

General overview

Corpus	Date	Language
hplt-v3-mkd_Cyrl	9/24/2025	Macedonian (mk)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
6,791,451	97,567,523	58,291,826 (59.75 %)	3B	16,333,470,805	27.26 GB

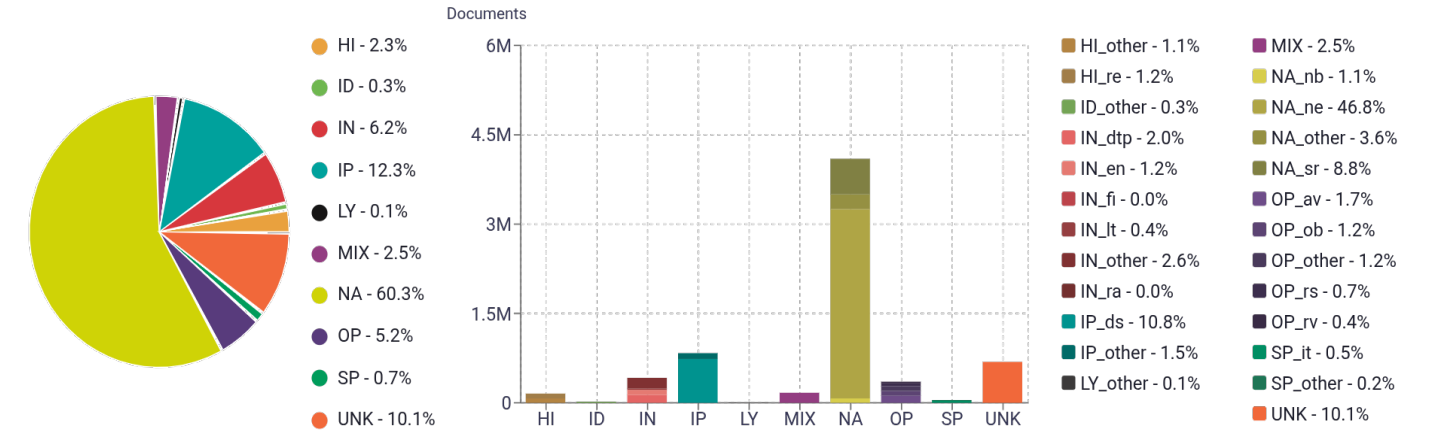
Top 10 domains

Domain	Docs	% of total
daily.mk	185K	2.72%
makfax.com.mk	110K	1.62%
slobodnaevropa.mk	93K	1.37%
netpress.com.mk	88K	1.29%
novamakedonija....	86K	1.26%
airbnb.com	72K	1.06%
republika.mk	71K	1.04%
kurir.mk	70K	1.03%
wikipedia.org	69K	1.02%
mkd.mk	59K	0.86%

Top 10 TLDs

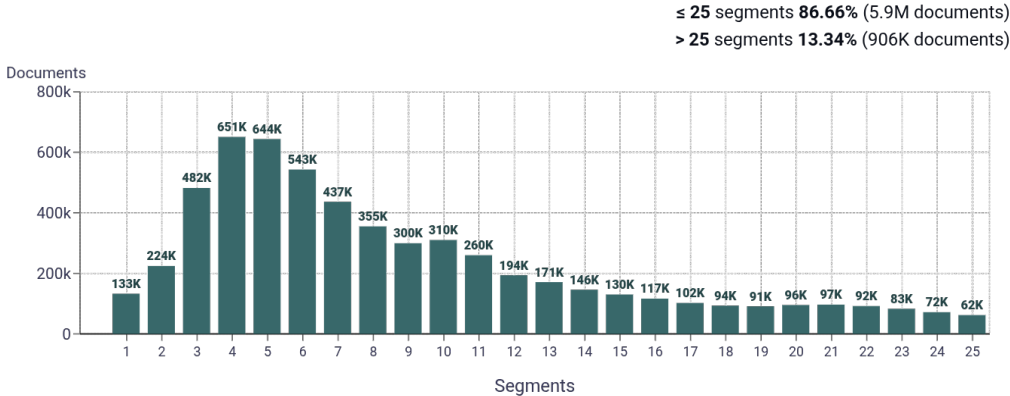
Domain	Docs	% of total
mk	3.9M	57.16%
com	1.1M	16.11%
com.mk	778K	11.46%
org	414K	6.09%
org.mk	108K	1.60%
gov.mk	90K	1.33%
net	49K	0.71%
edu.mk	47K	0.69%
co.uk	39K	0.57%
net.tr	34K	0.50%

