

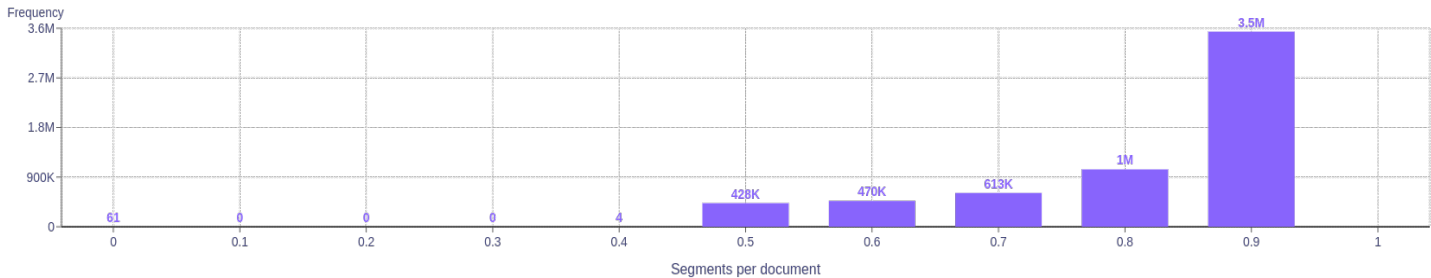
General overview

Corpus	Analytics date	Source language	Target language
HPLT.en-et	10/25/2023	English (en)	Estonian (et)

Volumes

Segments	Unique segments	Src tokens	Trg tokens	Src size	Trg size	Src characters	Trg characters
6,089,852	6,089,792 (100.00 %)	111M	90M	564.07 MB	579.75 MB		

