

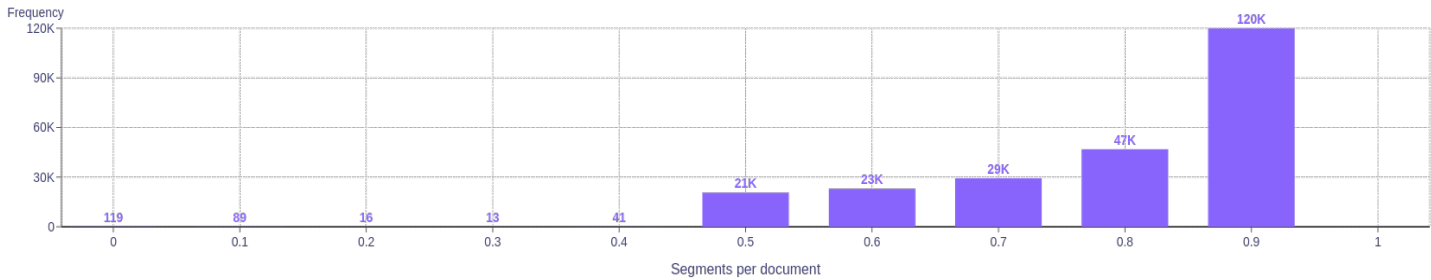
General overview

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

Volumes

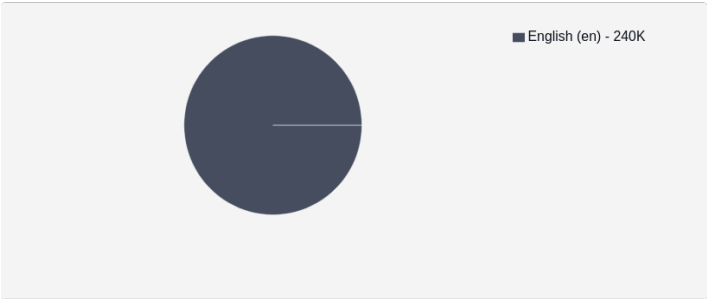
Segments	Unique segments	Src tokens	Trg tokens	Src size	Trg size	Src characters	Trg characters
240,015	240,013 (100.00 %)	3.2M	3.2M	16.89 MB	17.76 MB		

