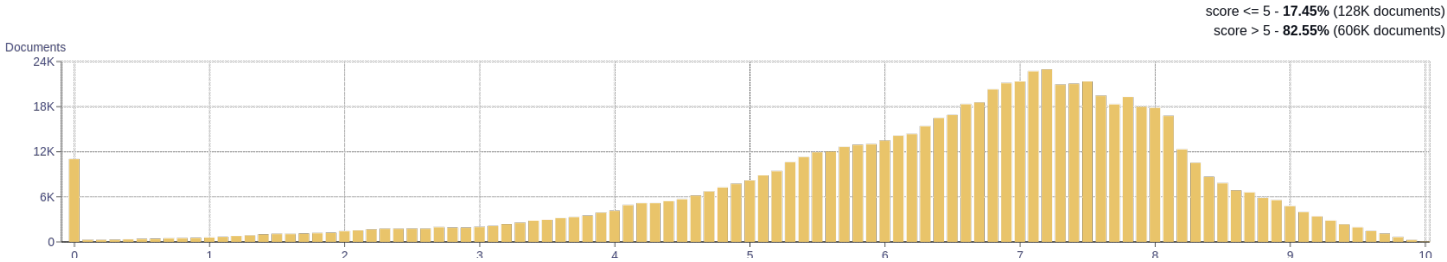
Number of segments



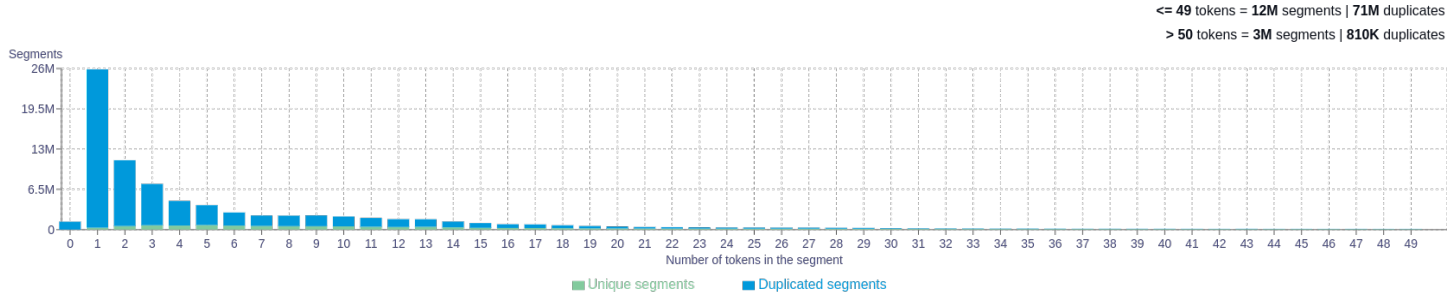
Percentage of segments in Macedonian (mk) inside documents



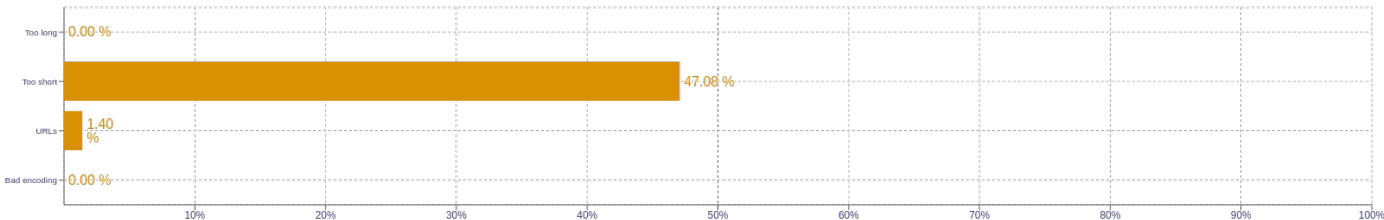
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>македонија 2546789</div> <div>вести 1408817</div> <div>година 1313749</div> <div>скопје 1187774</div> <div>видео 1138123</div>
2	<div>република македонија 237582</div> <div>најнови вести 171533</div> <div>of the 142448</div> <div>милиони евра 124160</div> <div>read more 122921</div>
3	<div>права се задржани 127603</div> <div>услови за користење 120153</div> <div>љубов и секс 96559</div> <div>all rights reserved 83602</div> <div>skip to content 68997</div>
4	<div>cookie is set by 47536</div> <div> 40187</div> <div>the user consent for 40018</div> <div>user consent for the 39935</div> <div>consent for the cookies 39935</div>
5	<div>user consent for the cookies 39935</div> <div>the user consent for the 39935</div> <div>for the cookies in the 39933</div> <div>consent for the cookies in 39933</div> <div>the cookies in the category 38982</div>

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.sitinfo.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>