

General overview

Corpus	Date	SL	TL
hplt-v2-en-et.tsv	1/26/2025	English (en)	Estonian (et)

Volumes

Segments	SL tokens	SL characters	SL size
8,797,574	207M	1,090,375,785	1.02 GB

TL tokens	TL characters	TL size
158M	1,043,287,714	1.0 GB

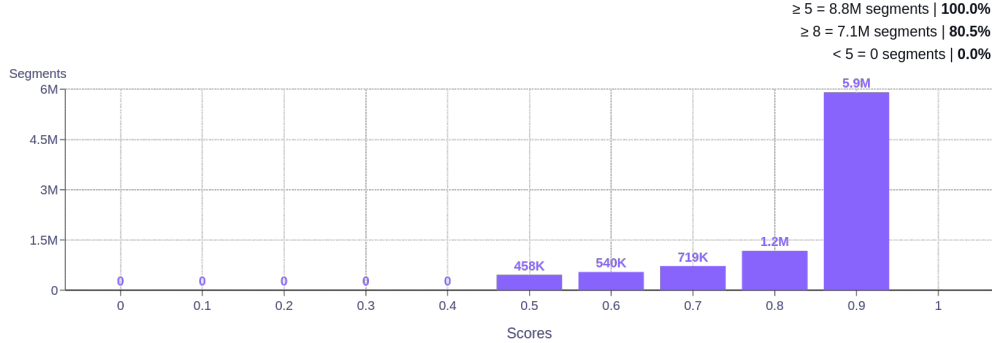
Dataset top 10 domains

SL domain	Segments	TL domain	Segments
europa.eu	15.5%	europa.eu	12.5%
hotels.com	12.5%	hotels.com	6.0%
google.com	6.7%	wikipedia.org	3.8%
agoda.com	4.9%	agoda.com	3.6%
wikipedia.org	4.5%	google.com	3.2%
booking.com	3.4%	riigiteataja.ee	3.0%
microsoft.com	3.2%	microsoft.com	2.1%
riigiteataja.ee	2.9%	booking.com	2.1%
office.com	2.2%	office.com	2.0%
err.ee	1.0%	err.ee	1.1%

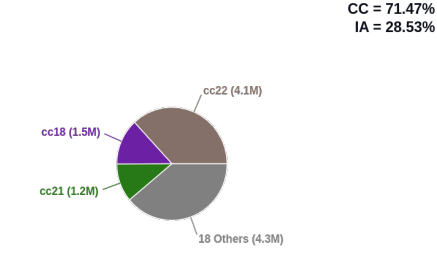
Dataset top 10 TLDs

SL domain	Segments	TL domain	Segments
com	87.4%	com	56.0%
ee	21.9%	ee	34.6%
eu	19.7%	eu	16.0%
org	10.8%	org	9.0%
net	3.5%	net	2.9%
info	2.4%	info	2.3%
co.uk	2.4%	fi	0.9%
de	1.2%	de	0.5%
fi	1.1%	ru	0.4%
ie	1.0%	lv	0.4%

