

General overview

Corpus	Analytics date	Source language	Target language
HPLT.en-ar	10/28/2023	English (en)	Arabic (ar)

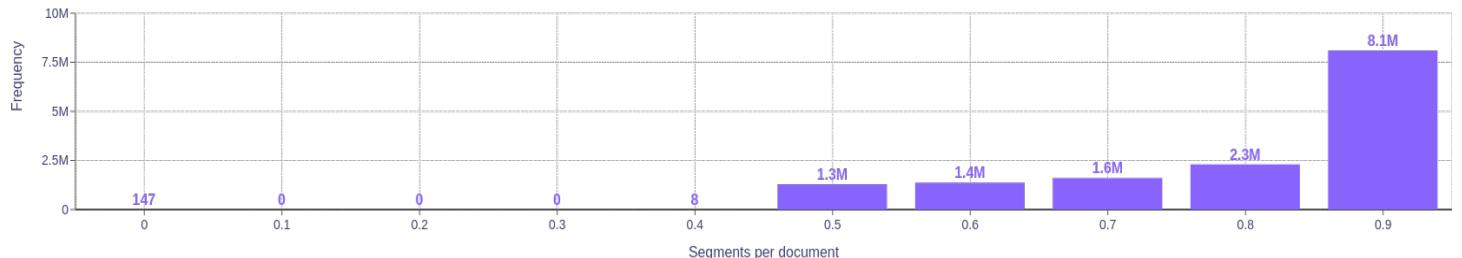
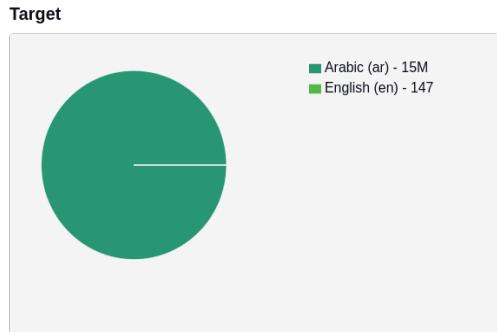
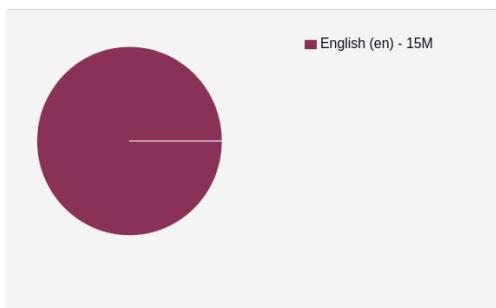
Volumes

Segments	Unique segments	Src tokens	Trg tokens	Src size	Trg size
14,645,275	3,652 (0.02 %)	282M	283M	1.4 GB	2.16 GB

