

General overview

Corpus	Date	Language
hplt-v3-yue_Hant	9/18/2025	Chinese (zh)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
217,261	4,618,489	3,461,753 (74.95 %)	158M	272,601,981	661.32 MB

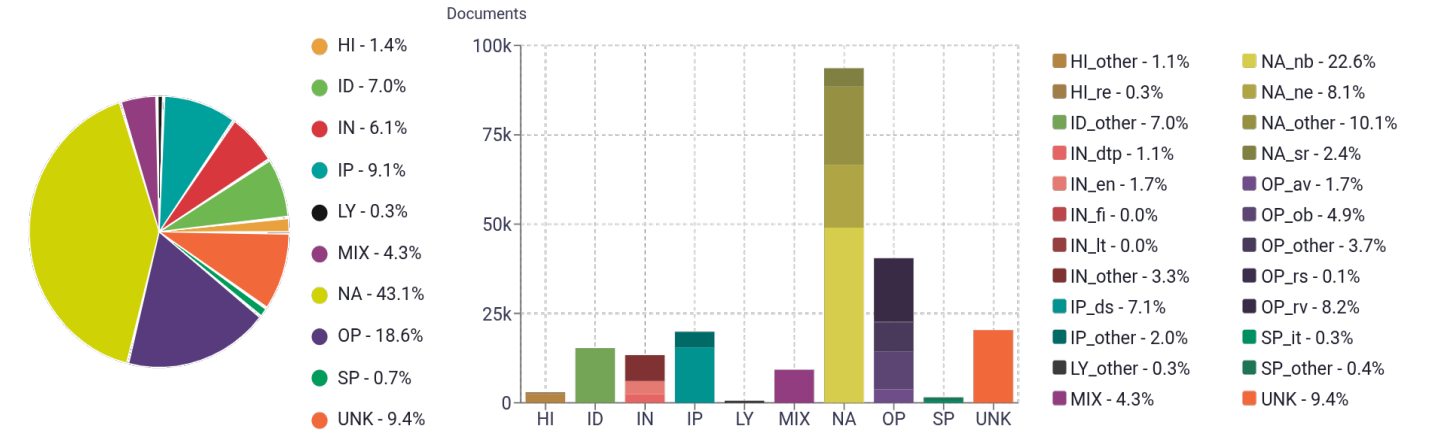
Top 10 domains

Domain	Docs	% of total
blogspot.hk	13K	6.04%
blogspot.com	11K	4.99%
openrice.com	10K	4.63%
on.cc	7.3K	3.38%
hotels.com	4.2K	1.92%
wikipedia.org	3.9K	1.78%
presslogic.com	3.7K	1.71%
fanpiece.com	3.6K	1.65%
yahoo.com	3.4K	1.57%
hkgolden.com	3.3K	1.53%

Top 10 TLDs

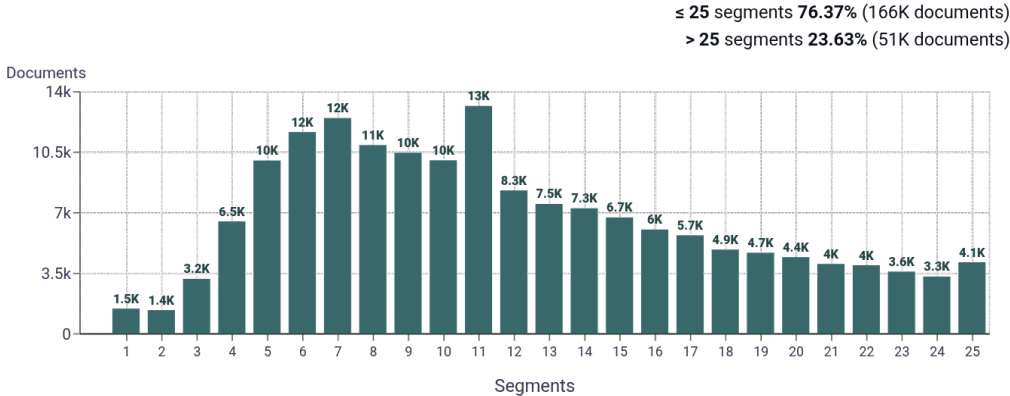
Domain	Docs	% of total
com	123K	56.83%
hk	32K	14.59%
com.hk	23K	10.73%
cc	7.7K	3.54%
org	5.9K	2.74%
net	5.8K	2.66%
tw	2.4K	1.09%
name	2.2K	1.02%
me	2.2K	0.99%
info	2K	0.92%

