

General overview

Corpus	Date	SL	TL
hplt-v2-en-ms.tsv	1/26/2025	English (en)	Malay (ms)

Volumes

Segments	SL tokens	SL characters	SL size
8,432,285	179M	933,359,950	893.65 MB

TL tokens	TL characters	TL size
174M	1,044,032,960	997.25 MB

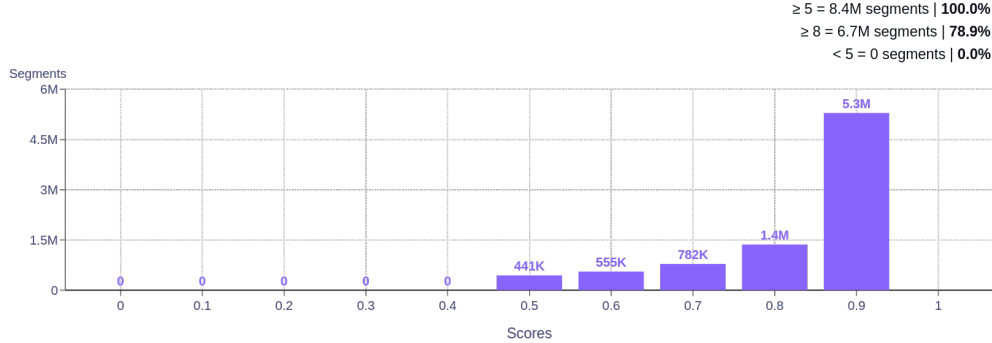
Dataset top 10 domains

SL domain	Segments	TL domain	Segments
hotels.com	63.9%	hotels.com	23.1%
google.com	19.8%	wikipedia.org	8.4%
wikipedia.org	10.4%	google.com	7.5%
agoda.com	8.3%	agoda.com	5.9%
booking.com	7.1%	booking.com	3.7%
orangesmile.com	1.6%	blogspot.com	2.3%
lacroix.com	1.5%	lacroix.com	1.5%
itsmygame.org	1.4%	hotelscombined.my	1.3%
masterstudies.com	1.4%	itsmygame.org	1.1%
airwise.com	1.4%	airwise.com	1.1%

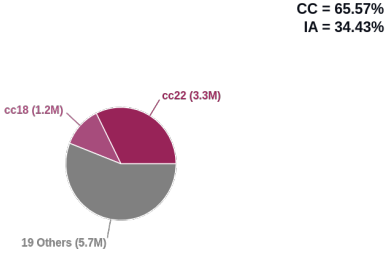
Dataset top 10 TLDs

SL domain	Segments	TL domain	Segments
com	182.2%	com	101.1%
org	19.5%	org	15.1%
net	5.3%	com.my	7.5%
com.my	3.1%	my	5.4%
gov.my	2.1%	net	4.1%
my	2.0%	gov.my	2.2%
co.uk	1.7%	ru	1.0%
ru	1.0%	info	0.7%
ca	1.0%	edu.my	0.6%
fm	0.8%	nu	0.5%

