

General overview

Corpus	Date	Language
hplt-v3-bjn_Arab	9/17/2025	Banjar (bjn)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
1,306	30,400	24,867 (81.80 %)	1.1M	4,582,094	7.88 MB

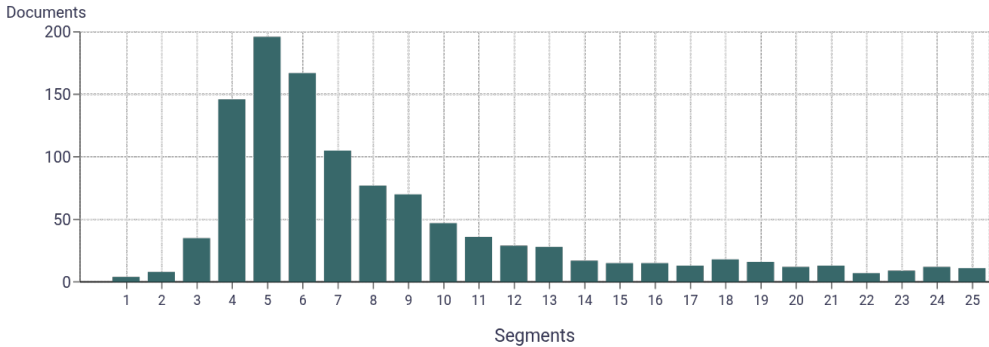
Top 10 domains

Domain	Docs	% of total
utusanmelayu.co...	662	50.69%
blogspot.com	153	11.72%
wordpress.com	147	11.26%
ahmadalikarim.com	54	4.13%
utusanv.com	33	2.53%
wikimedia.org	19	1.45%
harakahdaily.net	12	0.92%
ulamasedunia.org	11	0.84%
co.cc	11	0.84%
urusniaga.my	9	0.69%

Top 10 TLDs

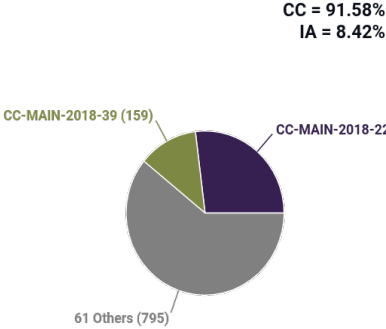
Domain	Docs	% of total
com.my	667	51.07%
com	475	36.37%
org	59	4.52%
my	30	2.30%
net	25	1.91%
cc	11	0.84%
moe	7	0.54%
org.my	6	0.46%
eu	4	0.31%
edu.my	4	0.31%

Documents size (in segments) ⓘ



≤ 25 segments **84.69%** (1.1K documents)  
> 25 segments **15.31%** (200 documents)