Register labels

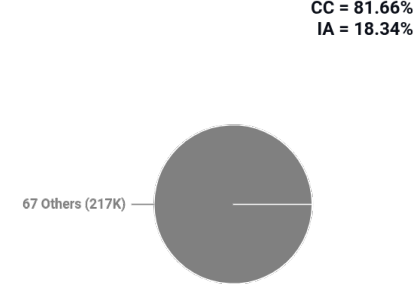


MT:0.2% | 399 Documents

Documents size (in segments)

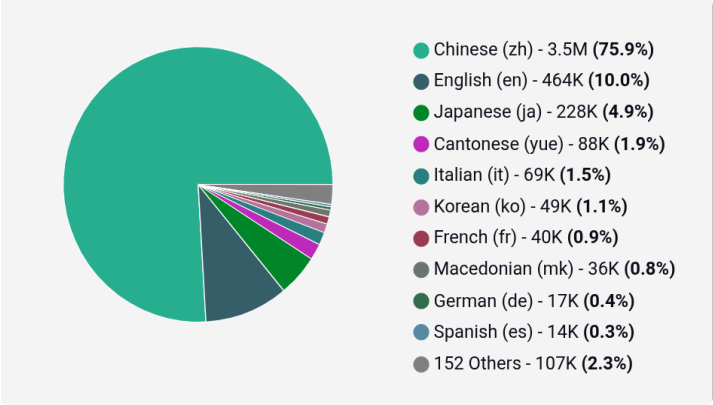


Document collections

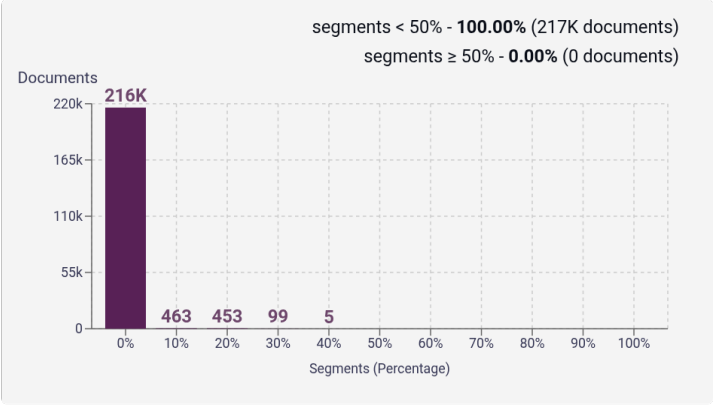


Language Distribution

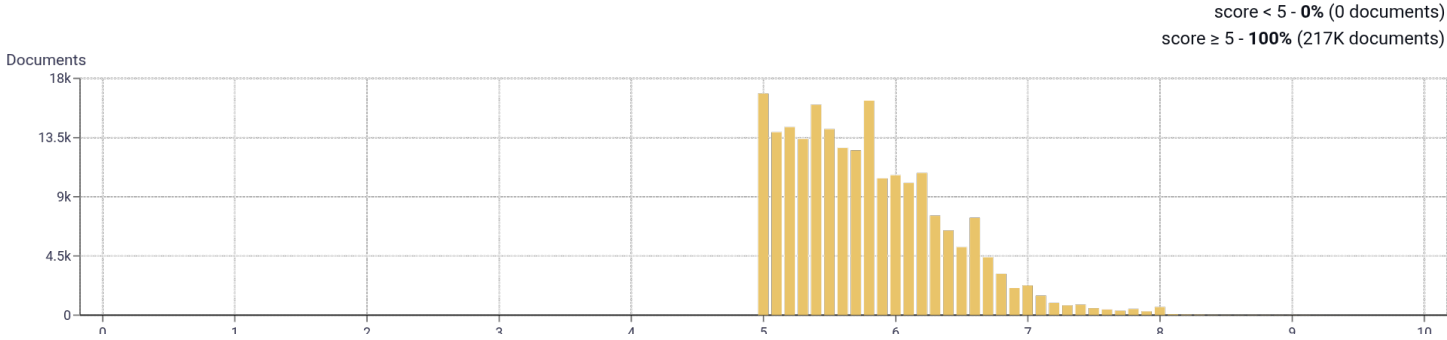
Number of segments in the Chinese (zh) corpus



