

General overview

Corpus	Date	SL	TL
hplt-v2-en-ar.tsv	1/29/2025	English (en)	Arabic (ar)

Volumes

Segments	SL tokens	SL characters	SL size
17,505,366	508M	2,711,238,177	2.53 GB

TL tokens	TL characters	TL size
497M	2,506,340,590	4.15 GB

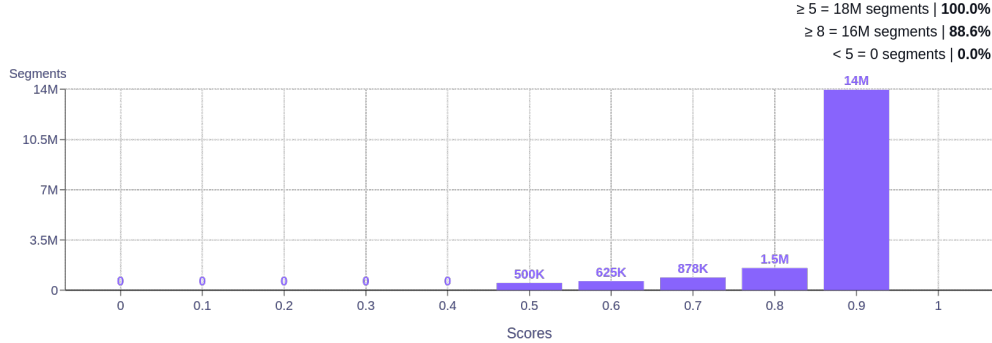
Dataset top 10 domains

SL domain	Segments	TL domain	Segments
hotels.com	9.0%	hotels.com	4.4%
wikipedia.org	3.5%	wikipedia.org	2.8%
alibaba.com	3.3%	alibaba.com	2.8%
microsoft.com	2.0%	microsoft.com	1.5%
booking.com	1.6%	ohchr.org	1.2%
ohchr.org	1.3%	booking.com	1.2%
office.com	0.9%	office.com	0.8%
airwise.com	0.9%	airwise.com	0.8%
wikihow.com	0.7%	ciwanekurd.net	0.7%
orangesmile.com	0.7%	wikihow.com	0.7%

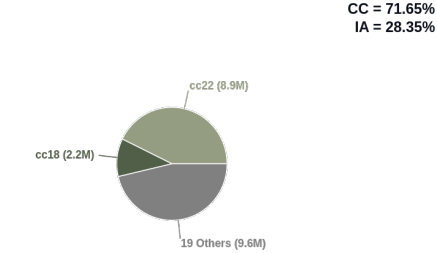
Dataset top 10 TLDs

SL domain	Segments	TL domain	Segments
com	101.3%	com	80.5%
org	21.0%	org	18.6%
net	7.5%	net	7.2%
ae	2.0%	ae	3.4%
co.uk	1.6%	com.eg	1.4%
info	0.8%	sa	1.0%
edu	0.8%	ma	0.8%
de	0.7%	info	0.8%
ru	0.6%	de	0.7%
ca	0.6%	ru	0.6%

