

General overview

Corpus	Date	Language
hplt-v3-kor_Hang	9/18/2025	Korean (ko)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
74,788,170	2,328,273,763	1,380,329,095 (59.29 %)	77B	162,051,985,657	347.45 GB

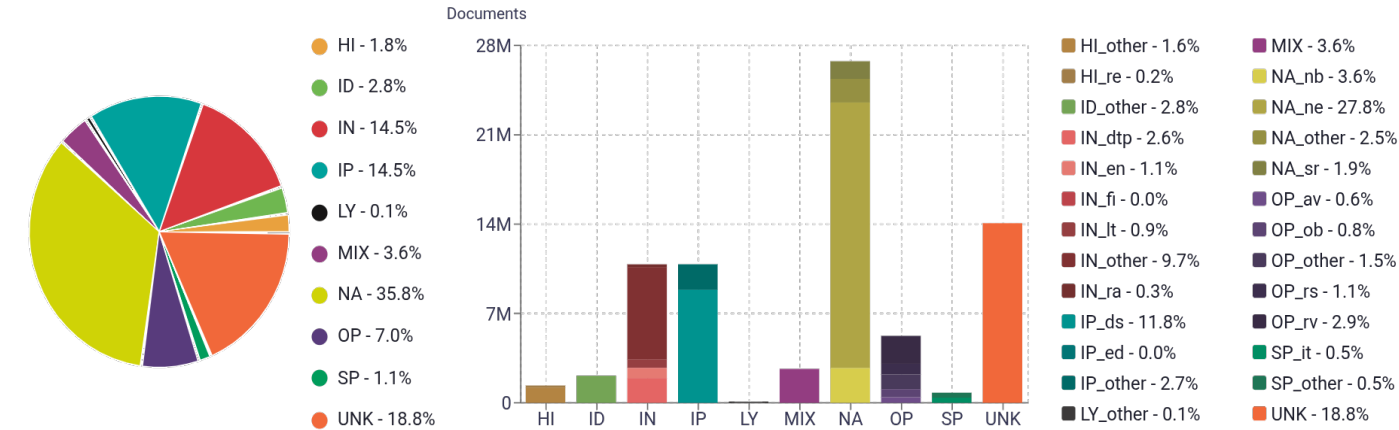
Top 10 domains

Domain	Docs	% of total
tistory.com	2.5M	3.28%
hankyung.com	1.3M	1.75%
happycampus.com	1.1M	1.42%
joins.com	933K	1.25%
donga.com	852K	1.14%
naver.com	827K	1.11%
daum.net	785K	1.05%
egloos.com	729K	0.97%
mt.co.kr	727K	0.97%
junfile.com	471K	0.63%

Top 10 TLDs

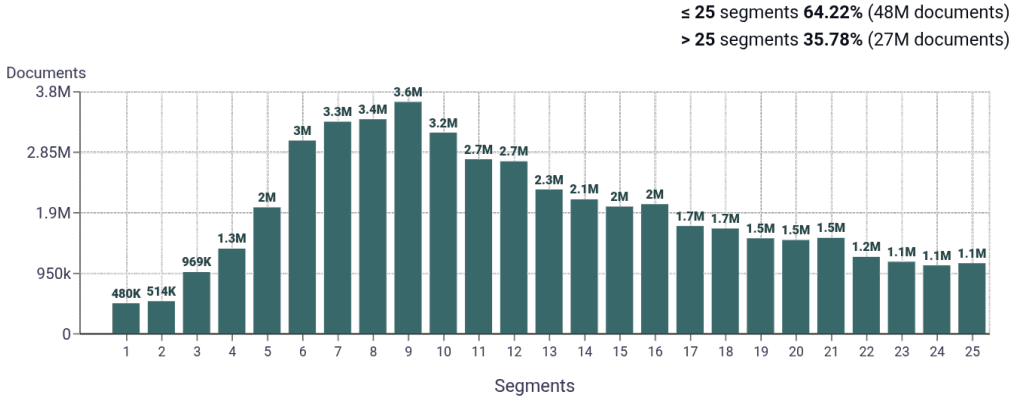
Domain	Docs	% of total
com	37M	49.07%
co.kr	18M	24.47%
kr	5.5M	7.33%
net	4.2M	5.55%
org	2M	2.69%
or.kr	1.5M	2.02%
go.kr	789K	1.05%
ac.kr	775K	1.04%
xyz	368K	0.49%
icu	272K	0.36%

