

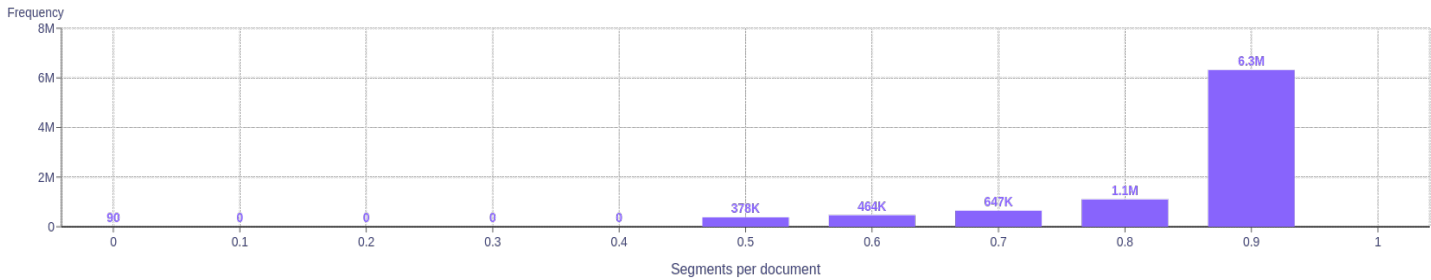
General overview

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

Volumes

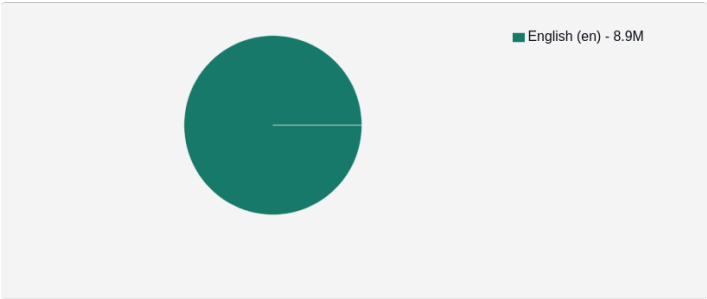
Segments	Unique segments	Src tokens	Trg tokens	Src size	Trg size	Src characters	Trg characters
8,905,979	8,905,890 (100.00 %)	165M	184M	834.16 MB	913.2 MB		