Register labels

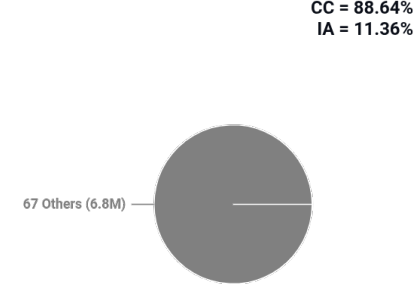


MT:8.9% | 605K Documents

Documents size (in segments) ⓘ

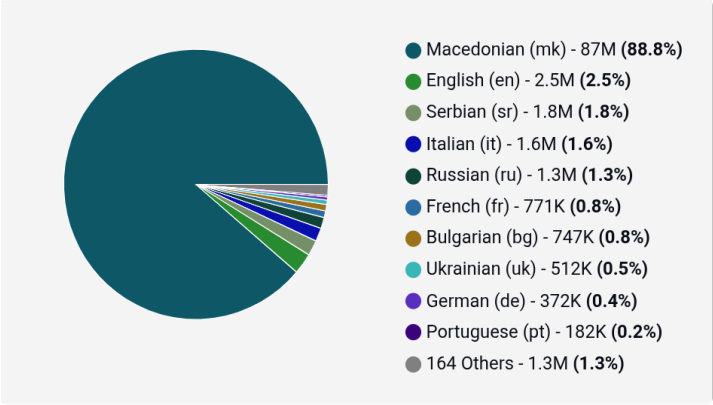


Document collections

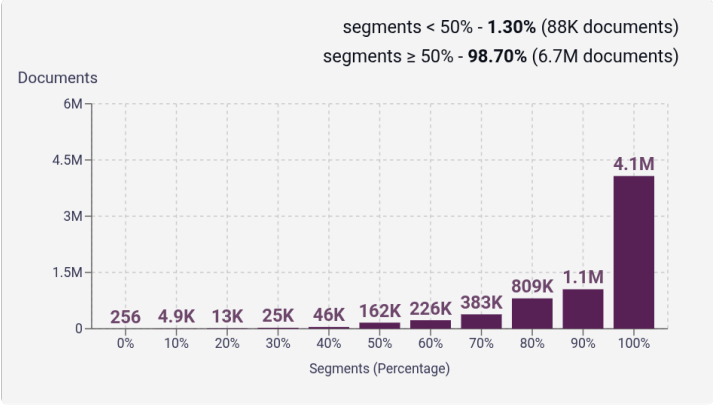


Language Distribution

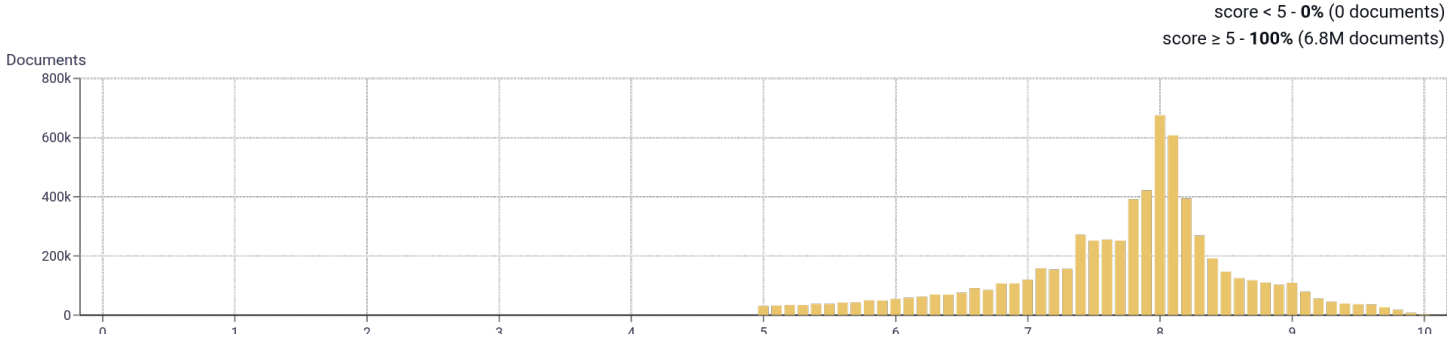
Number of segments in the Macedonian (mk) corpus



