

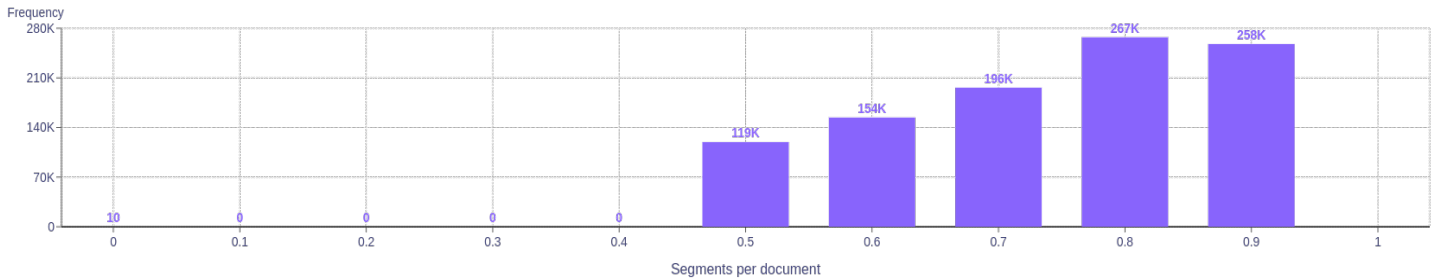
General overview

Corpus	Analytics date	Source language	Target language
HPLT.en-ga	10/23/2023	English (en)	Irish (ga)

Volumes

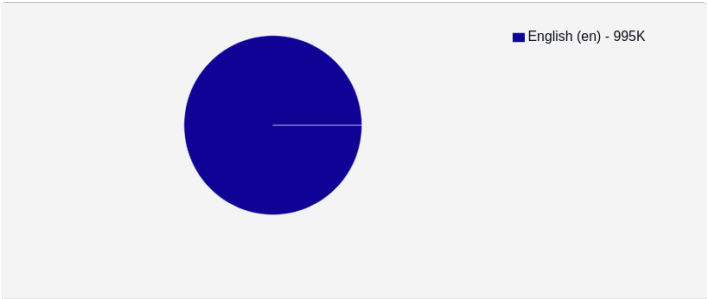
Segments	Unique segments	Src tokens	Trg tokens	Src size	Trg size	Src characters	Trg characters
994,756	994,747 (100.00 %)	18M	20M	94.79 MB	113.35 MB		