Translation likelihood



Language Distribution

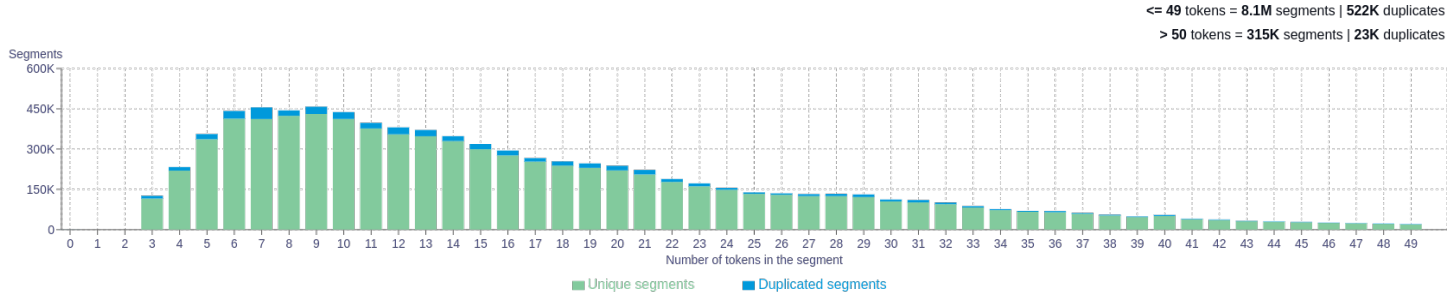
Source



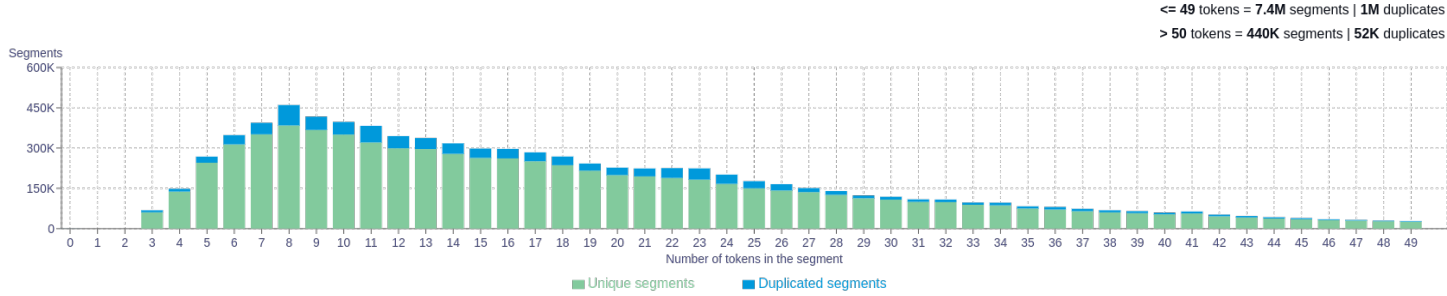
Target



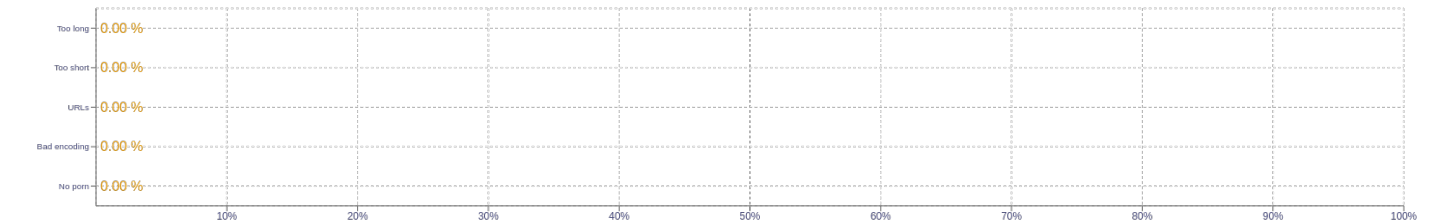
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	<div>hotel 576694</div> <div>car 509062</div> <div>best 474642</div> <div>airport 384841</div> <div>see 360587</div>
2	<div>car hire 259758</div> <div>remote control 183445</div> <div>best price 166453</div> <div>vat included 135162</div> <div>universal remote 98705</div>
3	<div>universal remote control 97614</div> <div>see available equivalences 90547</div> <div>rent a car 83577</div> <div>quickly and easily 79182</div> <div>see customer ratings 78522</div>
4	<div>models universal remote control 71234</div> <div>get the best price 61101</div> <div>find you the best 60960</div> <div>work hard to find 60954</div> <div>rent a car car 34191</div>
5	<div>amend your booking for free 65914</div> <div>rentalcars.com and you can amend 65913</div> <div>find you the best prices 60950</div> <div>book with us and get 60950</div> <div>android apps on google play 30939</div>

Target n-grams

Size	n-grams
1	<div>hotel 550673</div> <div>lloguer 487274</div> <div>hotels 468082</div> <div>millors 362903</div> <div>preus 297181</div>
2	<div>millors preus 146450</div> <div>millors descomptes 103549</div> <div>distància universal 96361</div> <div>lloc web 96255</div> <div>veure equivalències 90532</div>
3	<div>lloguer de cotxes 198059</div> <div>comandament a distància 156846</div> <div>descomptes en línia 103482</div> <div>veure equivalències disponibles 90532</div> <div>cotxe de lloguer 76050</div>
4	<div>millors descomptes en línia 103480</div> <div>hotels amb els millors 103480</div> <div>comandament a distància universal 96354</div> <div>models comandament a distància 71243</div> <div>forma fàcil i ràpida 71150</div>
5	<div>hotels amb els millors descomptes 103480</div> <div>models comandament a distància universal 71240</div> <div>reserveu online de forma fàcil 71114</div> <div>qualificacions dels clients i reserveu 71114</div> <div>consulteu les qualificacions dels clients 71114</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>