

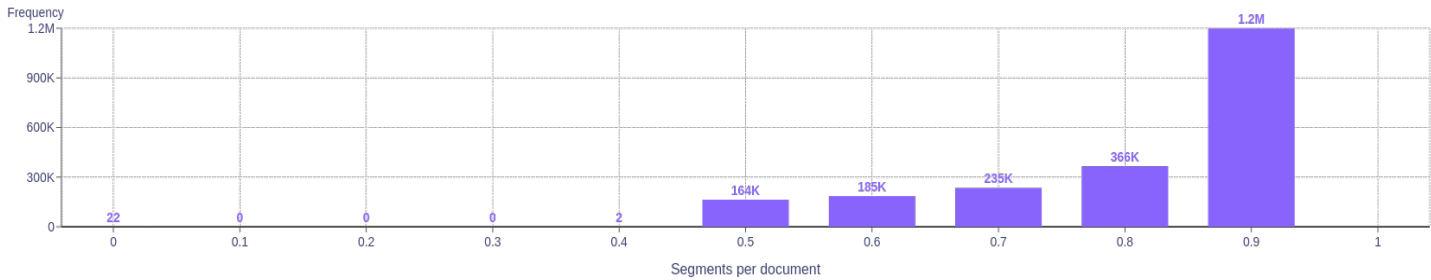
General overview

Corpus	Analytics date	Source language	Target language
HPLT.en-is	10/26/2023	English (en)	Icelandic (is)

Volumes

Segments	Unique segments	Src tokens	Trg tokens	Src size	Trg size	Src characters	Trg characters
2,148,876	2,148,855 (100.00 %)	33M	33M	167.14 MB	195.03 MB		