Translation likelihood

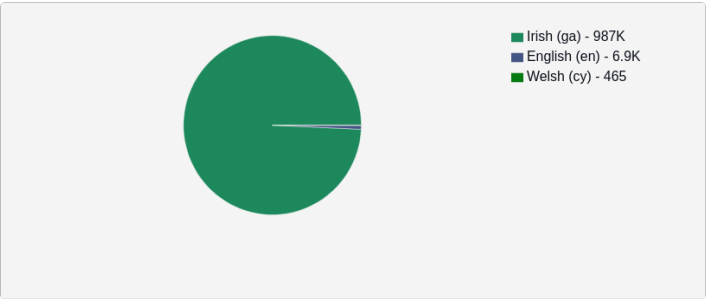


Language Distribution

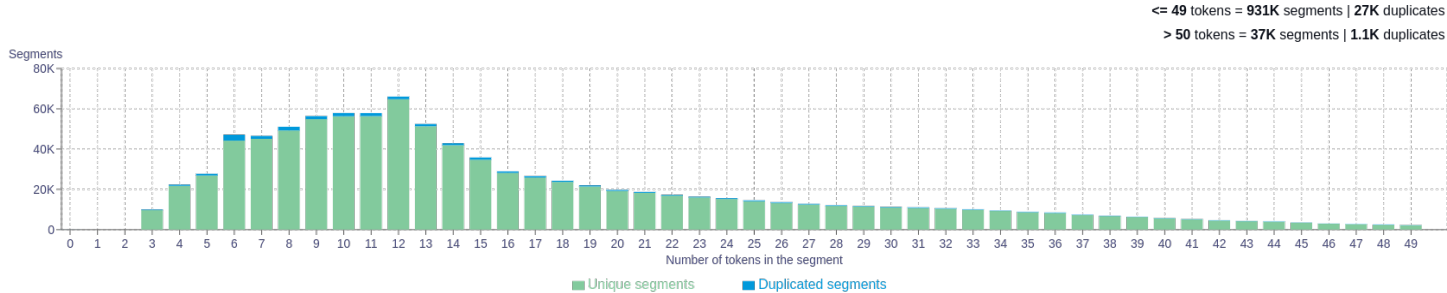
Source



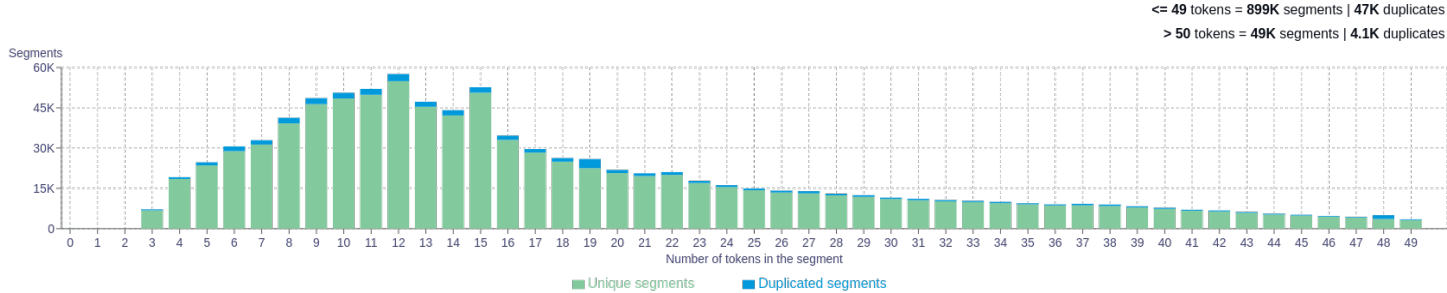
Target



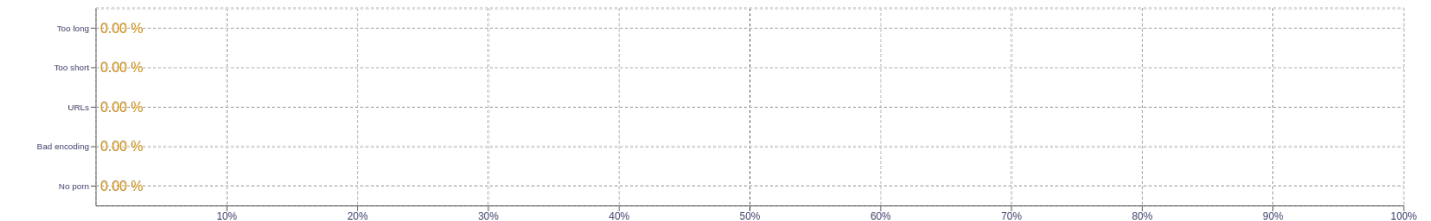
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	<div>porn 168845</div> <div>free 137838</div> <div>quality 119879</div> <div>hd 108388</div> <div>site 103134</div>
2	<div>excellent quality 38113</div> <div>hd excellent 36807</div> <div>good quality 35092</div> <div>quality hd 34752</div> <div>without registration 29869</div>
3	<div>hd excellent quality 36779</div> <div>video in good 29731</div> <div>good quality hd 22720</div> <div>free and without 17869</div> <div>hd great quality 12000</div>
4	<div>video in good quality 29727</div> <div>name of this movie 25821</div> <div>free and without registration 17454</div> <div>look free and without 16169</div> <div>eyes of the operator 11834</div>
5	<div>video in good quality hd 22412</div> <div>look free and without registration 16169</div> <div>video in high quality hd 11008</div> <div>video is in the categories 8618</div> <div>original name of this movie 8618</div>

Target n-grams

Size	n-grams
1	<div>porn 177398</div> <div>saor 137800</div> <div>aisce 135445</div> <div>hd 108810</div> <div>suíomh 105501</div>
2	<div>féidir leat 36056</div> <div>caighdeán maith 34370</div> <div>hd cáilíochta 26015</div> <div>maith hd 22398</div> <div>porn saor 22056</div>
3	<div>saor in aisce 134561</div> <div>fiseán i caighdeán 30133</div> <div>caighdeán maith hd 22396</div> <div>chaighdeán den scoth 22242</div> <div>cáilíochta den scoth 15817</div>
4	<div>fiseán i caighdeán maith 29987</div> <div>aisce agus gan chlárú 22824</div> <div>porn saor in aisce 21962</div> <div>hd chaighdeán den scoth 21563</div> <div>t-ainm ar an scannán 17310</div>
5	<div>fiseán i caighdeán maith hd 22342</div> <div>fiseán i chaighdeán ard hd 10683</div> <div>aisce agus gan chlárú catagóir 10305</div> <div>féach ar saor in aisce 9290</div> <div>féach ar agus a íoslódáil 8555</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>