

General overview

Corpus	Date	Language
hplt-v3-slv_Latn	9/18/2025	Slovenian

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
16,812,095	402,087,512	226,735,448 (56.39 %)	11B	62,072,495,644	59.29 GB

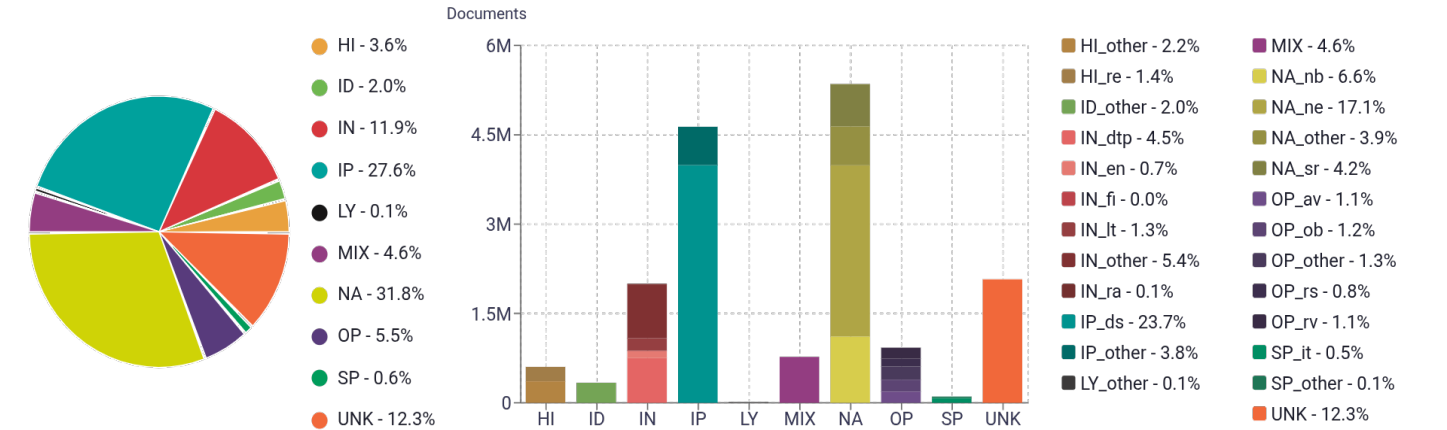
Top 10 domains

Domain	Docs	% of total
delo.si	264K	1.57%
rtvslo.si	227K	1.35%
sta.si	194K	1.16%
metropolitan.si	181K	1.08%
dnevnik.si	176K	1.05%
zurnal24.si	148K	0.88%
siol.net	148K	0.88%
blogspot.com	135K	0.80%
arnes.si	107K	0.63%
ognjisce.si	97K	0.58%

Top 10 TLDs

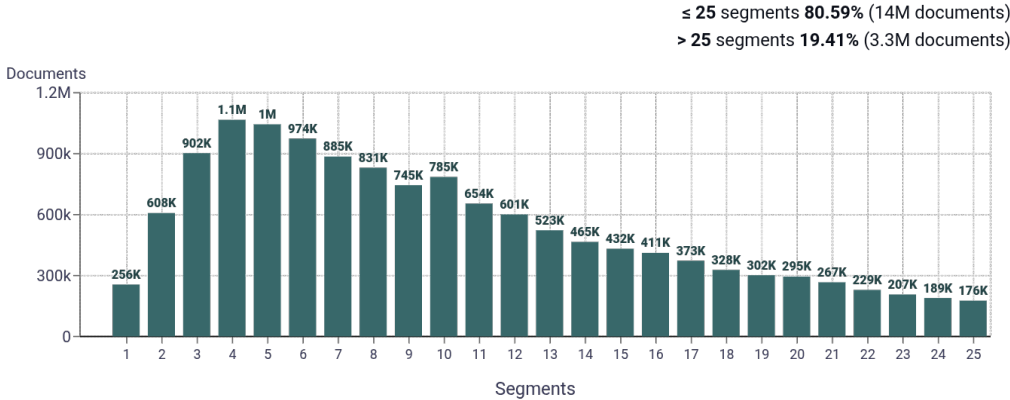
Domain	Docs	% of total
si	10M	60.50%
com	3.8M	22.60%
net	776K	4.62%
org	537K	3.20%
eu	407K	2.42%
info	188K	1.12%
cz	126K	0.75%
tv	50K	0.30%
pl	45K	0.27%
sk	41K	0.24%

