

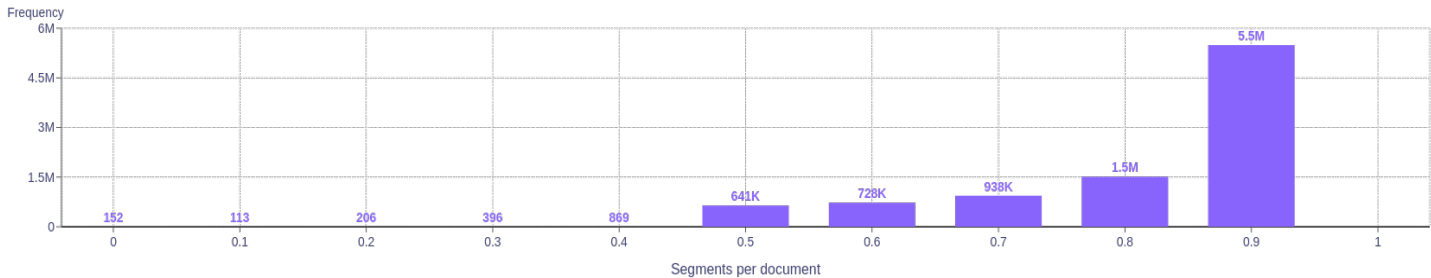
General overview

Corpus	Analytics date	Source language	Target language
HPLT.en-hr	10/27/2023	English (en)	Croatian (hr)

Volumes

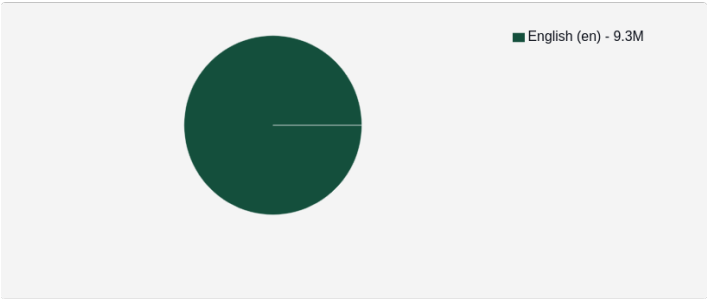
Segments	Unique segments	Src tokens	Trg tokens	Src size	Trg size	Src characters	Trg characters
9,310,369	9,310,276 (100.00 %)	162M	152M	815.9 MB	852.71 MB		

