

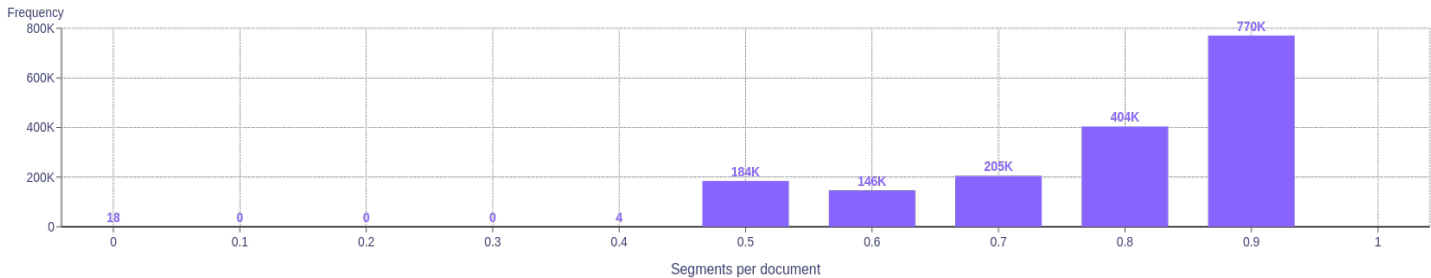
General overview

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

Volumes

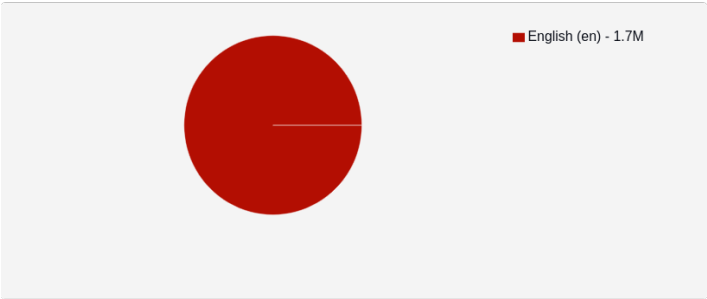
Segments	Unique segments	Src tokens	Trg tokens	Src size	Trg size	Src characters	Trg characters
1,710,223	1,710,206 (100.00 %)	25M	26M	112.98 MB	130.75 MB		