Translation likelihood

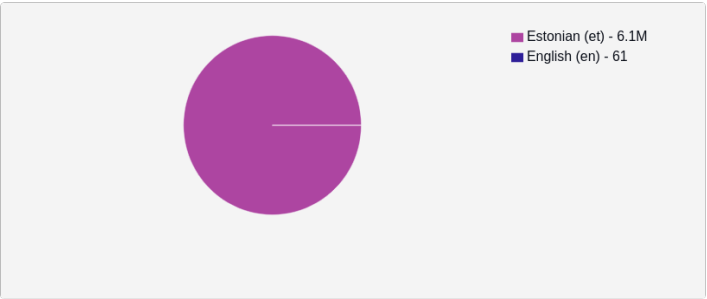


Language Distribution

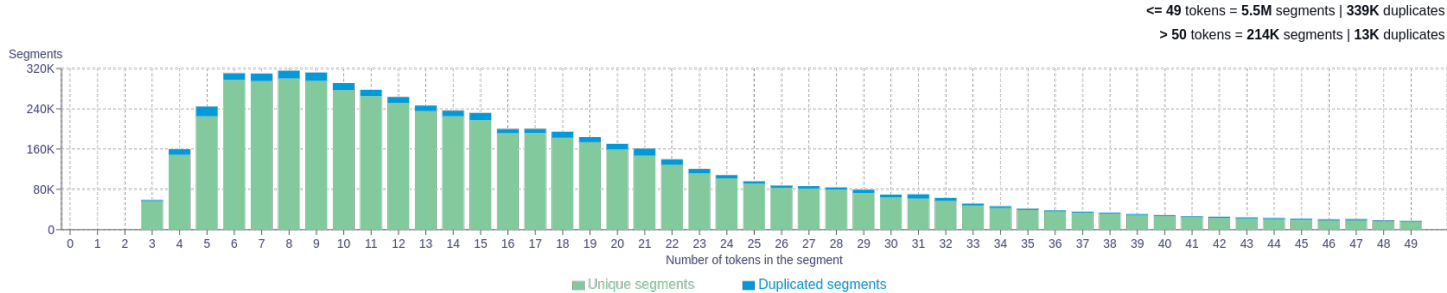
Source



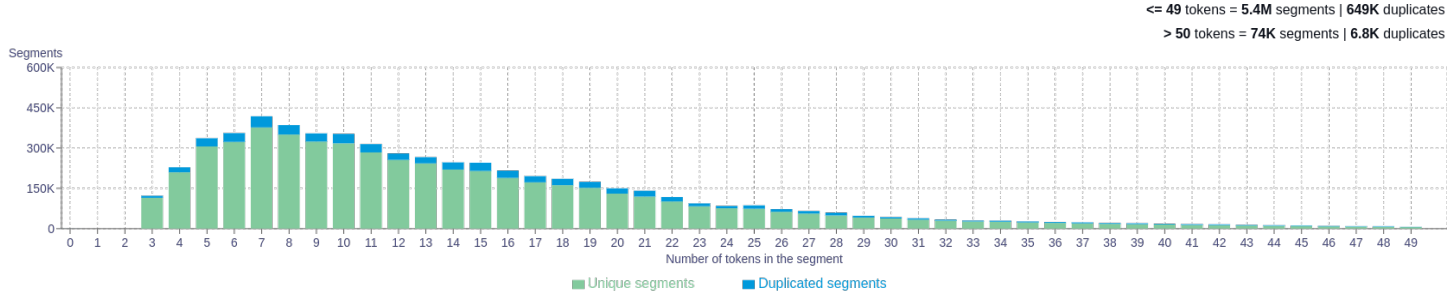
Target



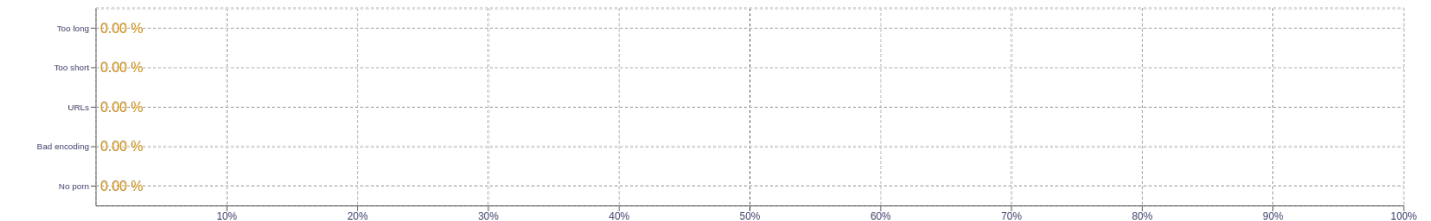
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	<div>best 319870</div> <div>airport 299809</div> <div>books 285365</div> <div>hotel 276390</div> <div>car 275518</div>
2	<div>car hire 118255</div> <div>best prices 103370</div> <div>best price 88195</div> <div>second hand 87976</div> <div>rare books 87927</div>
3	<div>year of manufacture 142172</div> <div>second hand books 87908</div> <div>books and second 87908</div> <div>available rare books 87908</div> <div>amend your booking 63983</div>
4	<div>used books and second 87908</div> <div>books of the title 87908</div> <div>books and second hand 87908</div> <div>find you the best 58984</div> <div>get the best price 58956</div>
5	<div>used books and second hand 87908</div> <div>hand books of the title 87908</div> <div>books and second hand books 87908</div> <div>amend your booking for free 63978</div> <div>rentalcars.com and you can amend 63977</div>

Target n-grams

Size	n-grams
1	<div>või 503390</div> <div>ning 343835</div> <div>kasutatud 234275</div> <div>asukohas 217472</div> <div>tasuta 212514</div>
2	<div>kasutatud raamatute 176436</div> <div>haruldased raamatud 88219</div> <div>saadaval haruldased 88218</div> <div>raamatute pealkiri 88218</div> <div>täielikult loetletud 85807</div>
3	<div>saadaval haruldased raamatud 88218</div> <div>raamatute ja kasutatud 88218</div> <div>kasutatud raamatute pealkiri 88218</div> <div>saate oma broneeringut 63981</div> <div>broneeringut tasuta muuta 63979</div>
4	<div>raamatute ja kasutatud raamatute 88218</div> <div>kasutatud raamatute ja kasutatud 88218</div> <div>saate oma broneeringut tasuta 63980</div> <div>kaudu ja te saate 63880</div> <div>leida teile parimad hinnapakkumised 58944</div>
5	<div>raamatute ja kasutatud raamatute pealkiri 88218</div> <div>kasutatud raamatute ja kasutatud raamatute 88218</div> <div>saate oma broneeringut tasuta muuta 63979</div> <div>veebilehe kaudu ja te saate 63880</div> <div>meie juures ja me garanteerime 58844</div>

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>