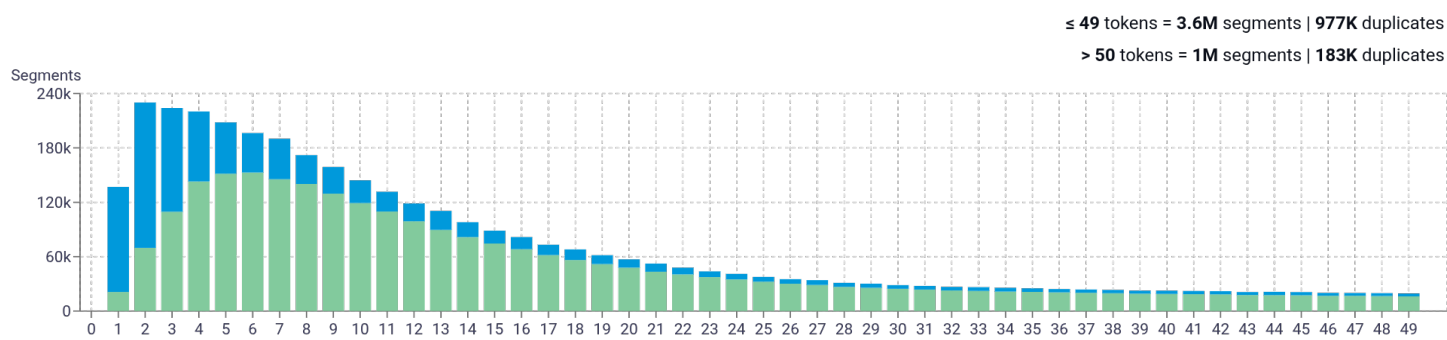
Percentage of segments in Chinese (zh) inside documents



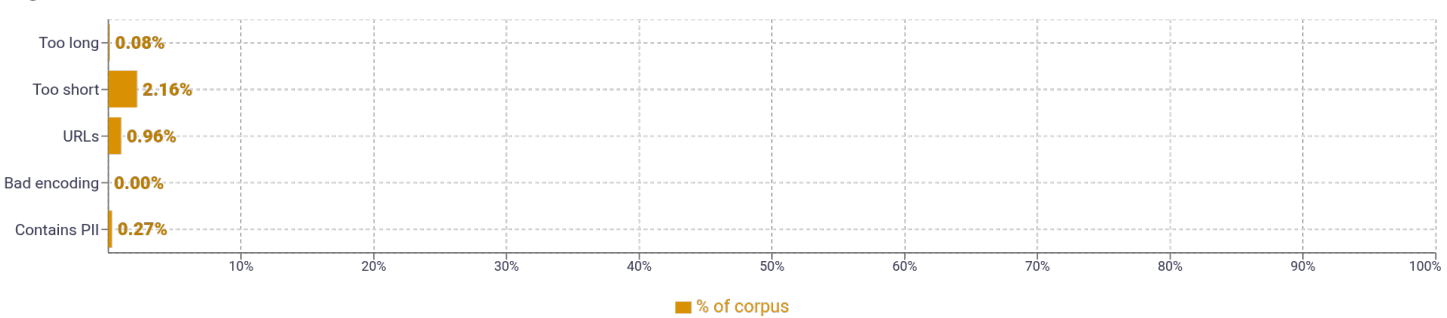
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	都 2,184,182 係 1,103,209 好 894,503 個 854,739 嘅 833,303	
2	都係 197,184 已經 124,719 都好 119,215 不過 115,423 都唔 93,760	
3	係一個 39,640 覺得好 14,592 好多人都 12,432 真係好 8,652 咁就唔會 8,078	
4	都係一個 6,302 pe n a n 5,313 仲有各種 3,911 讀者有任何關於 3,846 知~ 大家可以透過 3,844	
5	讀者有任何關於美食 3,846 知~ 大家可以透過 facebook 3,844 們知~ 大家可以透過 3,844 都歡迎話比小編 3,843 迎話比小編們知~ 3,842	

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>

Register labels

Name	Abbr.	Name	Abbr.	Name	Abbr.
Machine-translated	MT	How-to or instructions	HI	Description of a thing or person	dtp
Lyrical	LY	Recipe	re	FAQ	fi
Spoken	SP	Informational persuasion	IP	Legal terms & conditions	lt
Interview	it	Description with intent to sell	ds	Opinion	OP
Interactive discussion	ID	News & opinion blog or editorial	ed	Review	rv
Narrative	NA	Informational description	IN	Opinion blog	ob
News report	ne	Enciclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				