

General overview

Corpus	Date	SL	TL
hplt-v2-en-is.tsv	1/22/2025	English (en)	Icelandic (is)

Volumes

Segments	SL tokens	SL characters	SL size
2,694,541	59M	295,187,731	282.92 MB

TL tokens	TL characters	TL size
54M	289,932,912	305.16 MB

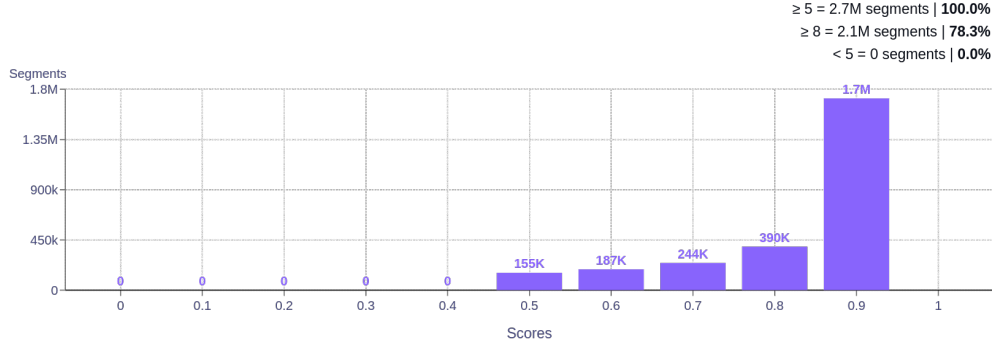
Dataset top 10 domains

SL domain	Segments	TL domain	Segments
hotels.com	78.7%	hotels.com	36.7%
biblegateway.com	10.9%	biblegateway.com	9.9%
booking.com	7.0%	booking.com	4.8%
wikipedia.org	3.8%	wikipedia.org	3.7%
vsaduidoma.com	2.6%	vsaduidoma.com	2.6%
eso.org	2.6%	collectiveray.com	2.3%
europa.eu	2.1%	eso.org	1.9%
lds.org	2.1%	jw.org	1.8%
collectiveray.com	2.1%	lds.org	1.7%
jw.org	1.9%	orangesmile.com	1.4%

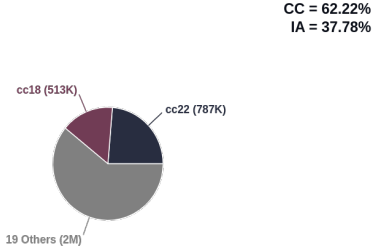
Dataset top 10 TLDs

SL domain	Segments	TL domain	Segments
com	160.8%	com	97.7%
org	20.5%	org	17.2%
is	9.0%	is	15.1%
net	5.2%	net	3.5%
eu	4.5%	eu	3.1%
nu	1.4%	nu	1.4%
co.uk	1.0%	de	0.8%
de	0.9%	info	0.7%
info	0.8%	top	0.6%
dk	0.6%	dk	0.6%

