

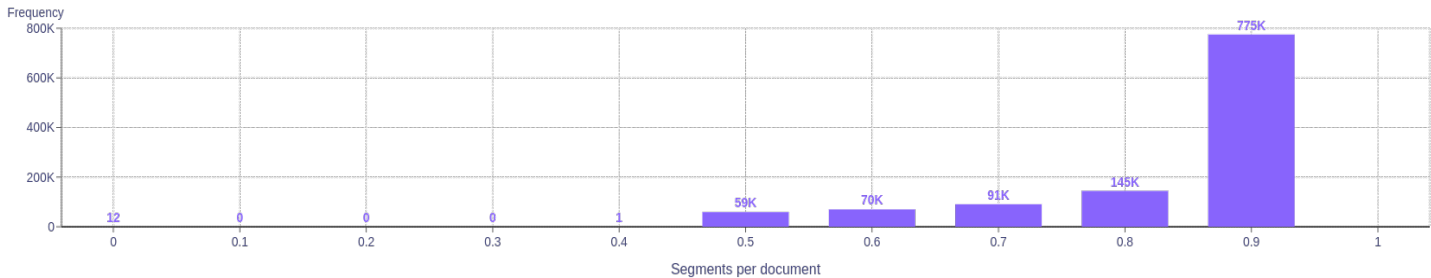
General overview

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

Volumes

Segments	Unique segments	Src tokens	Trg tokens	Src size	Trg size	Src characters	Trg characters
1,139,063	1,139,052 (100.00 %)	21M	21M	109.74 MB	200.95 MB		