Register labels



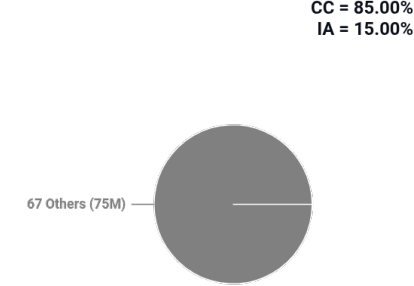
MT:13.2% | 9.9M Documents

Documents size (in segments) ⓘ



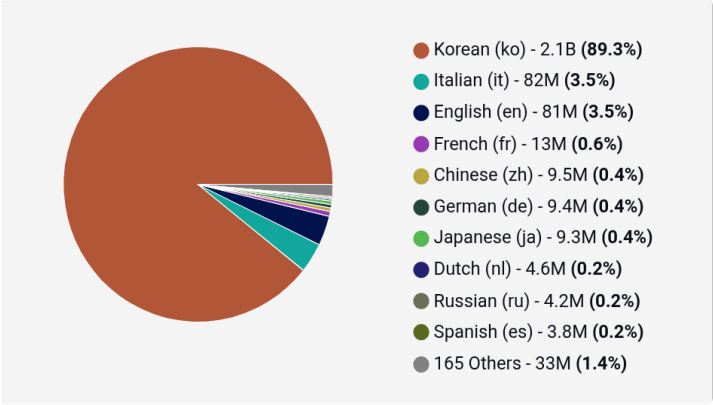
≤ 25 segments 64.22% (48M documents)
> 25 segments 35.78% (27M documents)

Document collections

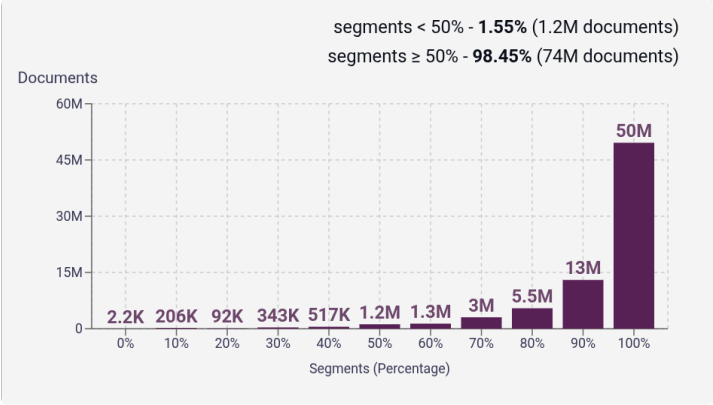


Language Distribution

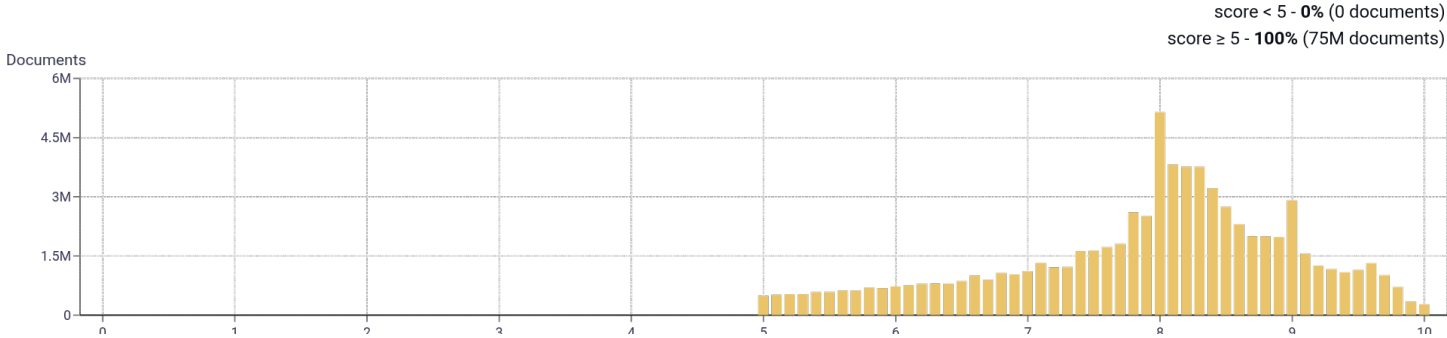
Number of segments in the Korean (ko) corpus