Register labels

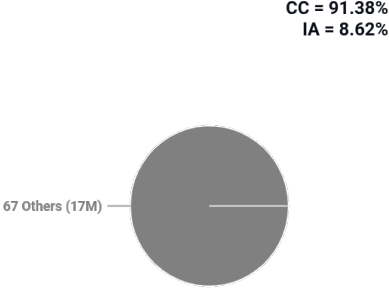


MT:9.0% | 1.5M Documents

Documents size (in segments) ⓘ

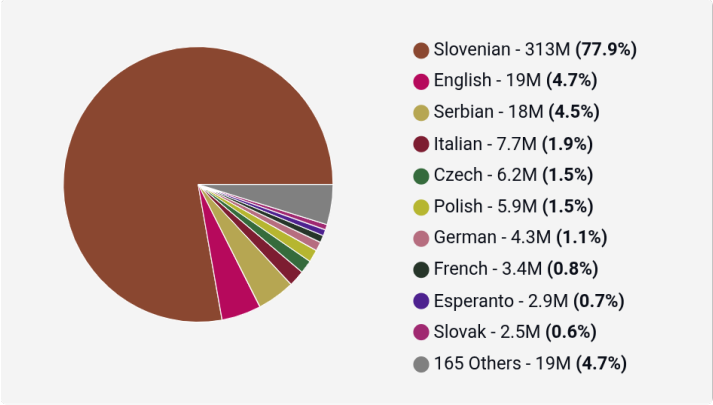


Document collections

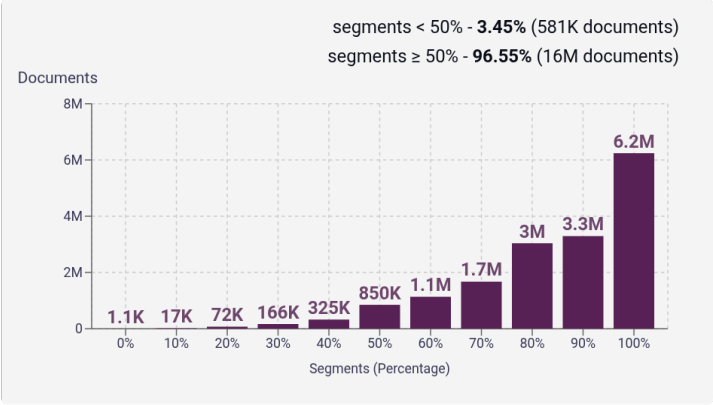


Language Distribution

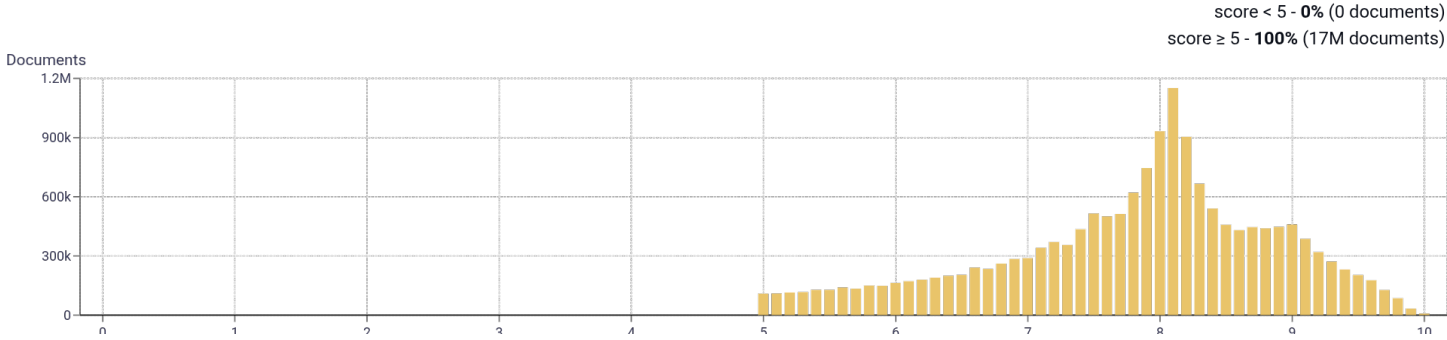
Number of segments in the Slovenian corpus