Translation likelihood

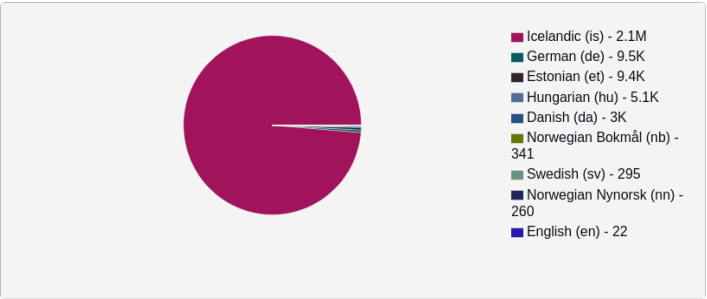


Language Distribution

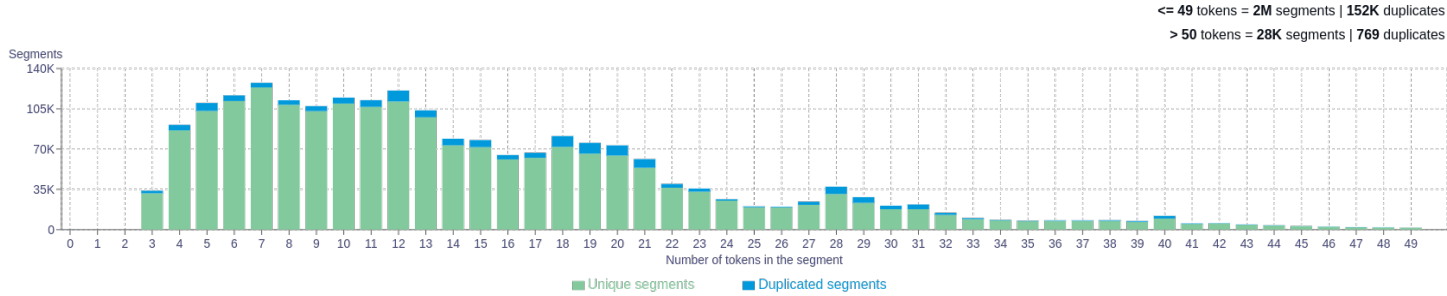
Source



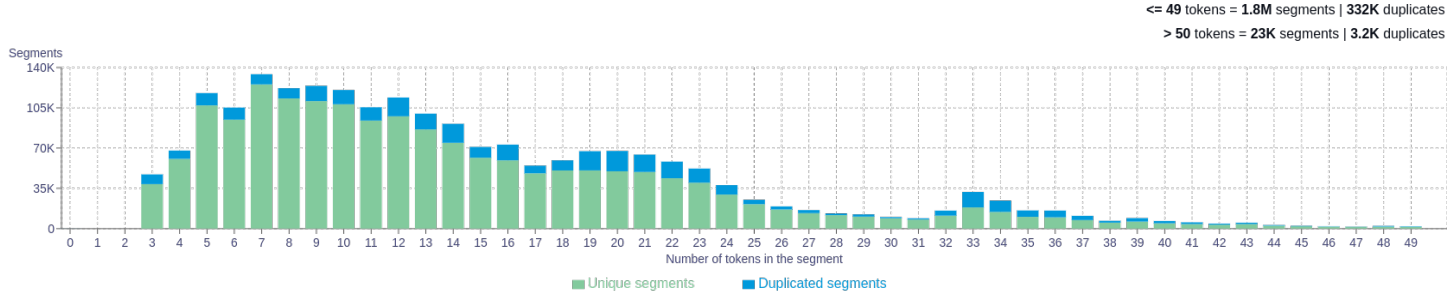
Target



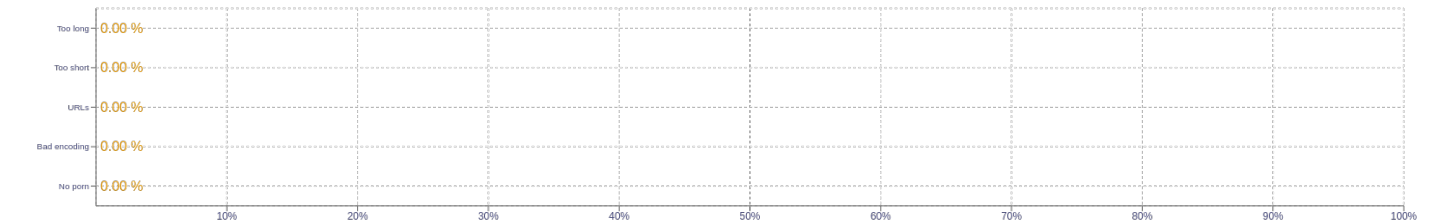
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	<div>car   263826</div> <div>airport   260383</div> <div>book   241892</div> <div>best   227192</div> <div>prices   198555</div>
2	<div>car hire   148479</div> <div>best prices   81098</div> <div>best price   78228</div> <div>find great   77495</div> <div>great prices   77359</div>
3	<div>quickly and easily   77351</div> <div>find great prices   77285</div> <div>see customer ratings   77284</div> <div>booking for free   76364</div> <div>amend your booking   76364</div>
4	<div>find you the best   75885</div> <div>get the best price   75857</div> <div>work hard to find   75854</div> <div>opens in new window   40326</div> <div>welcoming booking.com guests since   22128</div>
5	<div>amend your booking for free   76362</div> <div>rentalcars.com and you can amend   76361</div> <div>find you the best prices   75854</div> <div>book with us and get   75854</div> <div>rent a car car hire   19334</div>

Target n-grams

Size	n-grams
1	<div>bókaðu   235024</div> <div>bílaíleigubíl   177800</div> <div>airport   141736</div> <div>hotel   128558</div> <div>verð   119887</div>
2	<div>besta verðið   78427</div> <div>finndu fráðær   77517</div> <div>fráðær verð   77376</div> <div>getur breytt   77002</div> <div>breytt bókun   76871</div>
3	<div>bókaðu á netinu   77347</div> <div>finndu fráðær verð   77338</div> <div>verðið á bílaíleigubíl   76935</div> <div>getur breytt bókun   76870</div> <div>fáðu besta verðið   75875</div>
4	<div>besta verðið á bílaíleigubíl   76935</div> <div>rentalcars.com og þú getur   76870</div> <div>bókun þinni án endurgjalds   76870</div> <div>airport í gegnum rentalcars.com   62591</div> <div>opnast í nýjum glugga   40047</div>
5	<div>rentalcars.com og þú getur breytt   76870</div> <div>breytt bókun þinni án endurgjalds   76870</div> <div>fáðu besta verðið á bílaíleigubíl   75873</div> <div>tekið á móti gestum booking.com   22128</div> <div>verðs ef lítið er tilhlutfallsins   18199</div>

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 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>, <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>