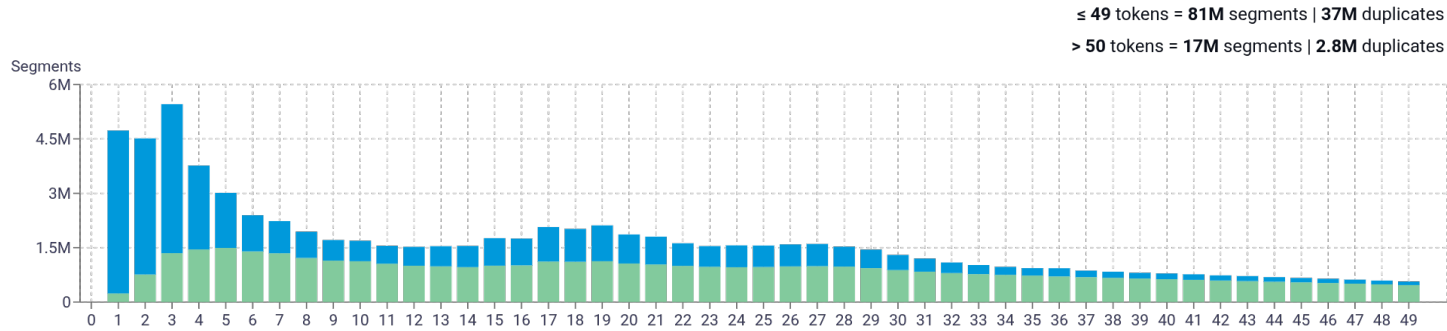
Percentage of segments in Macedonian (mk) inside documents



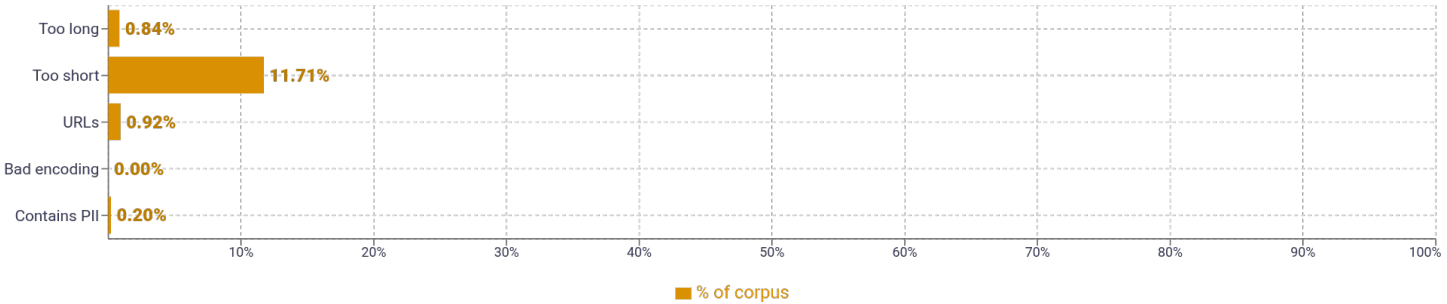
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	година 6,089,057 македонија 4,121,460 време 3,603,079 години 3,008,923 дел 2,789,299	
2	ве молиме 815,878 република македонија 644,988 станува збор 524,223 северна македонија 499,523 милиони евра 437,952	
3	република северна македонија 221,843 лигата на шампионите 136,497 можат да бидат 134,569 ве молиме контактирајте 131,144 министерството за здравство 115,663	
4	министерството за внатрешни работи 88,892 министерот за надворешни работи 71,335 министер за надворешни работи 66,752 are closed for this 66,507 comments are closed for 66,504	
5	comments are closed for this 66,504 are closed for this post 66,486 изнајмување на сместувања за одмор 50,024 view this post on instagram 40,348 сноси никаква одговорност за коментарите 37,881	

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>

Register labels

Name	Abbr.	Name	Abbr.	Name	Abbr.
Machine-translated	MT	How-to or instructions	HI	Description of a thing or person	dtp
Lyrical	LY	Recipe	re	FAQ	fi
Spoken	SP	Informational persuasion	IP	Legal terms & conditions	lt
Interview	it	Description with intent to sell	ds	Opinion	OP
Interactive discussion	ID	News & opinion blog or editorial	ed	Review	rv
Narrative	NA	Informational description	IN	Opinion blog	ob
News report	ne	Enciclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				