Translation likelihood

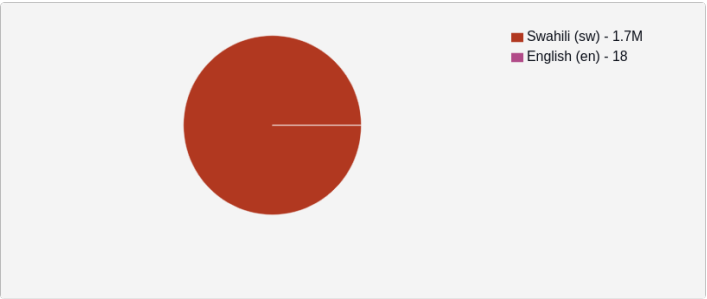


Language Distribution

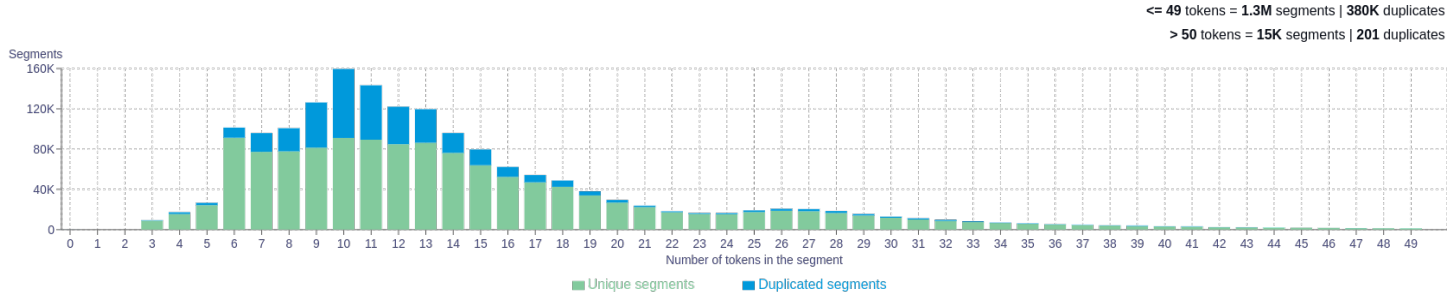
Source



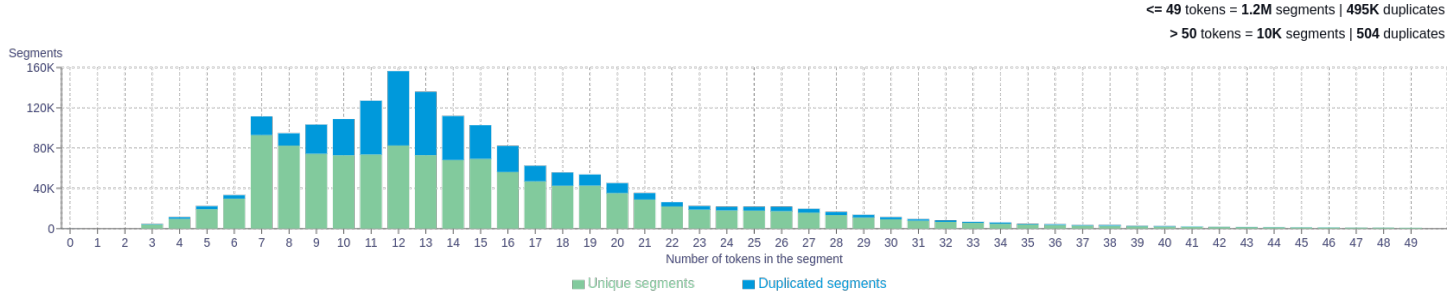
Target



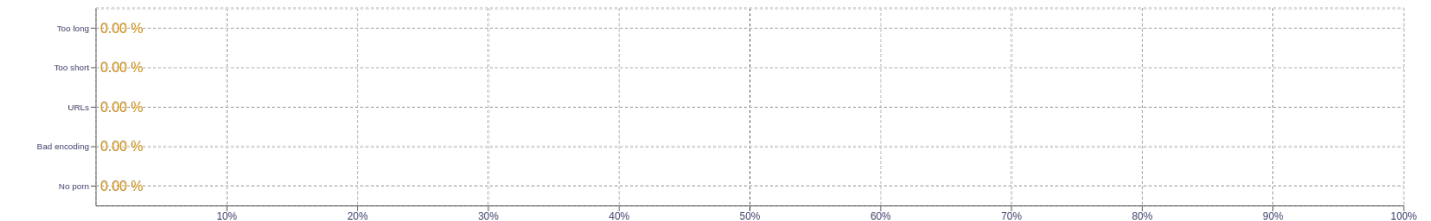
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	<div>english 545088</div> <div>used 408366</div> <div>books 318411</div> <div>year 295517</div> <div>ago 293229</div>
2	<div>year ago 290486</div> <div>rare books 84158</div> <div>used books 84146</div> <div>second hand 84140</div> <div>hand books 84121</div>
3	<div>second hand books 84121</div> <div>books and second 84121</div> <div>available rare books 84121</div> <div>apps on google 12118</div> <div>visit website email 10151</div>
4	<div>posted over a year 96028</div> <div>used books and second 84121</div> <div>books of the title 84121</div> <div>books and second hand 84121</div> <div>posts have been made 14153</div>
5	<div>posted over a year ago 95793</div> <div>used books and second hand 84121</div> <div>hand books of the title 84121</div> <div>books and second hand books 84121</div> <div>android apps on google play 12116</div>

Target n-grams

Size	n-grams
1	<div>lughā 594814</div> <div>kiingereza 568720</div> <div>kutumika 401071</div> <div>mwaka 313639</div> <div>uliopita 304802</div>
2	<div>zimeorodheshwa kabisa 85837</div> <div>vitabu vya 84577</div> <div>vya kichwa 84144</div> <div>vitabu kutumika 84126</div> <div>pili vitabu 84126</div>
3	<div>lughā ya kiingereza 566170</div> <div>mwaka mmoja uliopita 304319</div> <div>mkono wa pili 84135</div> <div>kutumika na mkono 84127</div> <div>vitabu vya kichwa 84126</div>
4	<div>ilitumwa zaidi ya mwaka 86552</div> <div>vitabu kutumika na mkono 84126</div> <div>pili vitabu vya kichwa 84126</div> <div>mkono wa pili vitabu 84126</div> <div>mwaka mmoja uliopita by 53485</div>
5	<div>kutumika na mkono wa pili 84127</div> <div>mkono wa pili vitabu vya 84126</div> <div>hakuna chapisho zilizowekwa kwa ukuta 14368</div> <div>programu za android kwenye google 12125</div> <div>tembelea tovuti tuma barua pepe 10242</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>