

General overview

Corpus	Date	Language
hplt-v3-bug_Latn	9/17/2025	Buginese (bug)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
1,173	32,290	28,275 (87.57 %)	1.7M	8,599,225	8.79 MB

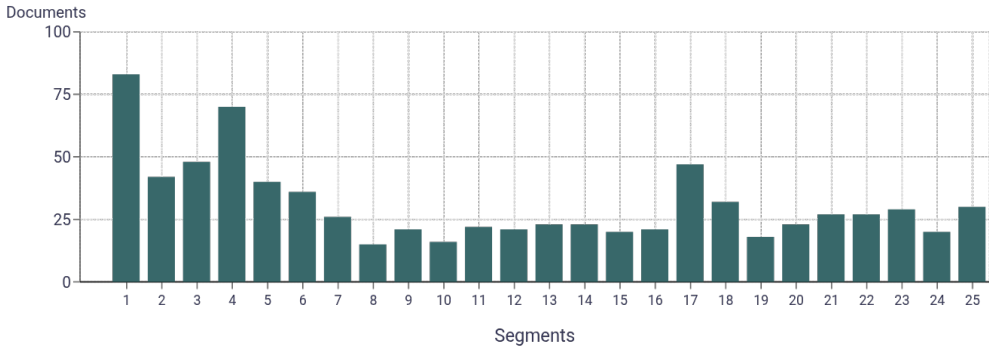
Top 10 domains

Domain	Docs	% of total
alkitab.mobi	535	45.61%
wordpress.com	110	9.38%
bible.is	46	3.92%
blogspot.com	37	3.15%
wikipedia.org	24	2.05%
indonesiachord.com	21	1.79%
chordpass.com	19	1.62%
basasulselwiki.org	17	1.45%
wikimedia.org	15	1.28%
ebible.org	15	1.28%

Top 10 TLDs

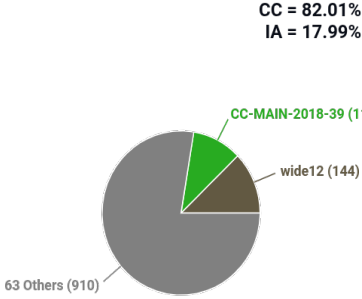
Domain	Docs	% of total
mobi	535	45.61%
com	369	31.46%
org	88	7.50%
is	46	3.92%
net	36	3.07%
id	25	2.13%
info	10	0.85%
pw	6	0.51%
de	6	0.51%
me	5	0.43%

Documents size (in segments) ⓘ



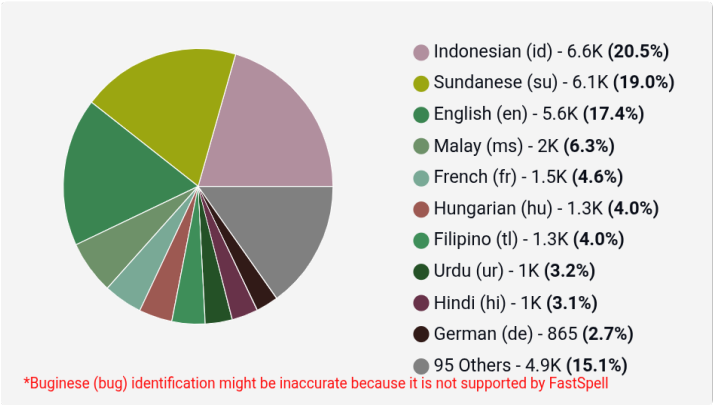
≤ 25 segments **66.5%** (780 documents)  
> 25 segments **33.5%** (393 documents)

