

General overview

Corpus	Analytics date	Language
ban_Latn.jsonl.tsv	10/3/2024	Balinese (ban)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
10,700	601,139	161,172 (26.81 %)	14M	73.75 MB	76,653,907

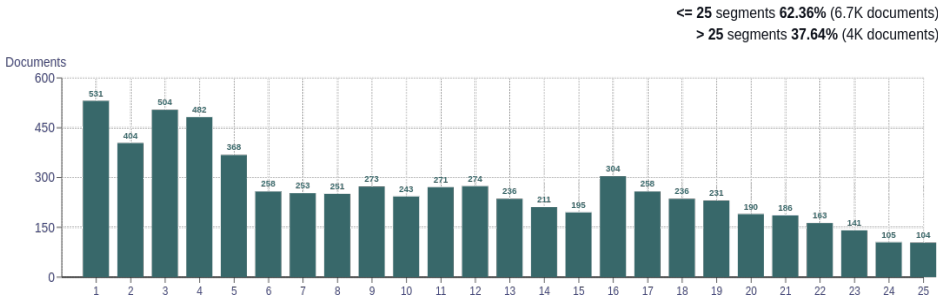
Top 10 domains

Domain	Docs	% of total
basabali.org	3.1K	28.95
wikipedia.org	1.6K	14.59
blogspot.com	1.5K	13.56
bible.is	674	6.30
wordpress.com	450	4.21
blogspot.co.id	409	3.82
alkitab.mobi	220	2.06
suarasakingballi.com	143	1.34
belajarbahasabali.com	123	1.15
scribd.com	90	0.84

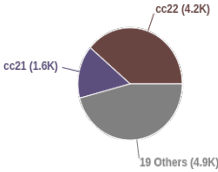
Top 10 TLDs

Domain	Docs	% of total
org	5K	46.33
com	3.6K	33.89
is	674	6.30
co.id	429	4.01
mobi	220	2.06
net	166	1.55
ac.id	56	0.52
asia	36	0.34
in	33	0.31
nl	31	0.29

Documents size (in segments)

