

General overview

Corpus	Date	Language
hplt-v3-bul_Cyrl	9/18/2025	Bulgarian (bg)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
42,967,942	978,252,958	539,310,730 (55.13 %)	27B	144,827,849,190	240.67 GB

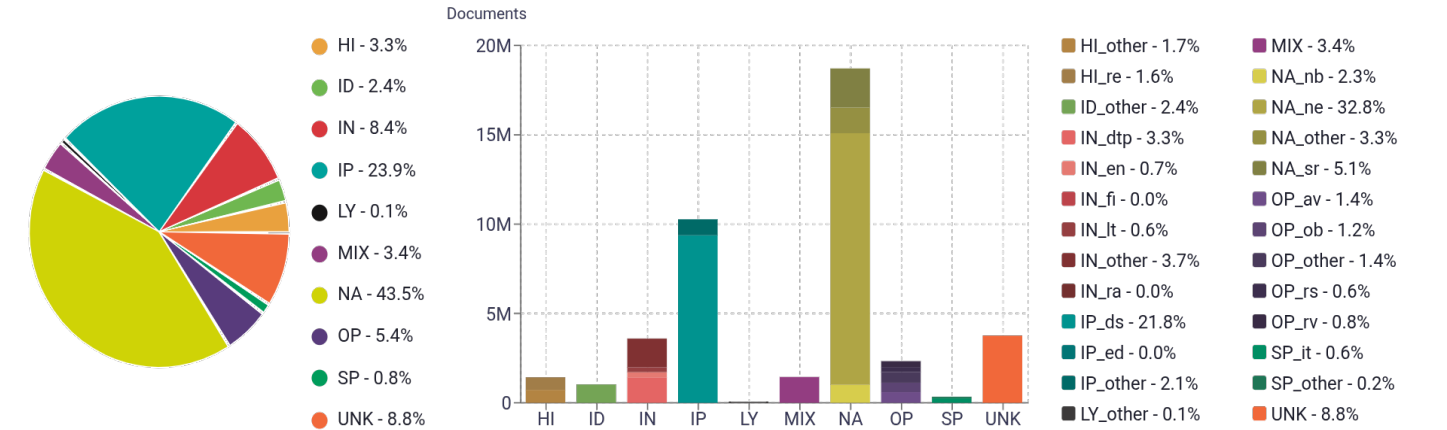
Top 10 domains

Domain	Docs	% of total
blogspot.com	391K	0.91%
offnews.bg	315K	0.73%
capital.bg	307K	0.71%
blog.bg	297K	0.69%
dnevnik.bg	286K	0.67%
dir.bg	277K	0.65%
blitz.bg	260K	0.61%
mediapool.bg	249K	0.58%
fakti.bg	228K	0.53%
sportal.bg	224K	0.52%

Top 10 TLDs

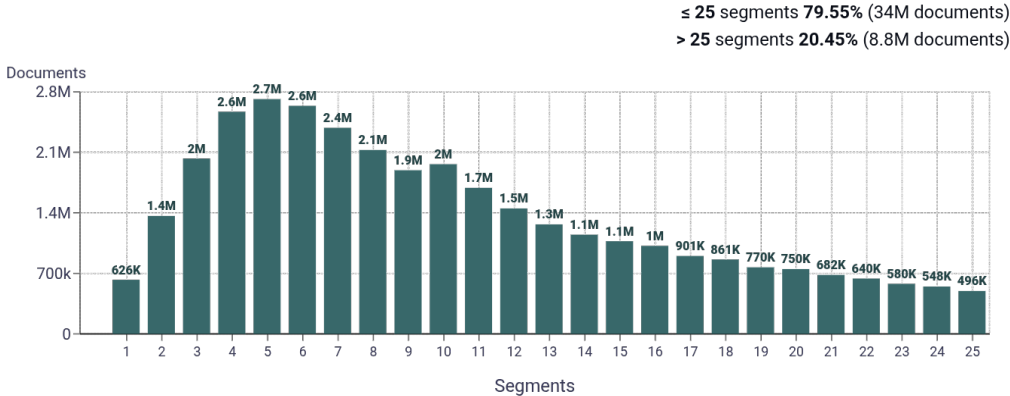
Domain	Docs	% of total
bg	20M	46.18%
com	15M	33.77%
net	2.1M	4.91%
org	2.1M	4.78%
eu	1.3M	2.94%
info	1.1M	2.55%
co.uk	186K	0.43%
news	185K	0.43%
ru	125K	0.29%
pt	115K	0.27%

