

General overview

Corpus	Date	SL	TL
hplt-v2-en-sr.tsv	1/25/2025	English (en)	Serbian (sr)

Volumes

Segments	SL tokens	SL characters	SL size
5,291,686	118M	607,110,021	581.59 MB

TL tokens	TL characters	TL size
105M	576,427,030	991.74 MB

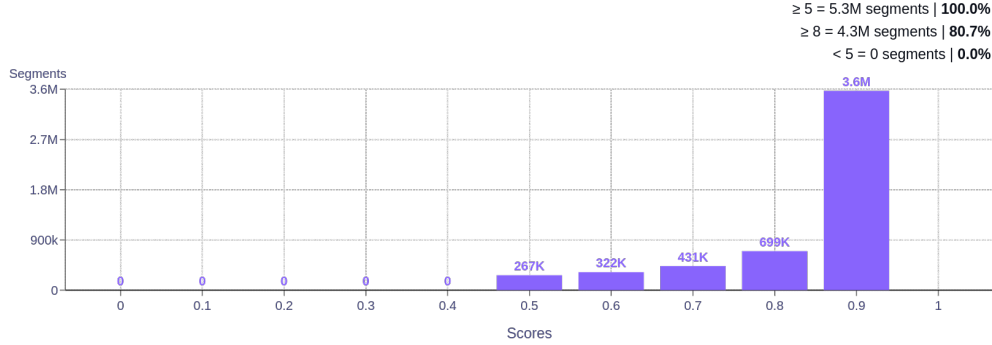
Dataset top 10 domains

SL domain	Segments	TL domain	Segments
wikipedia.org	15.9%	wikipedia.org	14.0%
google.com	12.0%	google.com	5.6%
biblegateway.com	1.9%	tripadvisor.rs	4.5%
erieaquariumsociety.com	1.9%	rbth.com	1.6%
rbth.com	1.7%	biblegateway.com	1.6%
jw.org	1.5%	erieaquariumsociety.com	1.5%
academiccourses.com	1.4%	jw.org	1.4%
bachelorstudies.com	1.3%	stealthsettings.com	1.2%
stealthsettings.com	1.3%	vsaduidoma.com	1.2%
educationbro.com	1.2%	stroifaq.com	0.9%

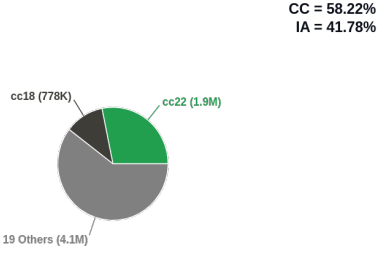
Dataset top 10 TLDs

SL domain	Segments	TL domain	Segments
com	88.2%	com	61.5%
org	30.3%	org	25.3%
net	6.8%	rs	14.2%
rs	6.4%	net	5.8%
co.uk	2.1%	gov.rs	2.0%
gov.rs	2.0%	info	1.2%
eu	1.4%	ru	1.1%
ru	1.2%	eu	1.0%
info	1.2%	nu	0.9%
de	1.0%	de	0.8%