Translation likelihood

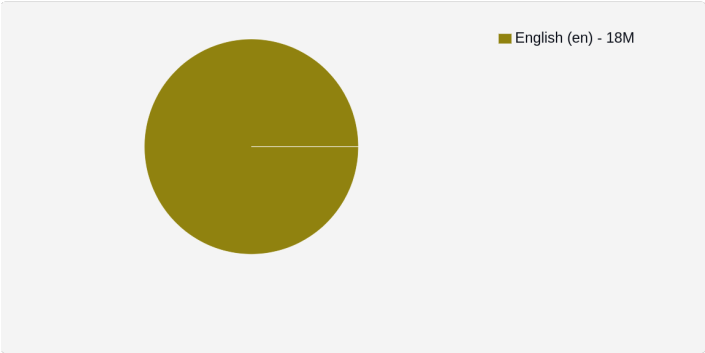


Collections

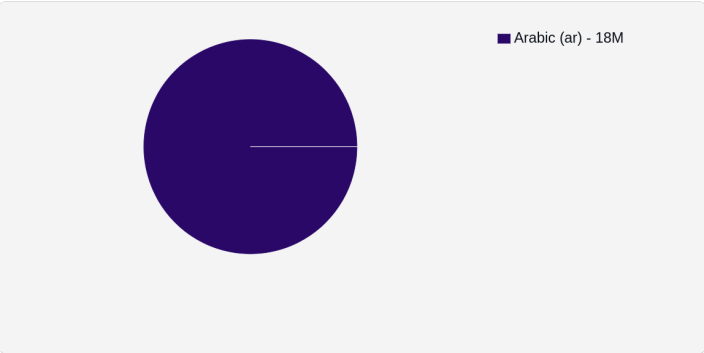


Language Distribution

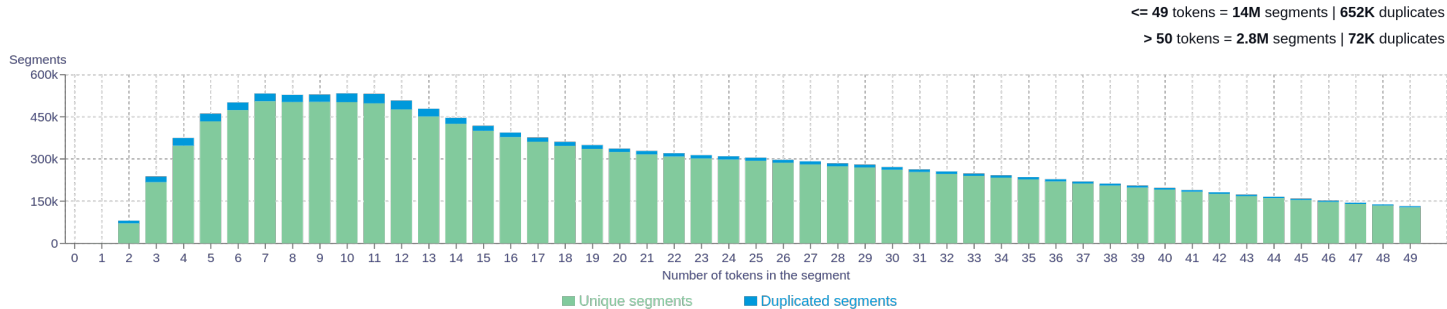
Source



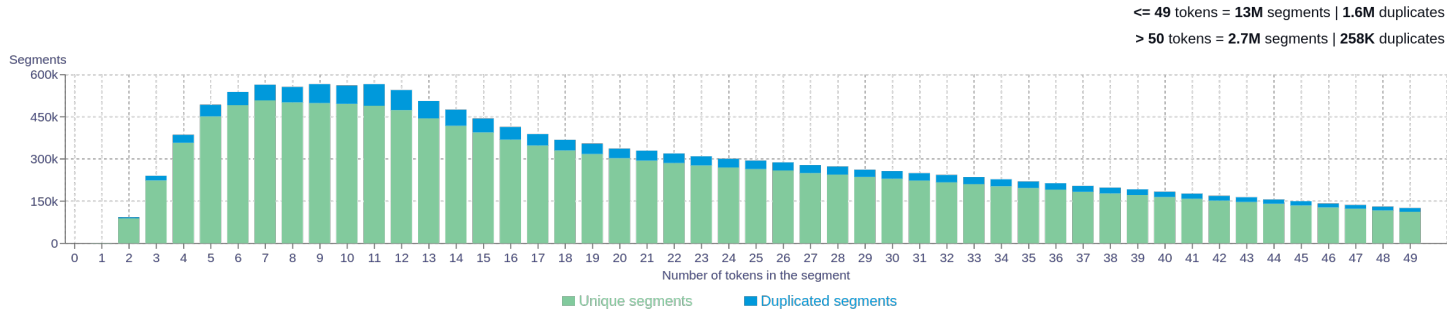
Target



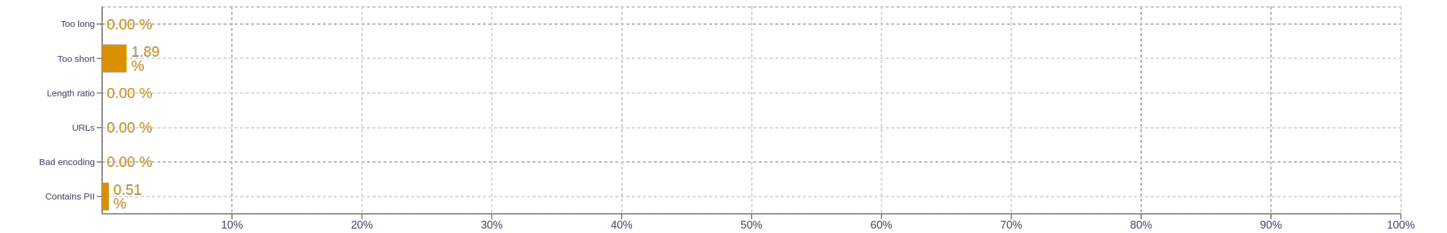
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	also 855982one 788031use 729131new 686901information 680604
2	show map 216877human rights 204723united states 143294personal data 94037united nations 91362
3	around the world 48639see on map 48424proud to partner 41906tripadvisor is proud 41755reservations with confidence 41754
4	discounts and special offers 39776find the perfect hotel 39712always with the best 39697best discounts and special 39694km from city centre 28378
5	tripadvisor is proud to partner 41754best discounts and special offers 39694month to find the perfect 39693always with the best discounts 39693travellers each month to find 20154

Target n-grams

Size	n-grams
1	981214 يمكن 930689 حلال 681969 بنم 676644 عام 667373 يمكنك
2	223972 الولايات المتحدة 187615 اعرض الخريطة 146184 سبيل المثال 127063 الأمم المتحدة 112328 غير الإنترنت
3	65445 الولايات المتحدة الأمريكية 64262 الإمارات العربية المتحدة 46755 يجب أن يكون 43982 المملكة العربية السعودية 42916 يمكن أن يكون
4	42026 يمكنك القيام بجورناك 41940 عمل على حد سواء 41940 العنور على الفندق المثالي 41935 أفضل الخصومات والعروض الخاصة 41934 وذلك مع القيام دائم
5	41934 والقيام برحلة عمل على حد 41934 لقضاء عطلة والقيام برحلة عمل 41934 شهر في العنور على الفندق 41934 برحلة عمل على حد سواء 41934 بتقديم أفضل الخصومات والعروض الخاصة

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>