Register labels

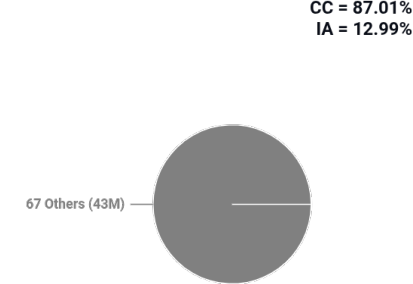


MT:5.7% | 2.4M Documents

Documents size (in segments) ⓘ

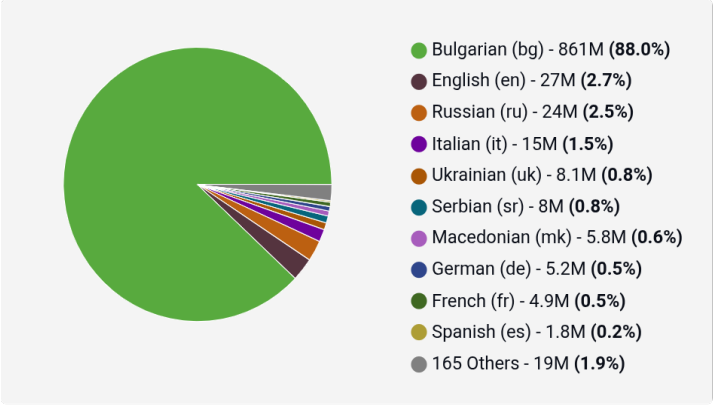


Document collections

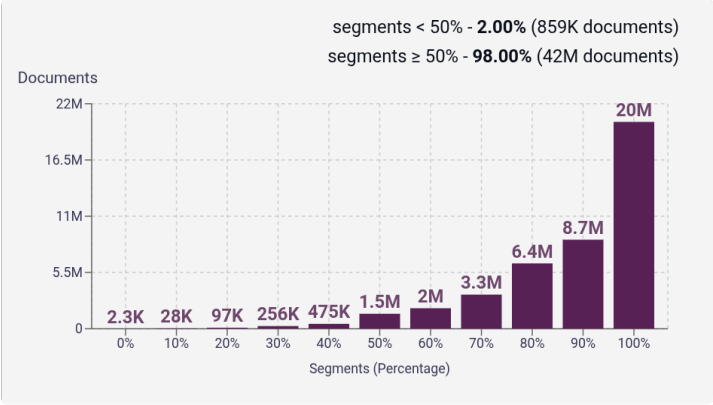


Language Distribution

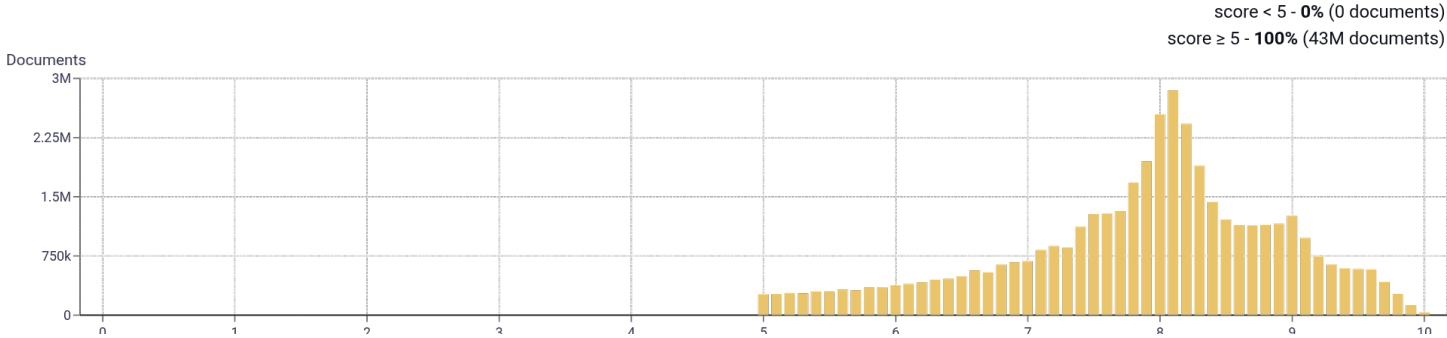
Number of segments in the Bulgarian (bg) corpus