Translation likelihood

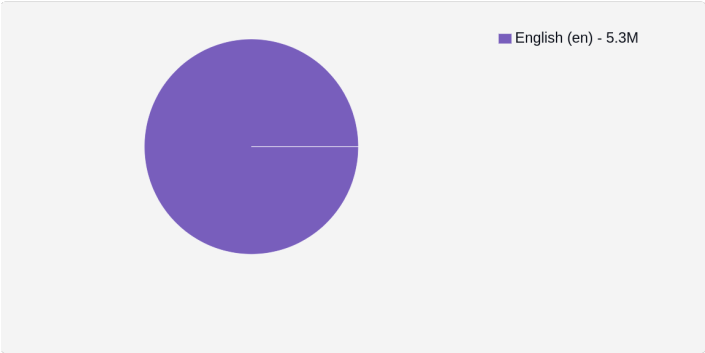


Collections

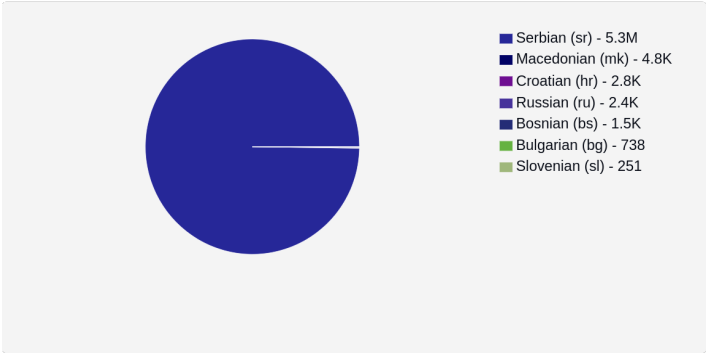


Language Distribution

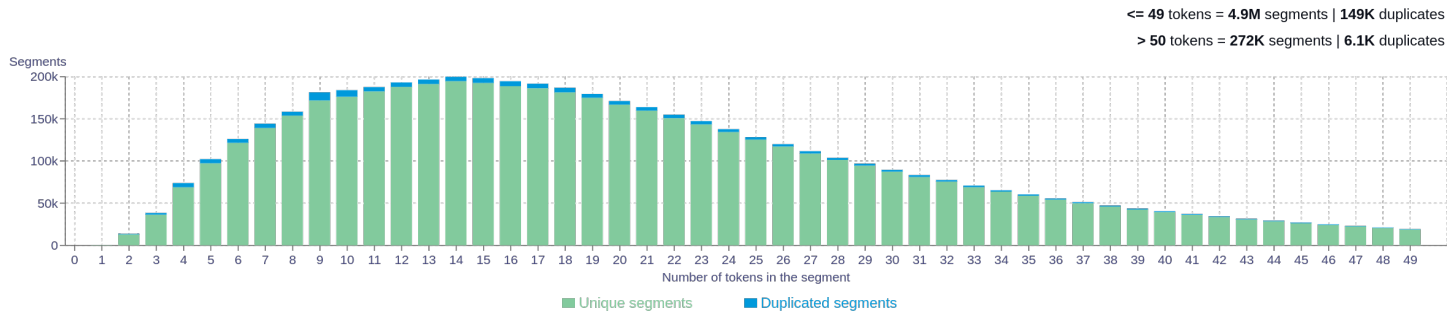
Source



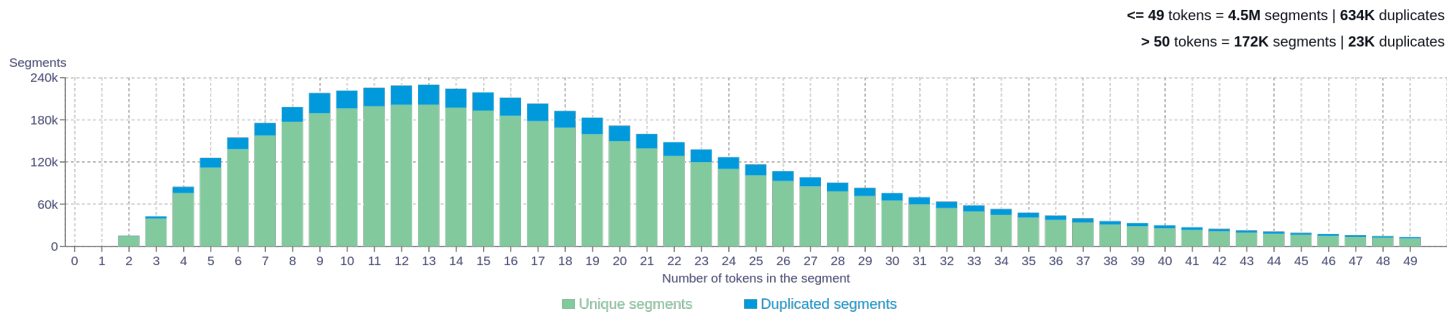
Target



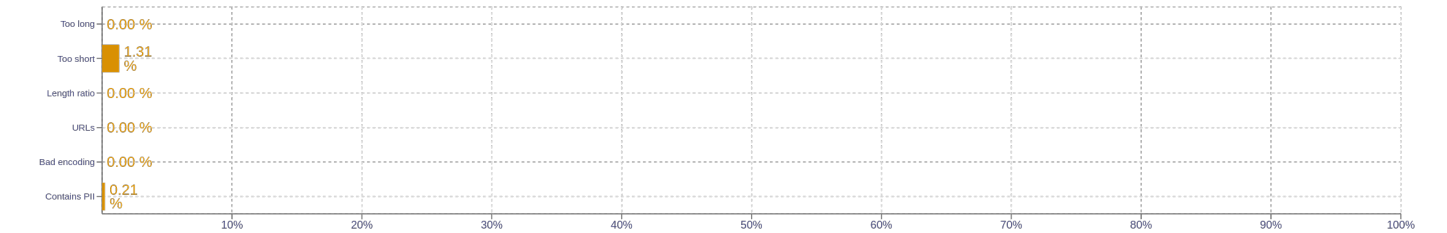
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	one   237634also   228866time   166318new   150877use   147535
2	personal data   22441united states   21667national assembly   13893online game   10540make sure   9909
3	republic of serbia   22280bosnia and herzegovina   10978republic of srpska   10445around the world   6076copy the code   4519
4	president of the republic   5168one of the best   4638assembly of the republic   4620code of your site   4496paste in the html   4494
5	national assembly of the republic   4548html code of your site   4495copy the code and paste   4495paste in the html code   4494email address is being protected   4018

Target n-grams

Size	n-grams
1	који   517003као   449279које   251468може   207248године   203753
2	може бити   33566пре него   19927републике србије   17519изговорио корисник   12639због тога   11708
3	босне и херцеговине   5742имајте на уму   5485још много тога   5081сједињене америчке државе   5071који је био   4989
4	копирајте код и налепите   4499хтмл код вашег сајта   4497налепите у хтмл код   4497висок квалитет мобител телефон   3773народне скупштине републике србије   2932
5	налепите у хтмл код вашег   4497код и налепите у хтмл   4497поште је заштићена од спамботова   3720заборавите да оцени ову игру   2012ову игру са својим најбољим   1819

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 or 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>