Translation likelihood

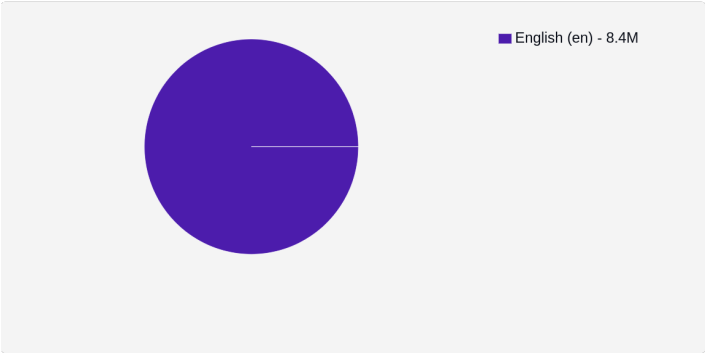


Collections

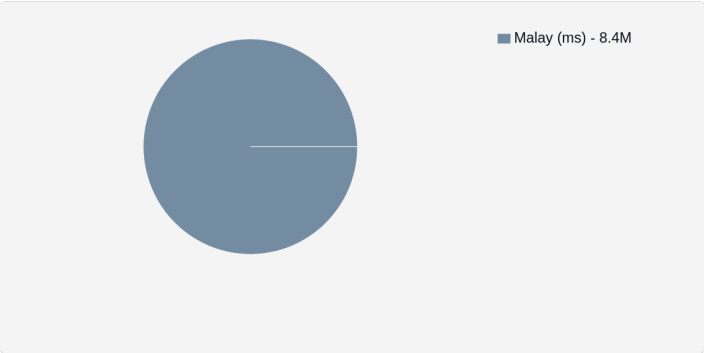


Language Distribution

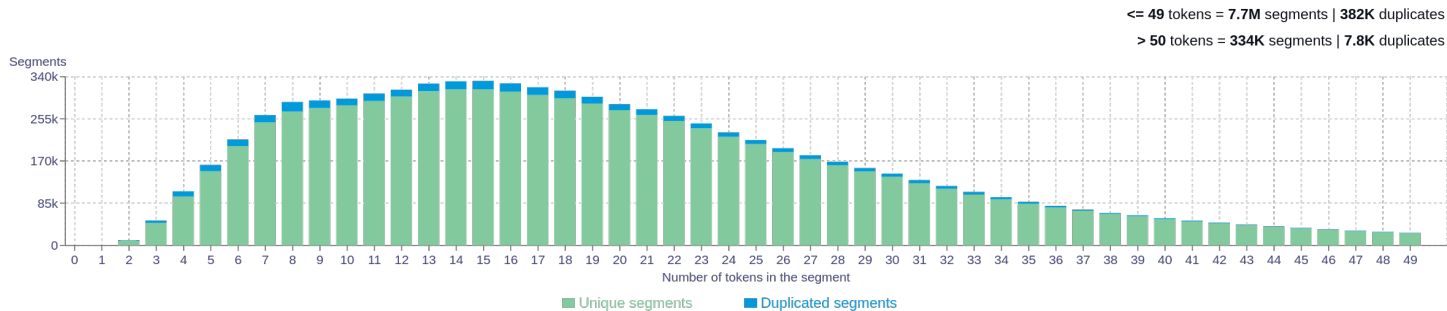
Source



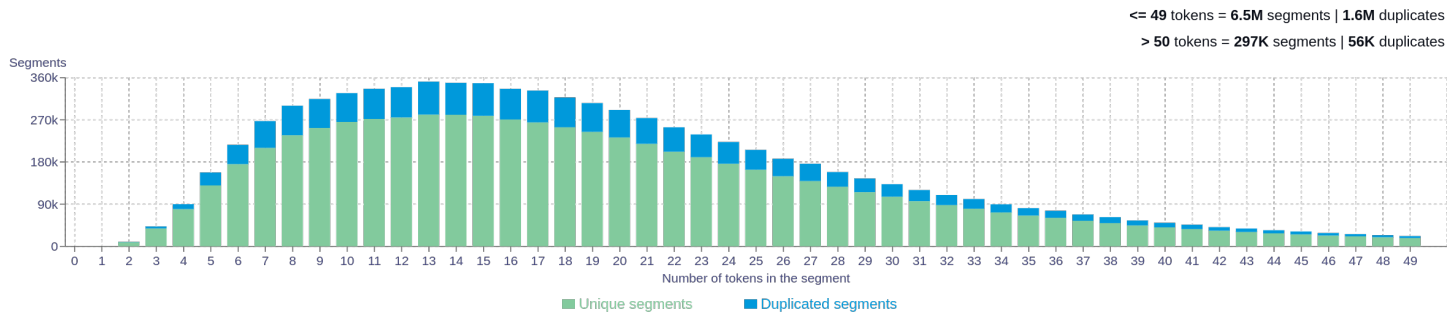
Target



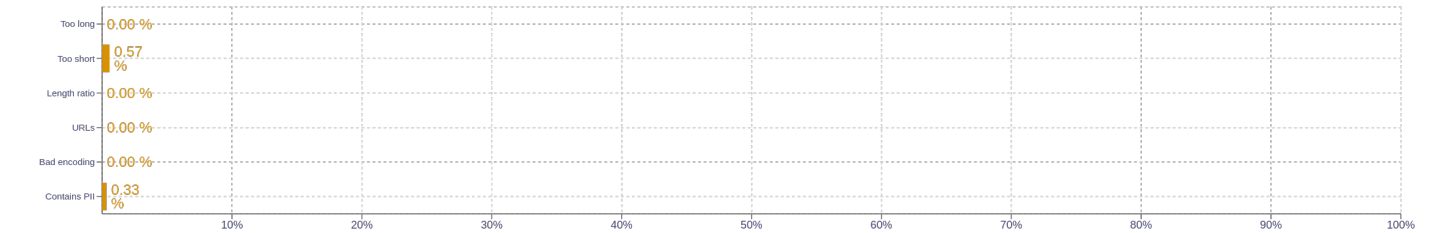
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	hotel 442938also 375479use 306657one 285506may 280405
2	personal data 70848personal information 61634public areas 40530privacy policy 36556kuala lumpur 31854
3	wi-fi in public 23435available to hotel 18578terms and conditions 18481choice for travelers 16476like the game 15652
4	hotels on a map 27564wi-fi in public areas 23412available to hotel guests 18552wi-fi available to hotel 18551wi-fi in all rooms 16759
5	wi-fi available to hotel guests 18551free wi-fi in all rooms 16749people looked at this hotel 12432hotel in the last hour 12432cards and reserves the right 9876

Target n-grams

Size	n-grams
1	anda 2924614hotel 681806permainan 327120bilik 295982laman 254012
2	laman web 188405peribadi anda 87062data peribadi 85495berjalan kaki 50491sekiranya anda 46497
3	data peribadi anda 49798minit berjalan kaki 39443laman web anda 31257terma dan syarat 19440permainan dalam talian 19342
4	fi untuk tetamu hotel 19033berkualiti tinggi bergerak telefon 18167percuma dalam semua bilik 10580kad kredit sebelum ketibaan 9921kredit sebelum ketibaan anda 9912
5	makan tengah hari dan makan 14794hotel ini dalam masa sejam 13163fi percuma dalam semua bilik 10579kad kredit sebelum ketibaan anda 9912menahan jumlah wang yang ditetapkan 9638

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>