

General overview

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

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
21,227	363,963	284,756 (78.24 %)	12M	66,693,924	64.32 MB

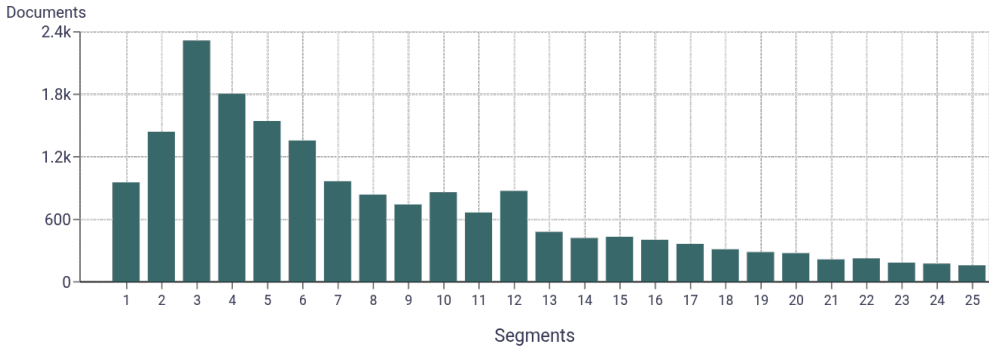
Top 10 domains

Domain	Docs	% of total
wikipedia.org	2.7K	12.48%
wordpress.com	1.1K	5.23%
blogspot.com	988	4.65%
tribunnews.com	640	3.02%
ruangguru.com	333	1.57%
banjarmasinbung...	302	1.42%
uinsby.ac.id	258	1.22%
web.app	255	1.20%
kanalkalimantan...	212	1.00%
d5d.org	183	0.86%

Top 10 TLDs

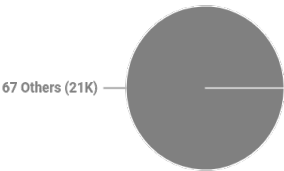
Domain	Docs	% of total
com	10K	47.89%
org	3.6K	17.18%
ac.id	2.2K	10.19%
co.id	731	3.44%
id	676	3.18%
net	514	2.42%
info	494	2.33%
go.id	482	2.27%
app	291	1.37%
co	215	1.01%

Documents size (in segments) ⓘ



≤ 25 segments **86.23%** (18K documents)
> 25 segments **13.77%** (2.9K documents)