Translation likelihood

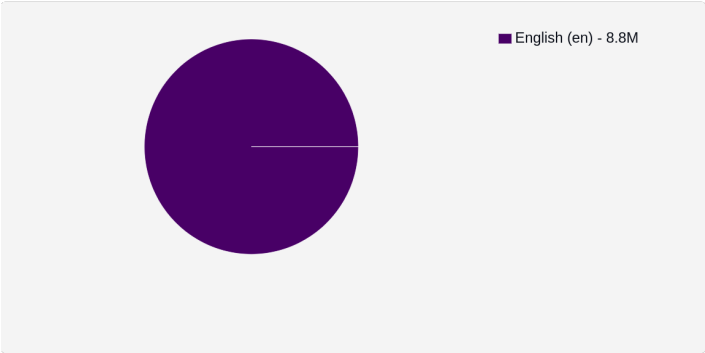


Collections

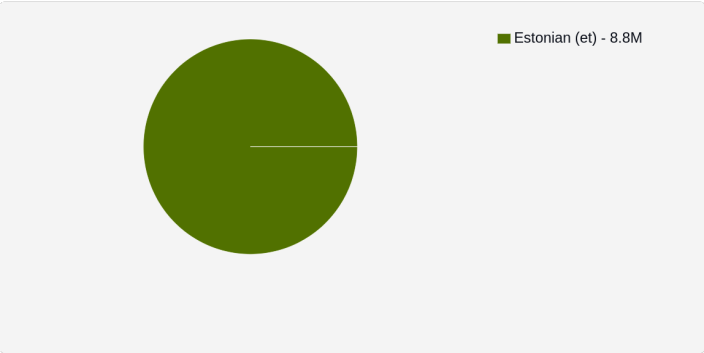


Language Distribution

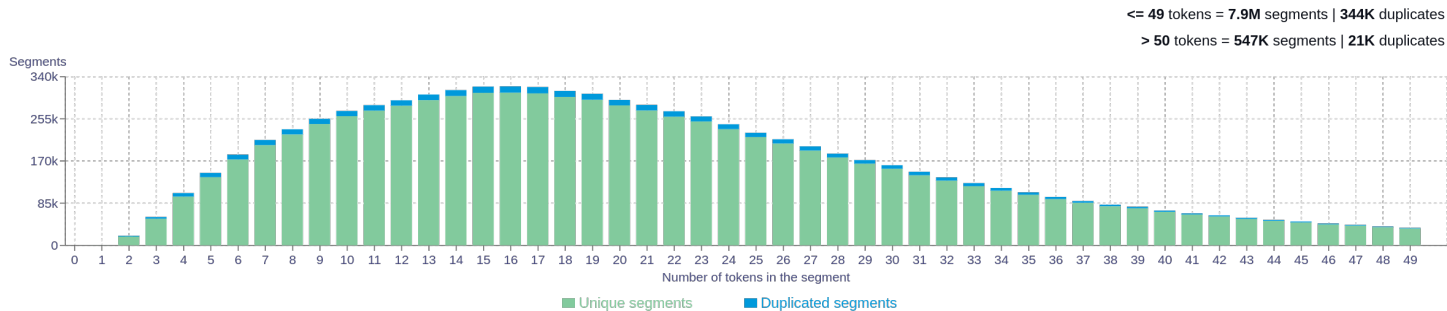
Source



Target



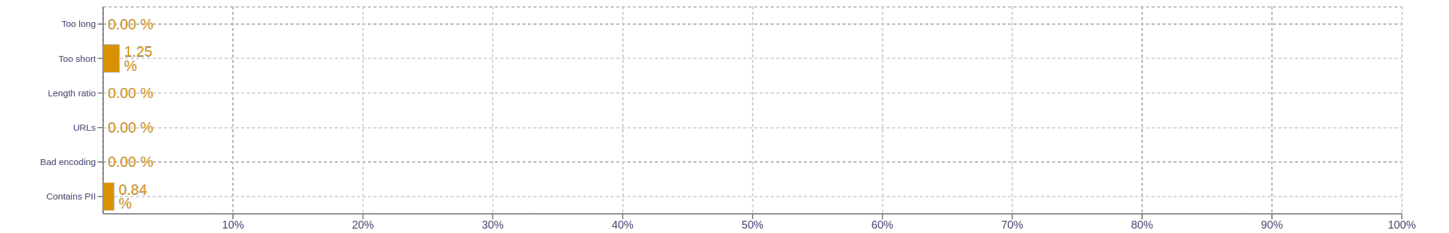
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	data 500981also 401557use 362407one 327699information 323378
2	personal data 209383member states 65018european union 45327member state 44082data protection 35424
3	processing of personal 23251terms and conditions 16851wi-fi in public 15512choice for travelers 11954right to object 10908
4	processing of personal data 23037processing of your personal 17271wi-fi in all rooms 16999wi-fi in public areas 15501referred to in article 13836
5	parliament and of the council 19222free wi-fi in all rooms 16988processing of your personal data 16146people looked at this hotel 7311hotel in the last hour 7311

Target n-grams

Size	n-grams
1	või 1169543ning 741931ka 567920teie 432300meie 368748
2	euroopa liidu 45900võib olla 40923teie isikuandmeid 35941igal ajal 32263ettevõttega ühendust 31607
3	parlamendi ja nõukogu 20766teil on õigus 17484wifi-ühendus kõigis tubades 17204tasuta wifi-ühendus kõigis 17204iganes on sinu 10752
4	euroopa parlamendi ja nõukogu 20639tasuta wifi-ühendus kõigis tubades 17204kuidas sinna kohale jõuda 8168vaadata vabade kohtade olemasolu 7717kohe vaadata vabade kohtade 7717
5	kohe vaadata vabade kohtade olemasolu 7717hindu ja kohe vaadata vabade 7400vaadanud seda hotelli viimase tunni 6623külastajat on vaadanud seda hotelli 6623hotels.com abiga on lihtne leida 6527

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>