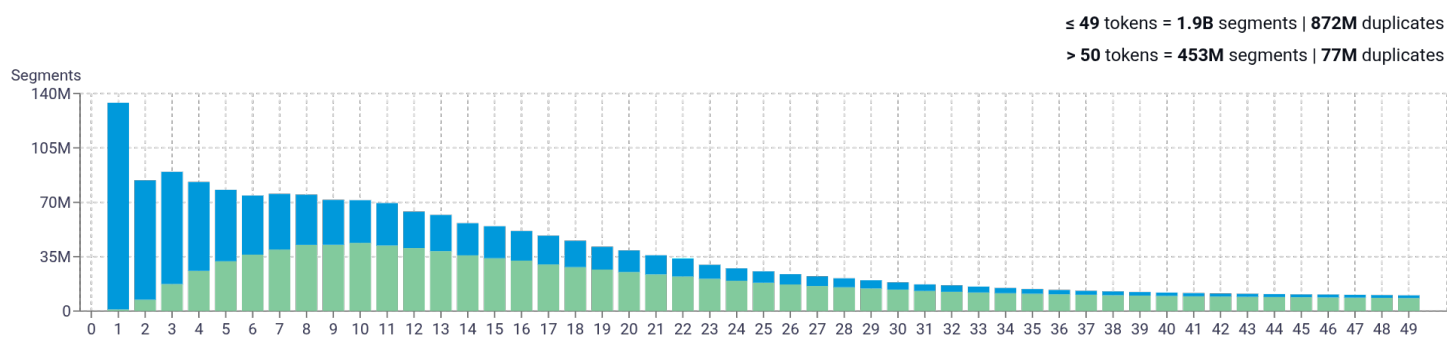
Percentage of segments in Korean (ko) inside documents



Distribution of documents by document score

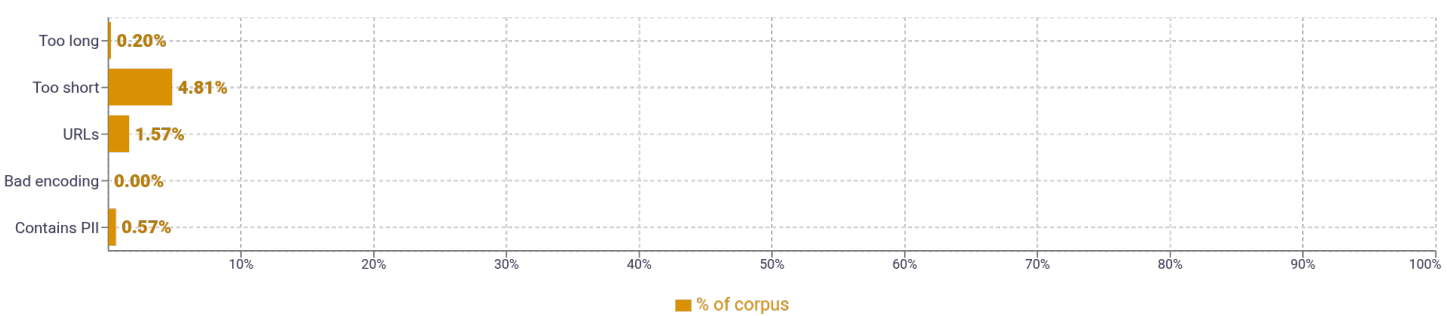


Segment length distribution by token



≤ 49 tokens = 1.9B segments | 872M duplicates
> 50 tokens = 453M segments | 77M duplicates

Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	는 1,775,783,166 다 989,722,927 은 978,066,093 고 823,156,830 있 801,191,099	
2	수 있 281,194,521 고 있 209,470,915 있는 170,749,396 했 다 160,559,206 할 수 149,405,973	
3	할 수 있 127,950,824 수 있는 64,442,701 고 있 다 64,262,223 고 있는 39,114,970 수 있 다 31,758,239	
4	할 수 있는 32,460,180 할 수 있 다 14,658,170 고 말 했 다 10,844,613 할 수 있 도록 10,310,465 뿐 만 아니 라 9,918,748	
5	동영상 보 기 영상 재생 4,321,282 는 것 이 중요 합니다 3,551,312 무단 전재 및 재 배포 2,812,200 전재 및 재 배포 금지 2,733,133 엔터 테 인 먼 트 2,056,941	

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				