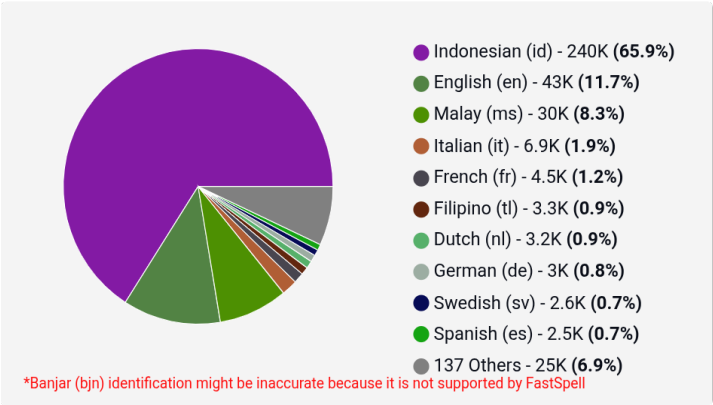
Document collections



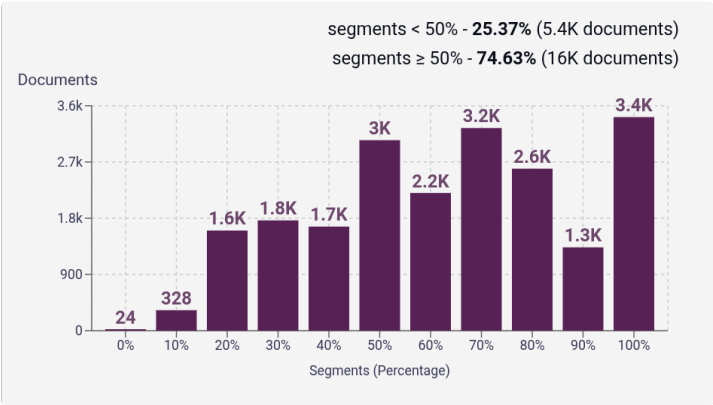
CC = 91.24%
IA = 8.76%

Language Distribution

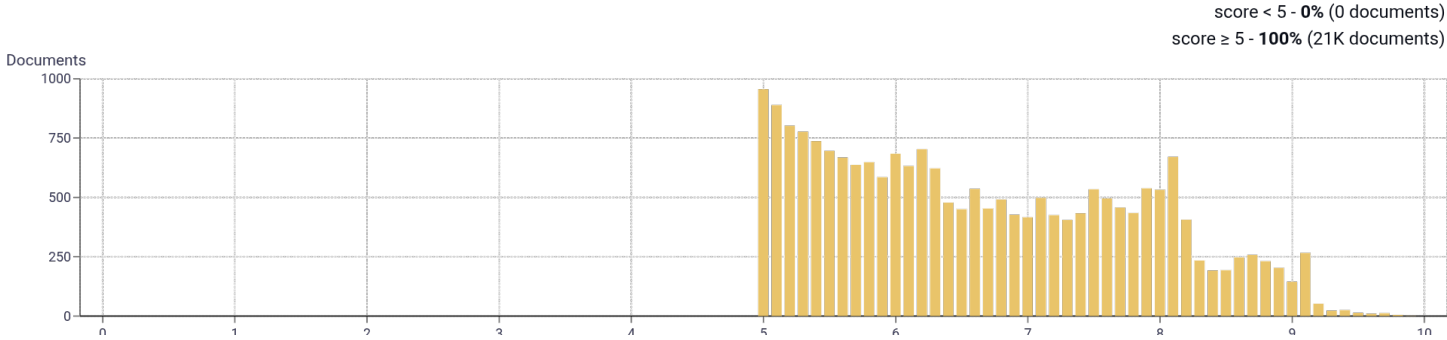
Number of segments in the Banjar (bjn) corpus



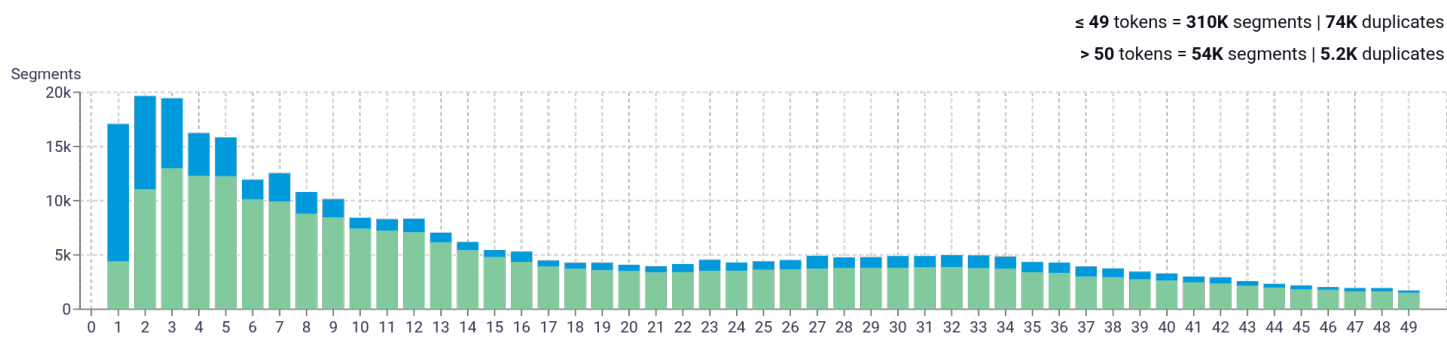
Percentage of segments in Banjar (bjn) inside documents



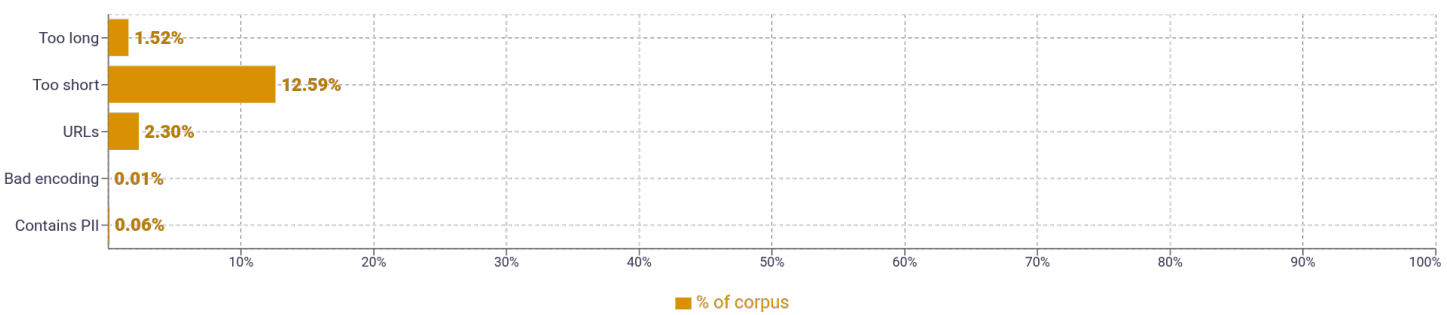
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	syair 57,021ban 36,765matriks 32,320surabaya 31,025sunan 30,058	
2	uin sunan 25,448sunan ampel 22,698ampel surabaya 21,698undergraduate thesis 20,573syair hk 8,186	
3	sunan ampel surabaya 21,674uin sunan ampel 19,931uin sunan kalijaga 5,360sunan kalijaga yogyakarta 4,212jarak mill laut 2,232	
4	uin sunan ampel surabaya 19,605uin sunan kalijaga yogyakarta 4,023iain sunan ampel surabaya 1,686fc liverpool fc liverpool 1,317liverpool fc liverpool fc 1,306	
5	liverpool fc liverpool fc liverpool 1,295fc liverpool fc liverpool fc 1,203mill jarak dari banjarmasin to 1,173berapa mill jarak dari banjarmasin 1,141jarak dari banjarmasin to halmahera 1,003	

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 $\frac{\text{number of types (uniques)}}{\text{number of tokens}}$, 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>, , , 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				