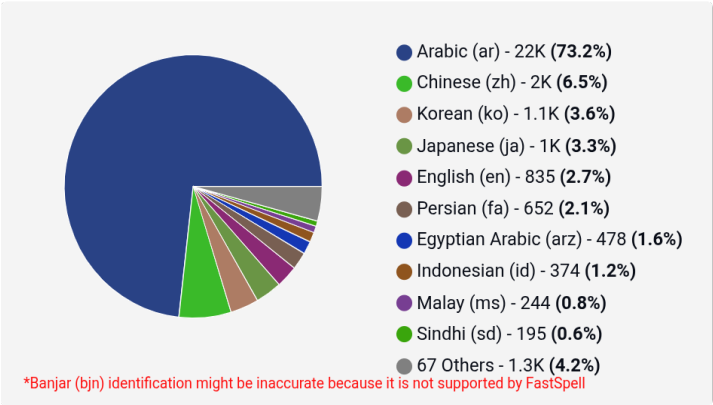
Document collections



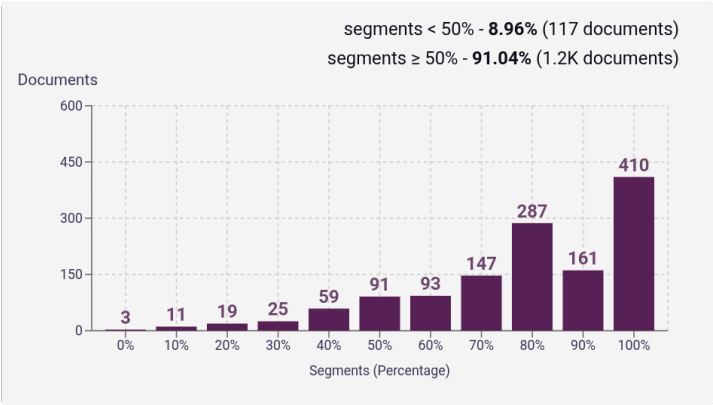
CC = **91.58%**  
IA = **8.42%**

Language Distribution

Number of segments in the Banjar (bjn) corpus



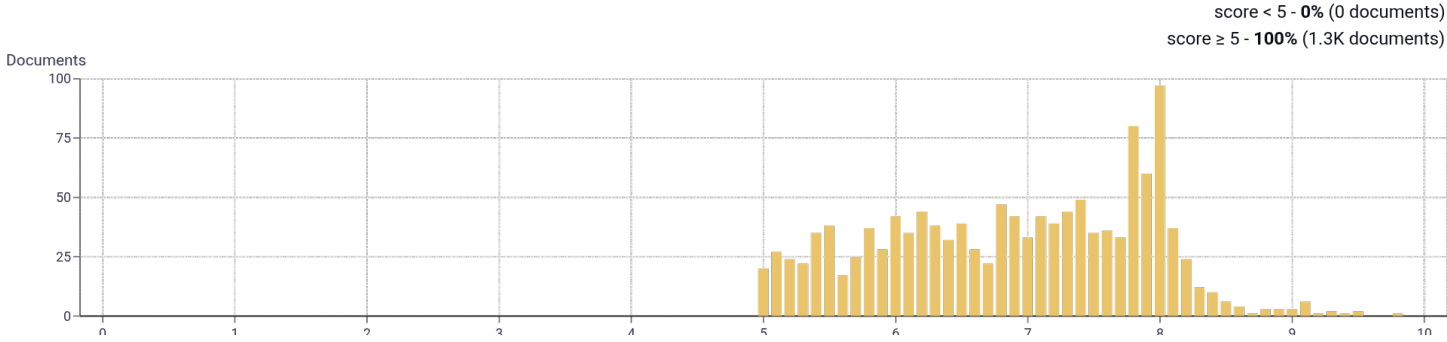
Percentage of segments in Banjar (bjn) inside documents



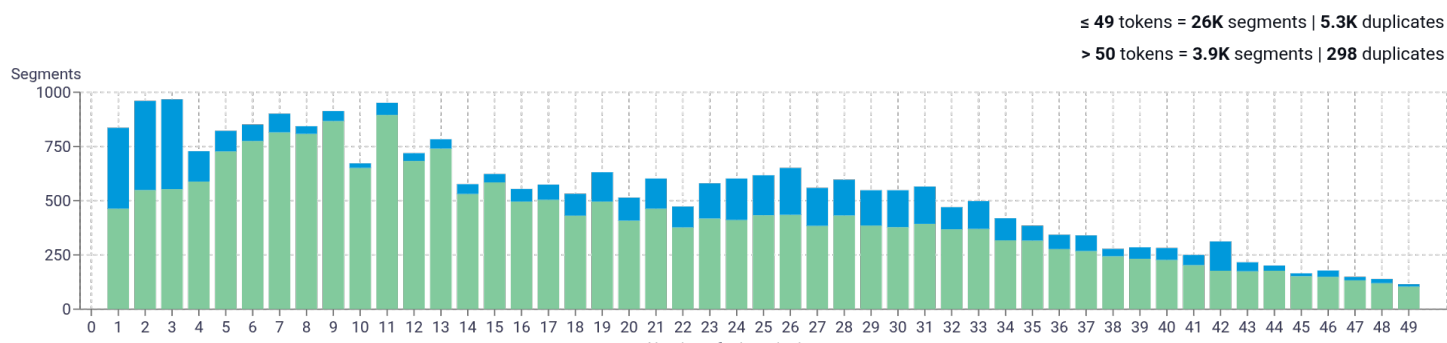
segments < 50% - **8.96%** (117 documents)  
segments ≥ 50% - **91.04%** (1.2K documents)

\*Banjar (bjn) identification might be inaccurate because it is not supported by FastSpell

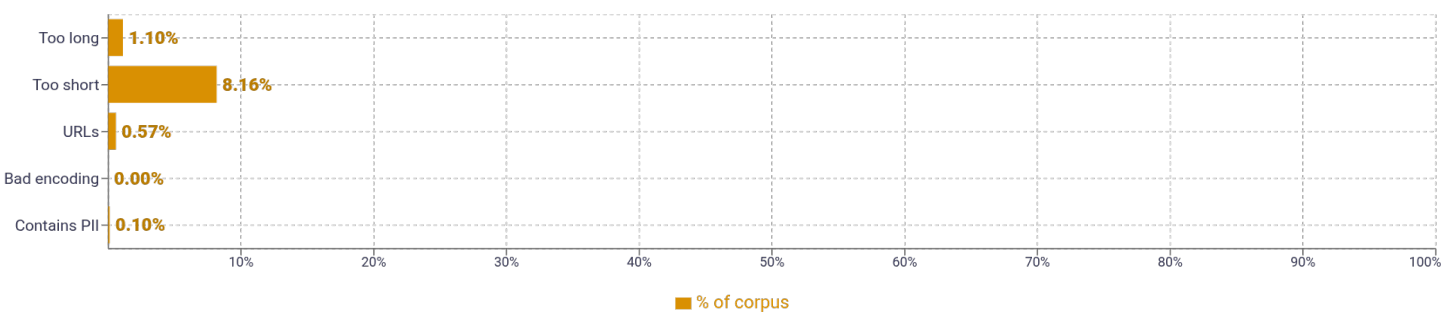
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	16,988   يغ   8,255   م   7,869   ن   7,212   و   6,597   الله	
2	arabic script   4,102   yue chinese   998   رسول الله   644   هاري اين   590   اورغ يغ   564	
3	صلى الله عليه   464   الله عليه وسلم   463   الله سبحانه وتعالى   399   الله صلى الله   183   رسول الله صلى   182	
4	صلى الله عليه وسلم   448   رسول الله صلى الله   179   الله صلى الله عليه   179   فنديديفن اسلم تيغكانن ساتو   123   سوكاتن بارو كريكولوم برسفادو   123	
5	رسول الله صلى الله عليه   175   محمد صلى الله عليه وسلم   112   رسول الله صلى الله عليه وسلم   170   رنخغن تاهونن فنديديفن اسلم تيغكانن   123   تاهونن فنديديفن اسلم تيغكانن ساتو   123	

# About HPLT Analytics

## Volumes - Segments

Segments correspond to paragraph and list boundaries as defined by HTML elements (<p>, <ul>, <ol>, 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  $\frac{\text{number of types (uniques)}}{\text{number of tokens}}$ , 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>, <ul>, <ol>, 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	ntp
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				