Translation likelihood

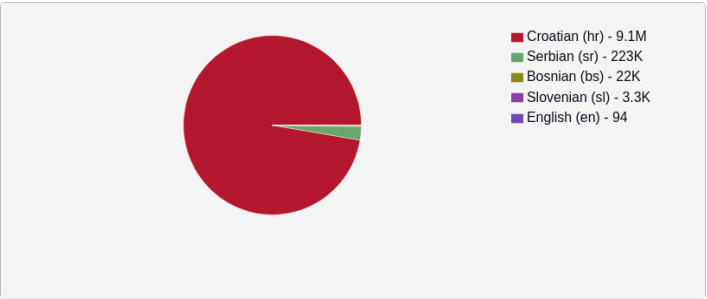


Language Distribution

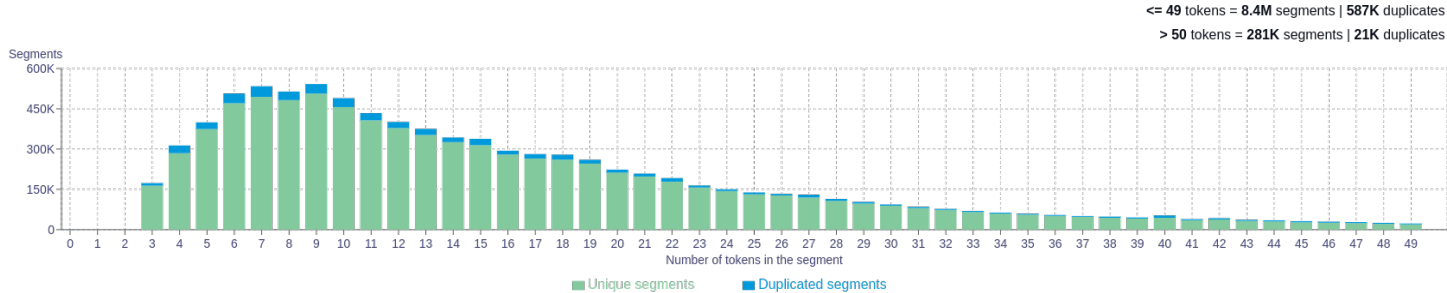
Source



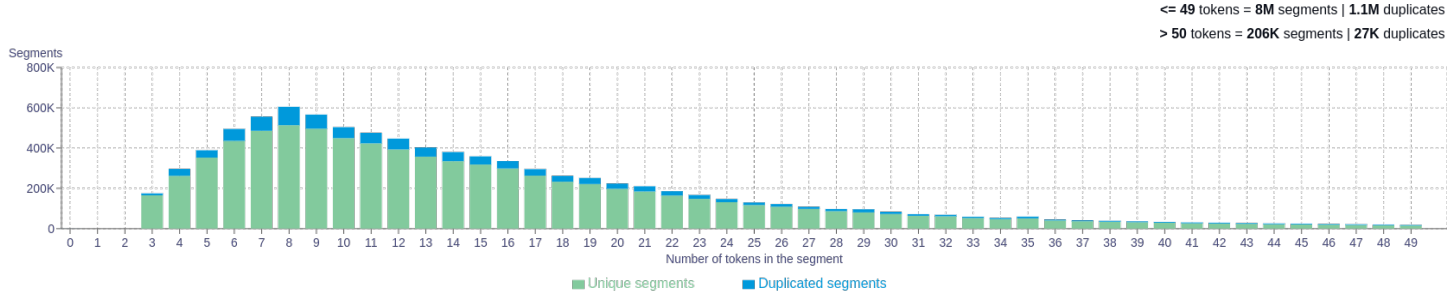
Target



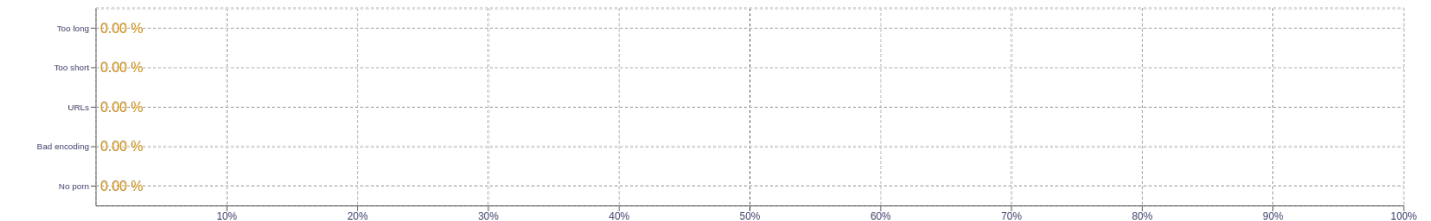
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	<div>hotel   370035</div> <div>weather   326122</div> <div>new   306112</div> <div>free   288669</div> <div>used   285819</div>
2	<div>personal data   74918</div> <div>local time   71682</div> <div>best prices   71369</div> <div>buy used   68441</div> <div>operating hours   68097</div>
3	<div>year of manufacture   153819</div> <div>prices from either   58907</div> <div>either machinery dealers   58907</div> <div>dealers or private   58907</div> <div>freemeteo hotel bookings   44971</div>
4	<div>prices from either machinery   58907</div> <div>machinery dealers or private   58907</div> <div>dealers or private sellers   58907</div> <div>best prices from either   58907</div> <div>clouds freemeteo hotel bookings   44597</div>
5	<div>prices from either machinery dealers   58907</div> <div>machinery dealers or private sellers   58907</div> <div>either machinery dealers or private   58907</div> <div>best prices from either machinery   58907</div> <div>snow clouds freemeteo hotel bookings   44570</div>

Target n-grams

Size	n-grams
1	<div>vrijeme   350121</div> <div>više   286272</div> <div>hotel   269926</div> <div>godina   223323</div> <div>može   182135</div>
2	<div>godina proizvodnje   156248</div> <div>zračna luka   90361</div> <div>vremenska prognoza   86402</div> <div>engleskom jeziku   85952</div> <div>radni sati   80648</div>
3	<div>dugoročna vremenska prognoza   66747</div> <div>freemeteo rezervacije hotela   44967</div> <div>više o smještajnom   44879</div> <div>oblaci freemeteo rezervacije   44574</div> <div>snijeg oblaci freemeteo   44537</div>
4	<div>strojeva bilo od privatnih   51508</div> <div>cijenama bilo od distributera   51508</div> <div>više o smještajnom objektu   44879</div> <div>oblaci freemeteo rezervacije hotela   44574</div> <div>snijeg oblaci freemeteo rezervacije   44537</div>
5	<div>strojeva bilo od privatnih prodavača   51508</div> <div>najboljim cijenama bilo od distributera   51508</div> <div>distributera strojeva bilo od privatnih   51508</div> <div>cijenama bilo od distributera strojeva   51508</div> <div>snijeg oblaci freemeteo rezervacije hotela   44537</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>