Document collections

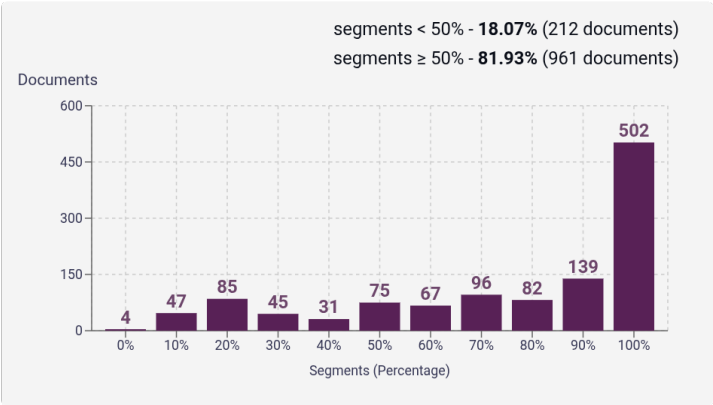


Language Distribution

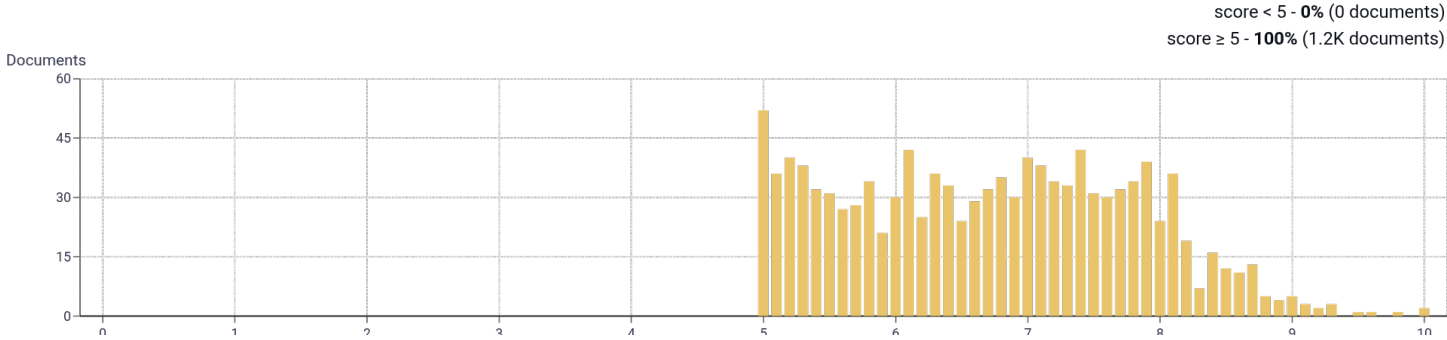
Number of segments in the Buginese (bug) corpus



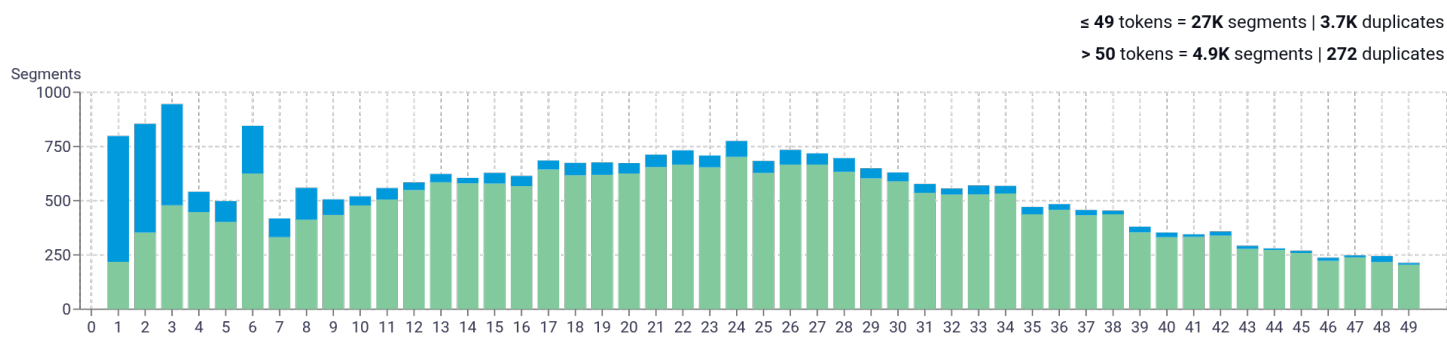
Percentage of segments in Buginese (bug) inside documents



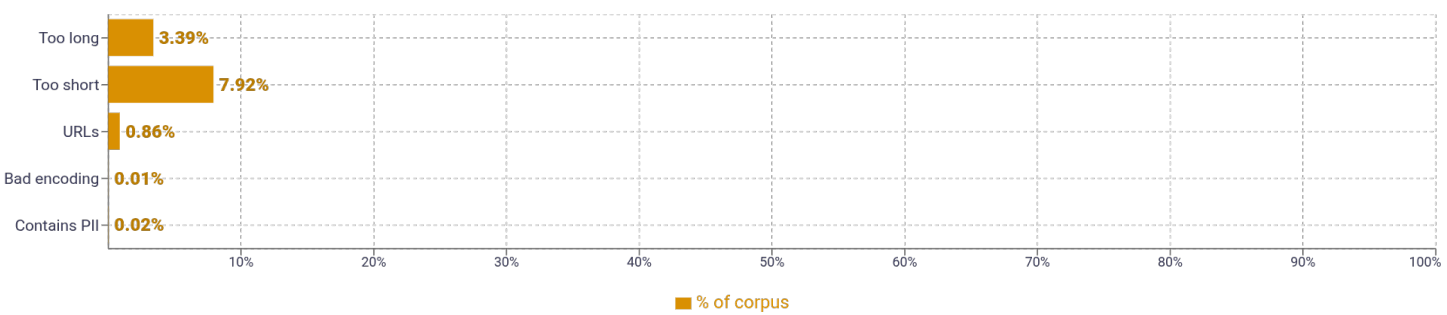
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	<div>sibawa   14,108i   8,080lao   8,014anjo   7,899tpi   7,082</div>	
2	<div>anjo taua   503kamma anne   497trans tv   426sibawa sining   388iyaro sining   376</div>	
3	<div>lao ri mennang   568lao ri iko   514lao ri iyya   327mae ri allata   295lao ri puwangnge   254</div>	
4	<div>mae ri kau ngaseng   173puwangnge lao ri musa   101lalang ri kittaka angkanaya   98yèsus lao ri mennang   88yèsus lao ri mennang   84</div>	
5	<div>nakkeda yèsus lao ri mennang   48name necklace with rhinestone letters   35custom name necklace with rhinestone   35moga saya tidak kena kutuk   34umma selleng malebbi engkae riamasei   33</div>	

# 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				