Translation likelihood

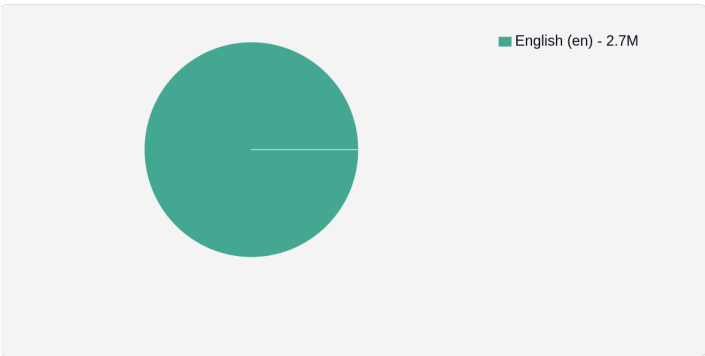


Collections

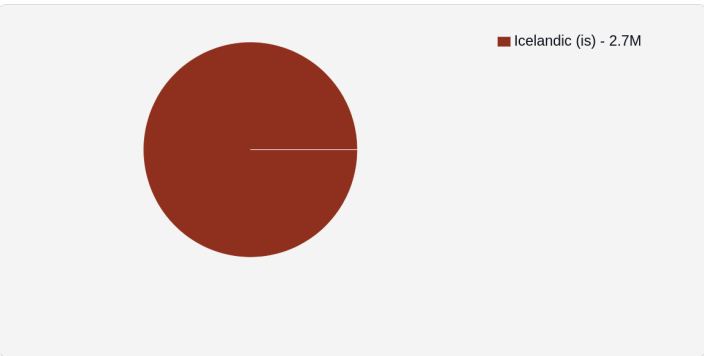


Language Distribution

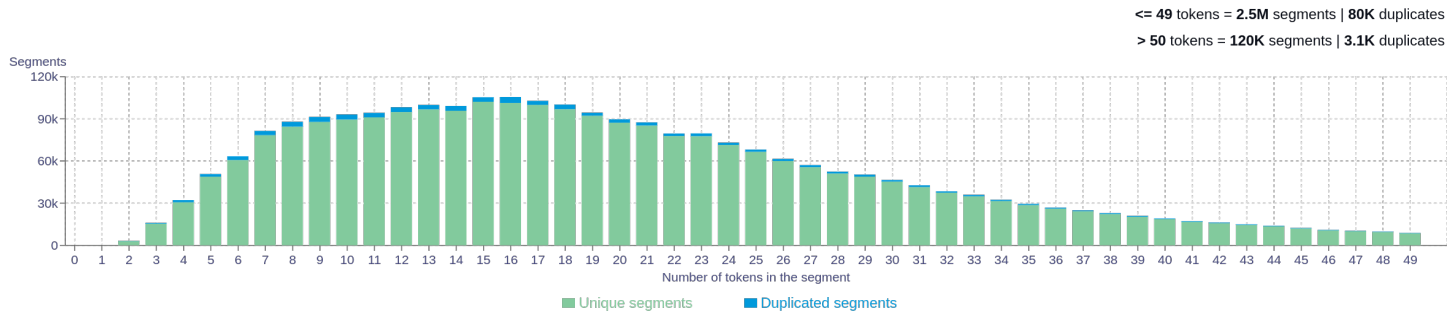
Source



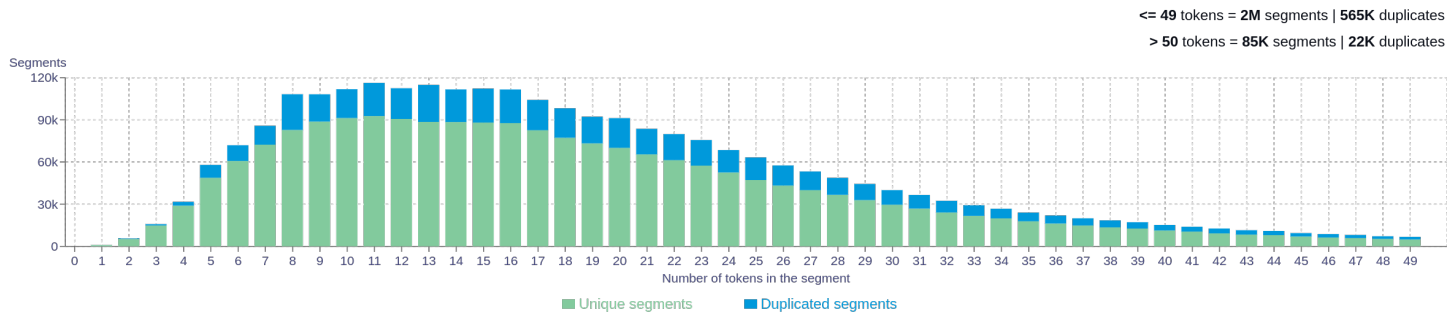
Target



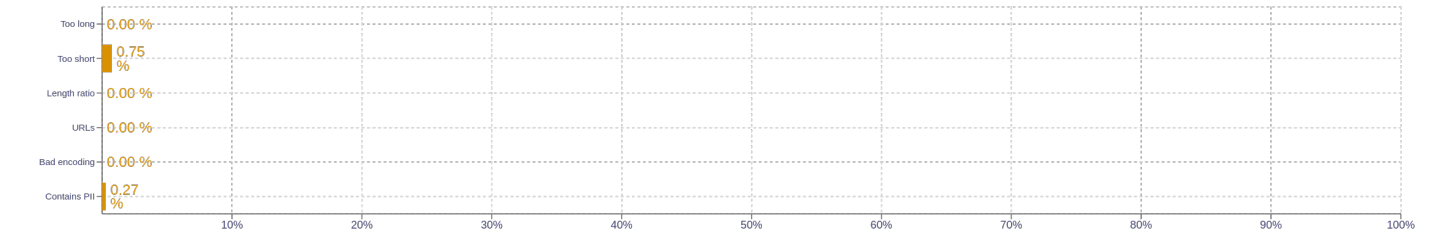
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	hotels 143231also 126529one 113504hotel 107412people 99859
2	last hour 27915people looked 27872hotels within 25971personal data 19165hotels near 14427
3	search for hotels 6989like the game 6245guest reviews yet 5480choice for travelers 5084last minute deals 4781
4	looked at this hotel 27849hotel in the last 27849hotels on a map 20321km from the city 17892great choice for travelers 5079
5	people looked at this hotel 27849hotel in the last hour 27849km from the city center 10433km from the city centre 7445great choice for travelers interested 5079

Target n-grams

Size	n-grams
1	var 161657getur 122212hafa 118042hótel 113560einnig 99824
2	siðustu klukkustund 27836einstaklingar skoðuðu 27828km fjarlægð 14402sjá hotel 12034getur verið 10452
3	skoðuðu þetta hotel 27828hotel á siðustu 27828km frá miðbænum 23270leitaðu að hótélum 15926hotel í nágrenninu 11841
4	hotel á siðustu klukkustund 27828einstaklingar skoðuðu þetta hotel 27828sjá hotel í nágrenninu 11838frábært val fyrir ferðalanga 7199býður upp á úrval 6399
5	skoðuðu þetta hotel á siðustu 27828hotel í nágrenninu á korti 11838fara á kort sem sýnir 6924hotels.com býður upp á úrval 6190samgöngurþegar flogið er á staðinn 5370

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>