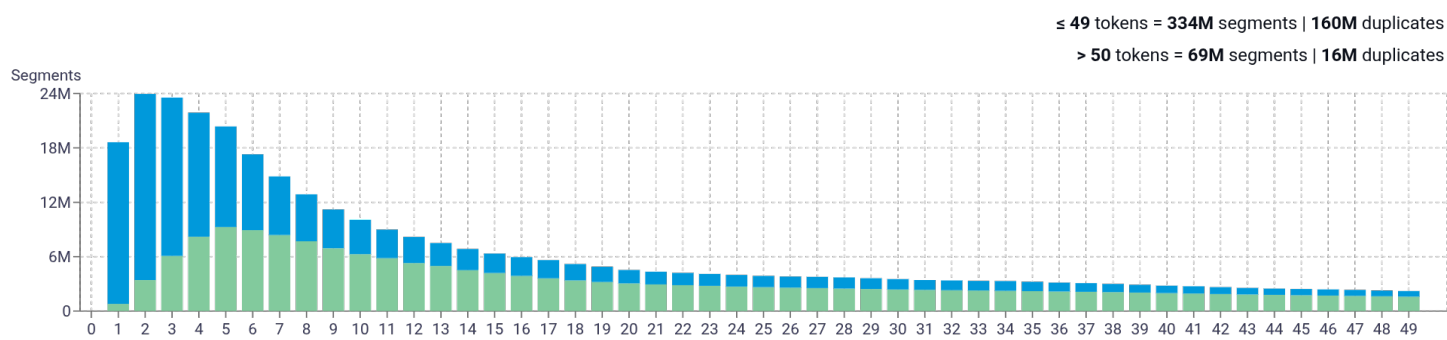
Percentage of segments in Slovenian inside documents



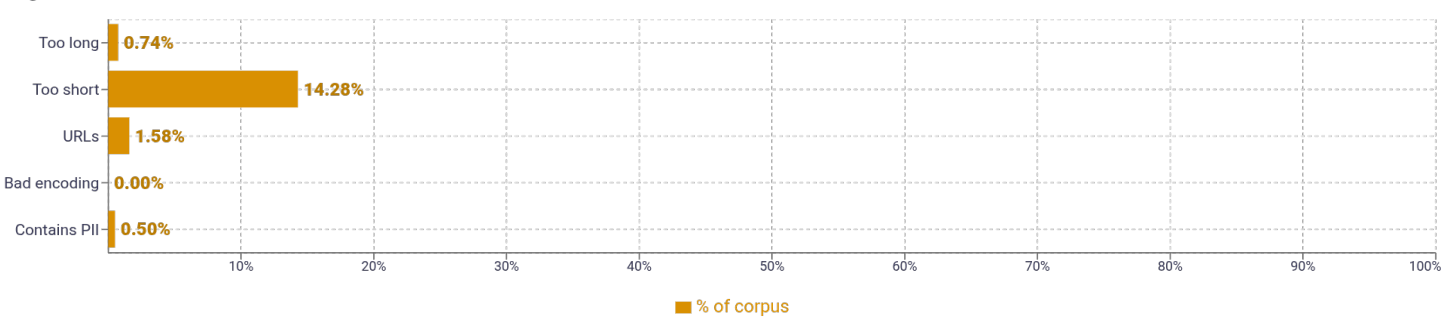
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	lahko 53,169,347 zelo 12,235,899 zato 11,475,857 ima 10,779,122 a 10,310,684	
2	spletni strani 1,254,060 igralni avtomat 1,026,230 free play 883,423 igranj igralni 843,597 play demo 836,406	
3	igranj igralni avtomat 841,336 free play demo 836,392 igre na srečo 494,237 no deposit bonus 442,458 iger na srečo 439,968	
4	free spins no deposit 238,033 spins no deposit bonus 237,314 casino free spins no 235,692 casino no deposit bonus 194,863 bonus and promo code 192,492	
5	free spins no deposit bonus 237,309 casino free spins no deposit 235,690 casino bonus and promo code 190,471 odgovoren za javno spodbujanje sovraštva 95,872 kazensko odgovoren za javno spodbujanje 95,239	

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 $\frac{\text{number of types (uniques)}}{\text{number of tokens}}$, 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				