Translation likelihood

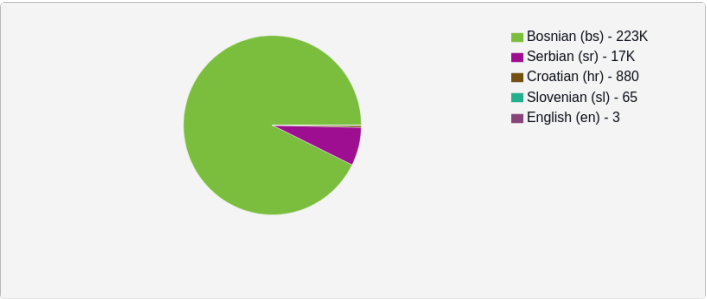


Language Distribution

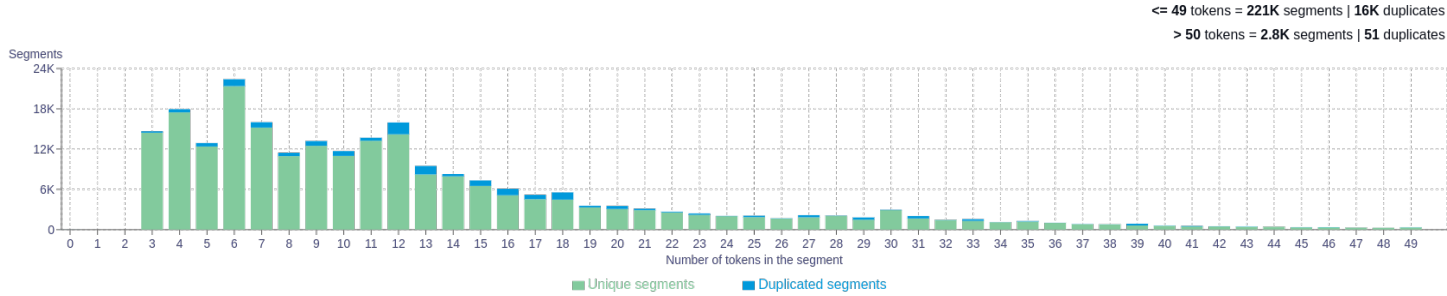
Source



Target

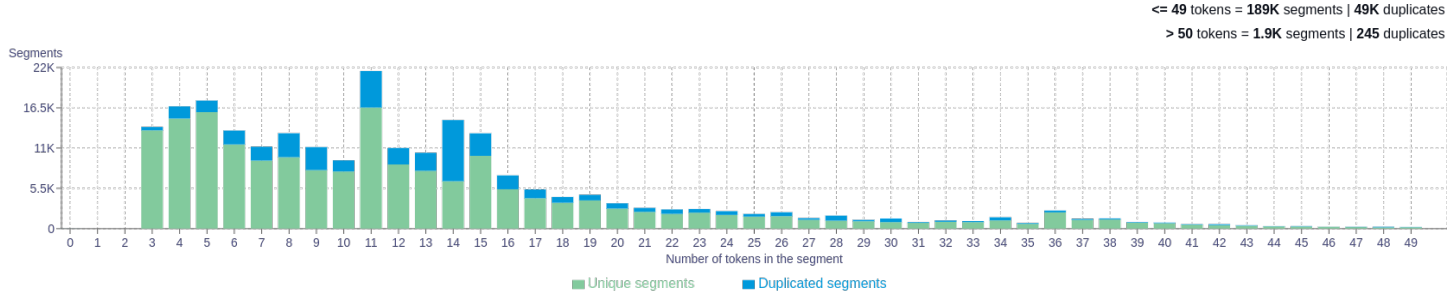


Source segment length distribution by token



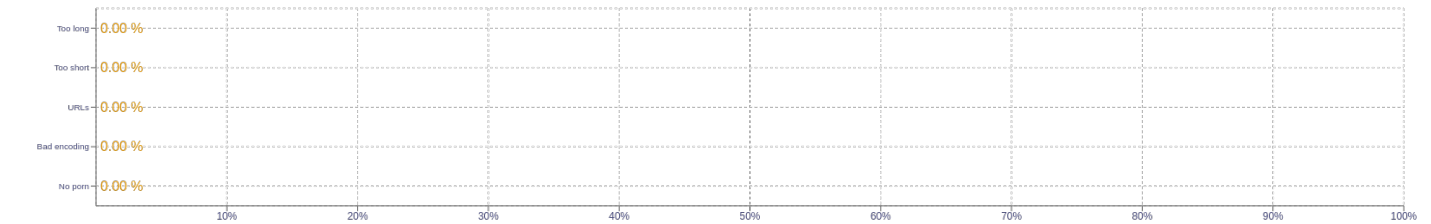
<= 49 tokens = 221K segments | 16K duplicates  
> 50 tokens = 2.8K segments | 51 duplicates

Target segment length distribution by token



<= 49 tokens = 189K segments | 49K duplicates  
> 50 tokens = 1.9K segments | 245 duplicates

Segment pair noise distribution



Source n-grams

Size	n-grams
1	<div>international   22926</div> <div>climate   15856</div> <div>used   13578</div> <div>united   12963</div> <div>usa   12768</div>
2	<div>international loads   6858</div> <div>international transportation   6338</div> <div>subtropical climate   6037</div> <div>humid subtropical   6037</div> <div>postal address   5181</div>
3	<div>humid subtropical climate   6037</div> <div>condition not indicated   4584</div> <div>cities and villages   4473</div> <div>climate humid subtropical   4023</div> <div>köppen climate classification   4020</div>
4	<div>nearby cities and villages   4464</div> <div>climate humid subtropical climate   4023</div> <div>transport cargoagent.net freight offers   3609</div> <div>cargoagent.net freight offers summary   3609</div> <div>offers summary international loads   3531</div>
5	<div>transport cargoagent.net freight offers summary   3609</div> <div>freight offers summary international loads   3531</div> <div>cargoagent.net freight offers summary international   3531</div> <div>get full analysis of name   2635</div> <div>exchange- international transportation and spedition   1687</div>

Target n-grams

Size	n-grams
1	<div>države   29219</div> <div>sjedinjene   28318</div> <div>američke   27756</div> <div>međunarodni   20832</div> <div>prevoz   14905</div>
2	<div>američke države   27600</div> <div>sjedinjene američke   27576</div> <div>međunarodni transport   9971</div> <div>međunarodni prevoz   9282</div> <div>nije navedeno   7062</div>
3	<div>sjedinjene američke države   27566</div> <div>vrućim ljetnim mjesecima   6959</div> <div>vlažna suptropska klima   6959</div> <div>klima s vrućim   6959</div> <div>tereti za međunarodni   6917</div>
4	<div>suptropska klima s vrućim   6959</div> <div>klima s vrućim ljetnim   6959</div> <div>okolnih gradova i sela   4470</div> <div>tereti za međunarodni transport   3090</div> <div>tereti za međunarodni prevoz   2983</div>
5	<div>vlažna suptropska klima s vrućim   6959</div> <div>suptropska klima s vrućim ljetnim   6959</div> <div>klima s vrućim ljetnim mjesecima   6959</div> <div>berza za međunarodni transport robe   2056</div> <div>kamiona za međunarodni prevoz robe   1630</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>