Translation likelihood

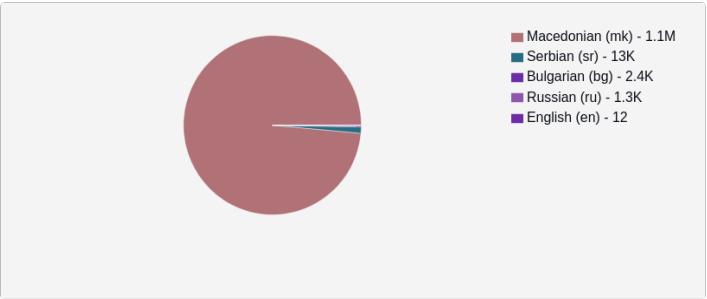


Language Distribution

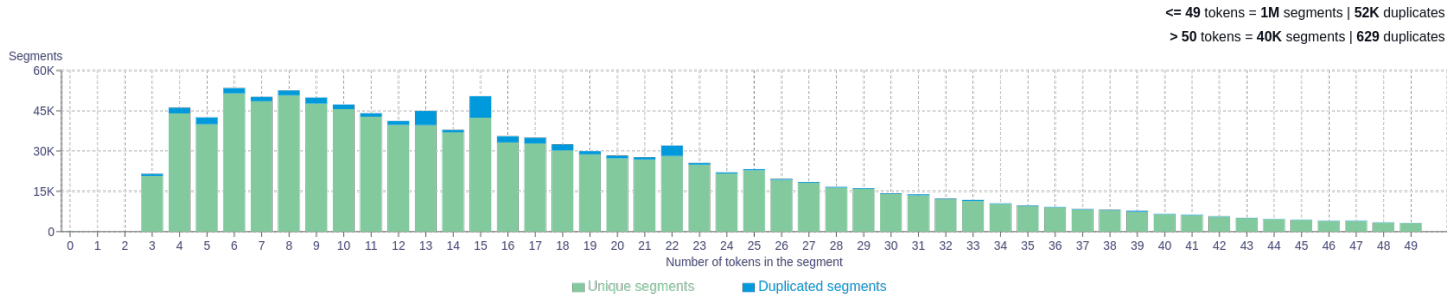
Source



Target

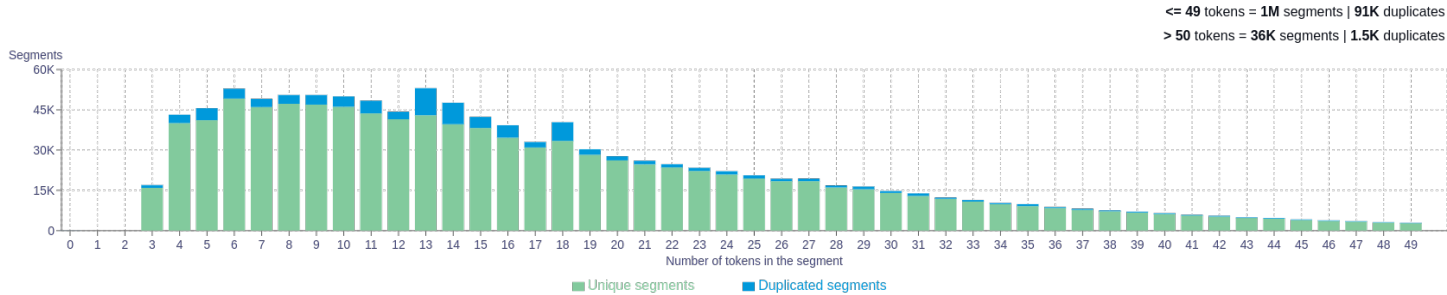


Source segment length distribution by token



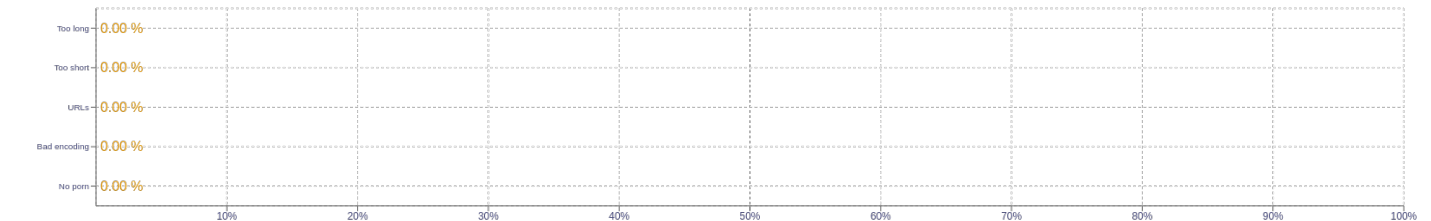
<= 49 tokens = 1M segments | 52K duplicates
> 50 tokens = 40K segments | 629 duplicates

Target segment length distribution by token



<= 49 tokens = 1M segments | 91K duplicates
> 50 tokens = 36K segments | 1.5K duplicates

Segment pair noise distribution



Source n-grams

Size	n-grams
1	<div>new 54656</div> <div>one 31981</div> <div>also 31829</div> <div>offers 30943</div> <div>first 30210</div>
2	<div>new offers 23628</div> <div>immeasurably grateful 9446</div> <div>working hours 8582</div> <div>search term 8172</div> <div>important details 8083</div>
3	<div>sure to appreciate 9446</div> <div>appreciate this gesture 9446</div> <div>clicking the button 7271</div> <div>republic of macedonia 6394</div> <div>call new offers 5597</div>
4	<div>gesture and be immeasurably 9446</div> <div>position on the main 7261</div> <div>machines offered at machineseeker 5140</div> <div>visit epoch and segpay 4927</div> <div>price not including vat 4289</div>
5	<div>sure to appreciate this gesture 9446</div> <div>gesture and be immeasurably grateful 9446</div> <div>position on the main page 7261</div> <div>first position on the main 7261</div> <div>please visit epoch and segpay 4927</div>

Target n-grams

Size	n-grams
1	<div>година 49110</div> <div>нови 32875</div> <div>македонија 29512</div> <div>понуди 27777</div> <div>време 26443</div>
2	<div>нови понуди 23657</div> <div>оцени гестот 9457</div> <div>бескрајно благодарна 9457</div> <div>работни часови 8444</div> <div>важни детали 8099</div>
3	<div>збоорот за пребарување 8171</div> <div>позиција на главната 7262</div> <div>кликнеш на копчето 7261</div> <div>резултати од пребарувањето 6900</div> <div>категирија на фирмата 6206</div>
4	<div>дефинитивно ќе го оцени 9457</div> <div>позиција на главната страница 7262</div> <div>првата позиција на главната 7261</div> <div>машини понудени на machineseeker 5278</div> <div>овластен застапник за продажба 4928</div>
5	<div>дефинитивно ќе го оцени гестот 9457</div> <div>првата позиција на главната страница 7261</div> <div>посетете ги ероч и сеграу 4928</div> <div>нашиот овластен застапник за продажба 4928</div> <div>фиксна цена не вклучувајќи ддв 2726</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>