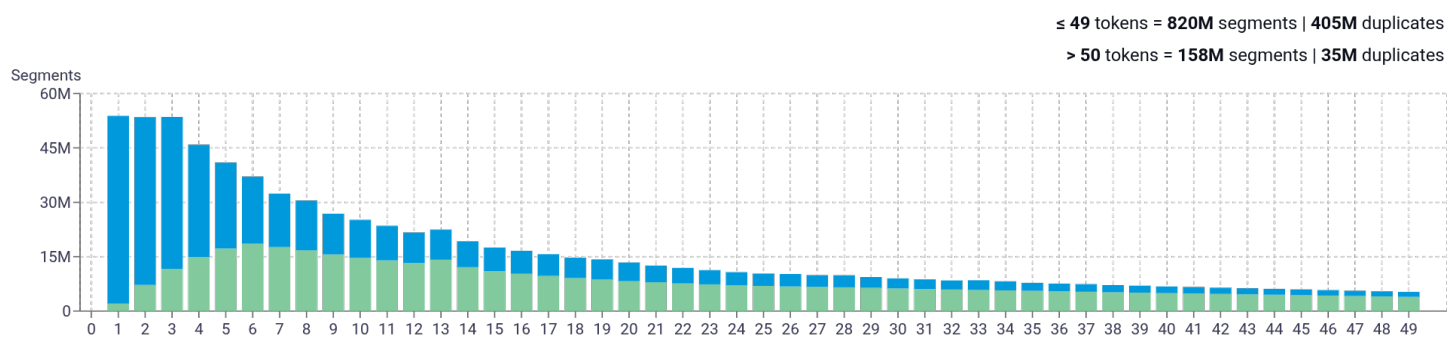
Percentage of segments in Bulgarian (bg) inside documents



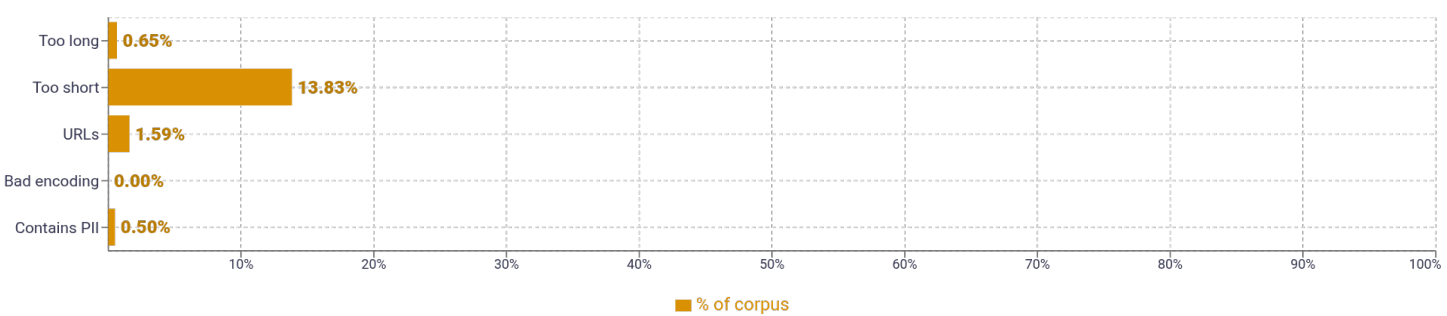
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	г. 38,753,143 България 22,182,363 част 19,052,047 хора 16,800,290 лв 15,984,244	
2	т. н. 2,727,646 т. е. 2,048,043 крайна сметка 1,741,286 стара загора 1,724,884 неутралнодо коментар 1,704,238	
3	що се отнася 778,655 течение на времето 487,531 загуба на тегло 466,811 коментарът беше изтрит 417,839 изтрит от модераторите 406,146	
4	етническа или верска основа 389,987 верска основа или призови 276,760 адрес на конкретни лица 267,420 регистрирайте се от тук 261,405 обидни или нецензурни квалификации 258,654	
5	коментарът беше изтрит от модераторите 405,817 призови към насилие по адрес 276,792 основа или призови към насилие 276,758 насилие по адрес на конкретни 267,963 съдържаше обидни или нецензурни квалификации 258,582	

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	ntp
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				