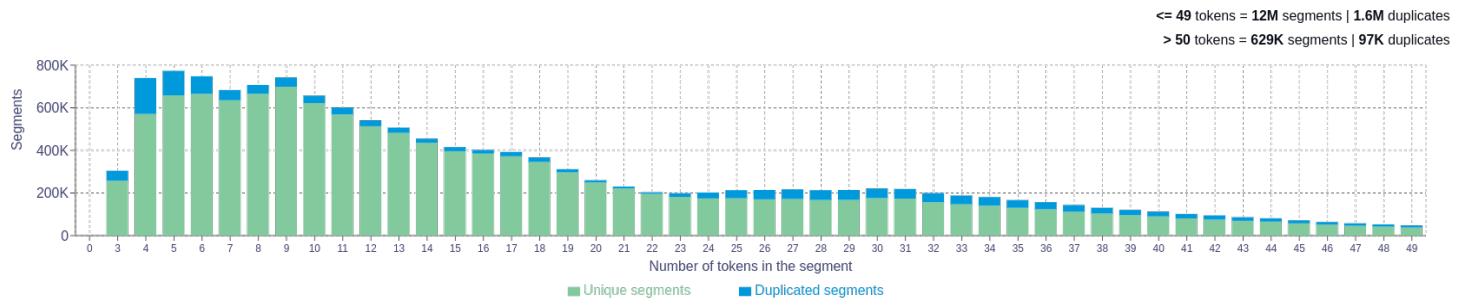
Type-Token Ratio

Source	Target
0.01	0.01

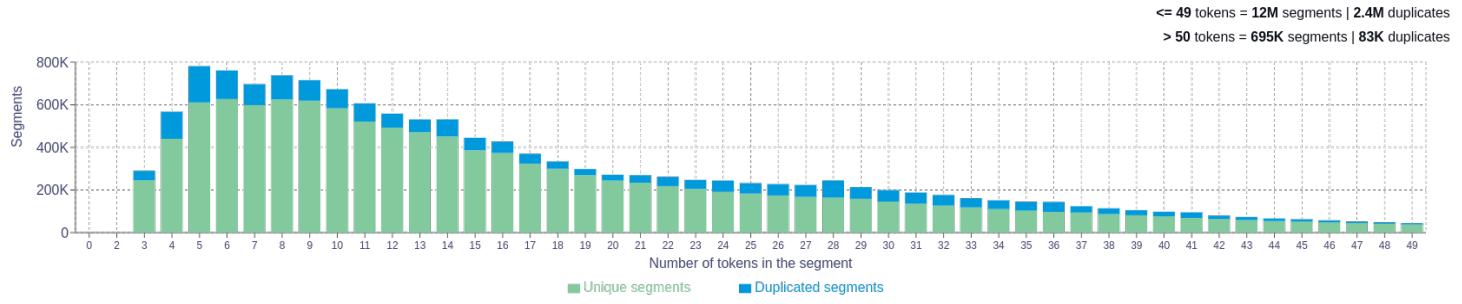
Translation likelihood

Language Distribution
Source

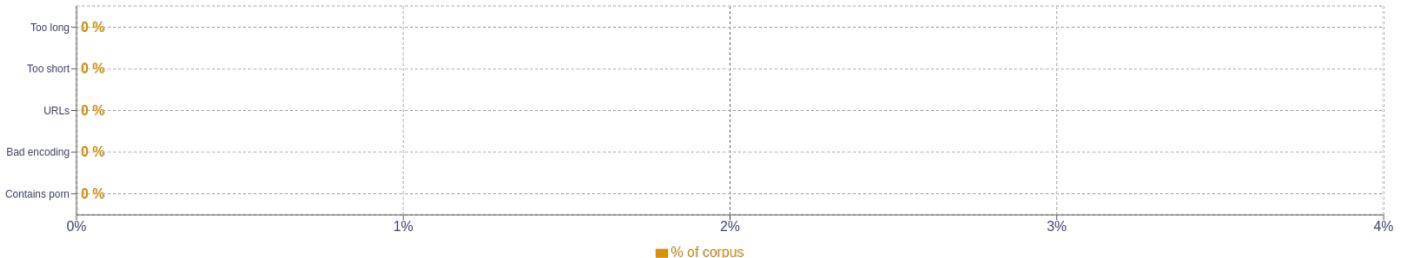
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	books 6081325 used 2113903 available 2022226 second 1929174 hand 1904473
2	second hand 1884215 rare books 1884038 used books 1883979 hand books 1883949 available rare 1883945
3	second hand books 1883949 books and second 1883945 available rare books 1883945 compare every offer 74397 offer archive entry 74145
4	used books and second 1883945 books of the title 1883945 books and second hand 1883945 things to do near 80302 compare every offer archive 74145
5	used books and second hand 1883945 hand books of the title 1883945 books and second hand books 1883945 compare every offer archive entry 74145 tripadvisor is proud to partner 45679

Target n-grams

Size	n-grams
1	(المستخدمة 3804053) (الكتب 2269116) (المناحة 2089717) (يتم 1927928) (1917519) (والكتب 1901642)
2	(ناتية للعنوان 1901470) (والتاريخ المنشورة 1901471) (الكتب المنشورة 1901470) (جهه ناتية 1901470) (المناحة 1901470) (والتاريخ المنشورة 1901470)
3	(المناحة والكتب المستخدمة 1901470) (والتاريخ المنشورة والكتب 1901470) (جهه ناتية للعنوان 1901470) (والتاريخ المنشورة والكتب 1901470) (المناحة والكتب المستخدمة والكتب 1901470)
4	(المناحة والكتب المستخدمة والكتب من جهة ناتية 1901470) (والتاريخ المنشورة والكتب المستخدمة 1901470) (والكتب من جهة ناتية 1901470) (المناحة والكتب المستخدمة والكتب من جهة ناتية 1901470) (والتاريخ المنشورة والكتب 1901470)
5	(المنشورة والكتب من جهة ناتية 1901470) (والتاريخ المنشورة والكتب من جهة ناتية للعنوان 1901470) (والكتب المستخدمة والكتب من جهة ناتية 1901470) (المنشورة والكتب المستخدمة والكتب من جهة ناتية 1901470) (والتاريخ المنشورة والكتب المستخدمة 1901470)

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/hpltt-project/data-analytics-tool/blob/main/tokenizers-info.md>

Type-Token Ratio

Lexical variety computed as *number or 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 distribution

Distribution of segments by fluency score

Distribution of segments by fluency score

Distribution of documents by average flux

Distribution of documents by average fluency score

Obtained with Monocleanel (<https://glu>)

Segment length distribution by token

Tokenized with https://github.com

Segment noise distribution

Obtained With Bicle

Frequent n-grams

Tokenized with <https://github.com/hplb-project/data-analytics-tool/blob/main/tokenizers-info.md>, after removing n-grams starting or ending in a stopword. Stopwords from <https://github.com/hplb-project/data-analytics-tool/blob/main/scripts/resources/README.txt>