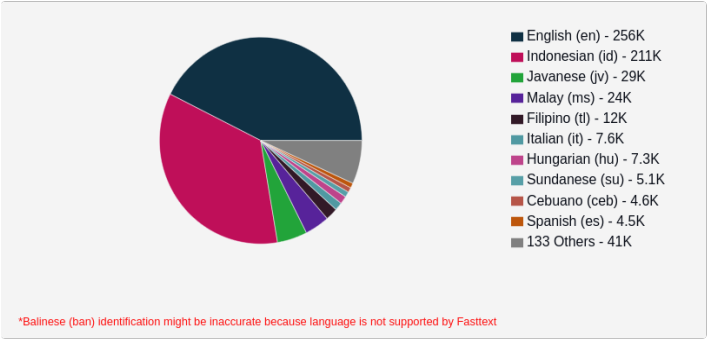


Documents by collection

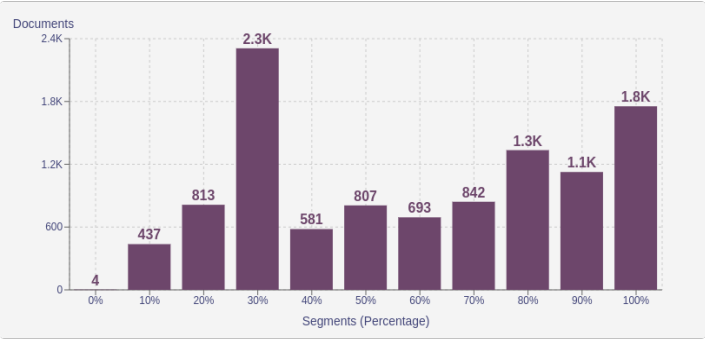


Language Distribution

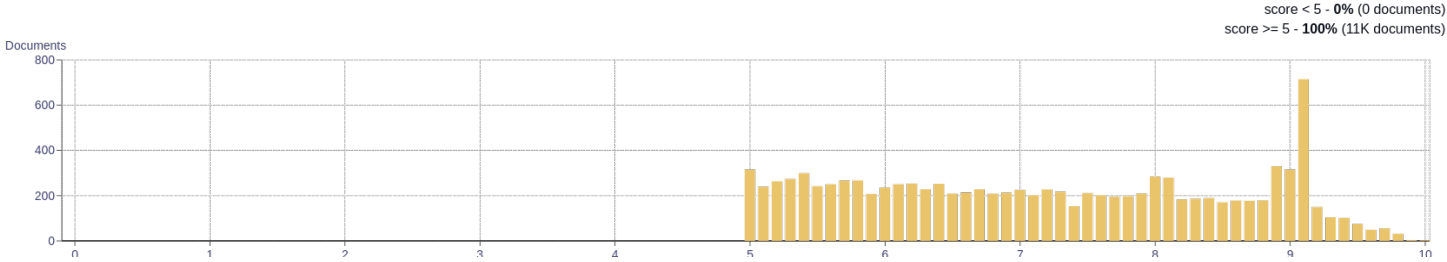
Number of segments



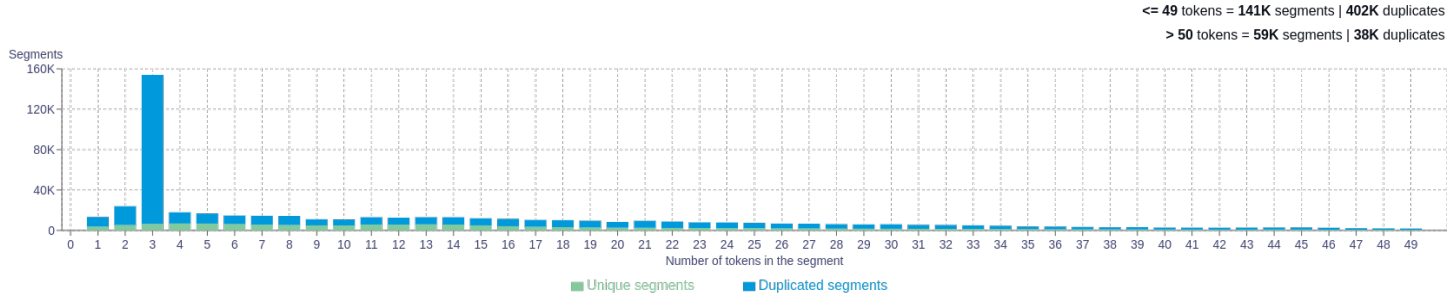
Percentage of segments in Balinese (ban) inside documents



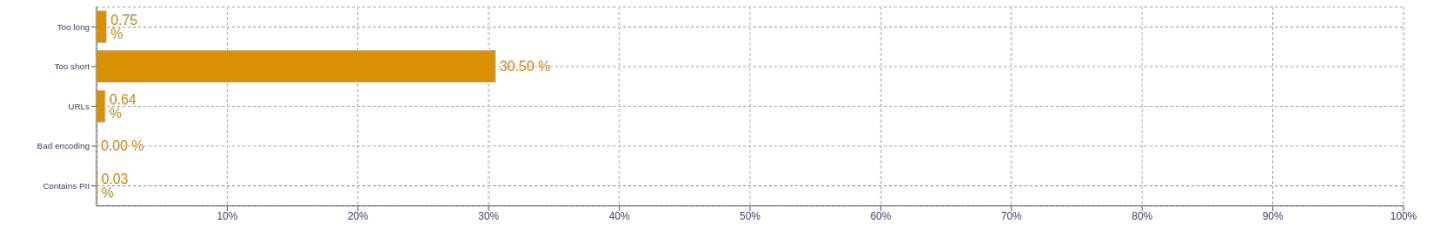
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>in 303017</div> <div>balinese 99435</div> <div>indonesian 97372</div> <div>english 97267</div> <div>bali 80618</div>
2	<div>in balinese 97070</div> <div>in english 96934</div> <div>in indonesian 96516</div> <div>hyang widi 12563</div> <div>widi wasa 12169</div>
3	<div>hyang widi wasa 12122</div> <div>usage examples pulled 1971</div> <div>pulled from the 1971</div> <div>examples pulled from 1971</div> <div>duwur nyane susunin 1730</div>
4	<div>usage examples pulled from 1971</div> <div>examples pulled from the 1971</div> <div>pulled from the virtual 1323</div> <div>from the virtual library 1318</div> <div>susunin antuk alad sesayut 1244</div>
5	<div>usage examples pulled from the 1971</div> <div>examples pulled from the virtual 1323</div> <div>pulled from the virtual library 1318</div> <div>nyane susunin antuk alad sesayut 1237</div> <div>duwur nyane susunin antuk alad 1237</div>

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.stinfo.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>