

General overview

Corpus	Analytics date	Language
bjn_Latn.jsonl.tsv	11/27/2024	Banjar (bjn)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
18,764	366,336	154,702 (42.23 %)	10M	53.33 MB	55,627,091

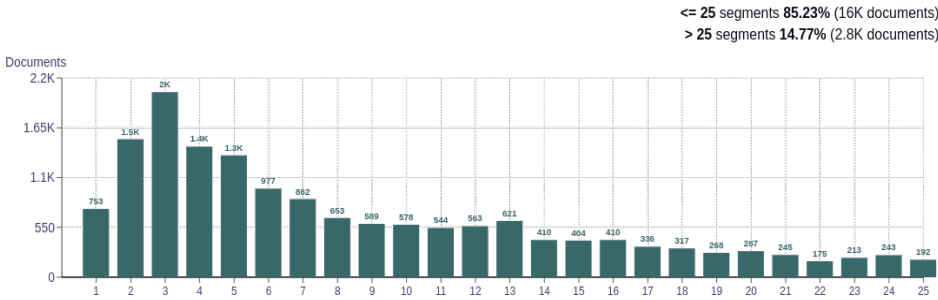
Top 10 domains

Domain	Docs	% of total
wikipedia.org	7.1K	37.59
blogspot.com	1.4K	7.62
wordpress.com	988	5.27
banjarmasinbungas.com	637	3.39
bible.is	535	2.85
sciencegraph.net	397	2.12
petalokasi.org	244	1.30
tribunnews.com	225	1.20
forvo.com	180	0.96
blogspot.co.id	167	0.89

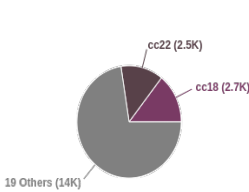
Top 10 TLDs

Domain	Docs	% of total
org	7.6K	40.67
com	7.4K	39.68
net	721	3.84
is	535	2.85
co.id	462	2.46
ac.id	292	1.56
go.id	181	0.96
id	177	0.94
asia	164	0.87
info	139	0.74

Documents size (in segments)

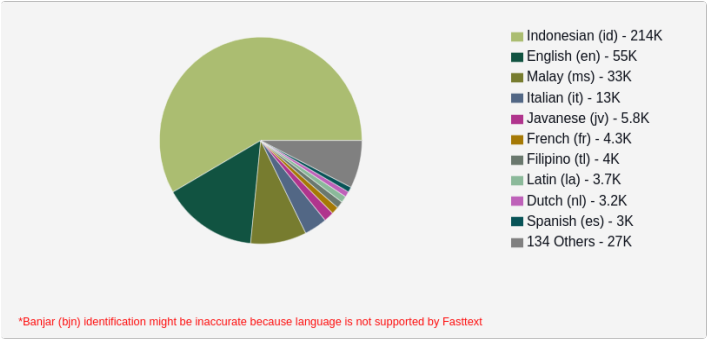


Documents by collection

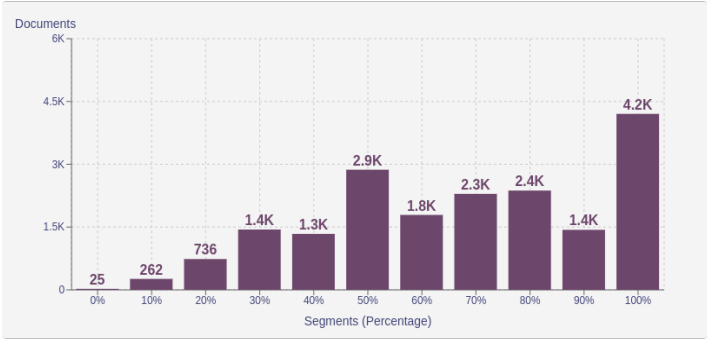


Language Distribution

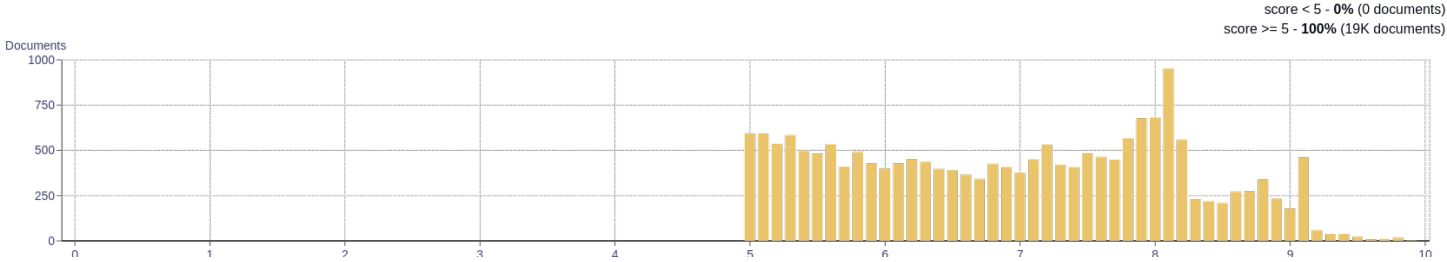
Number of segments



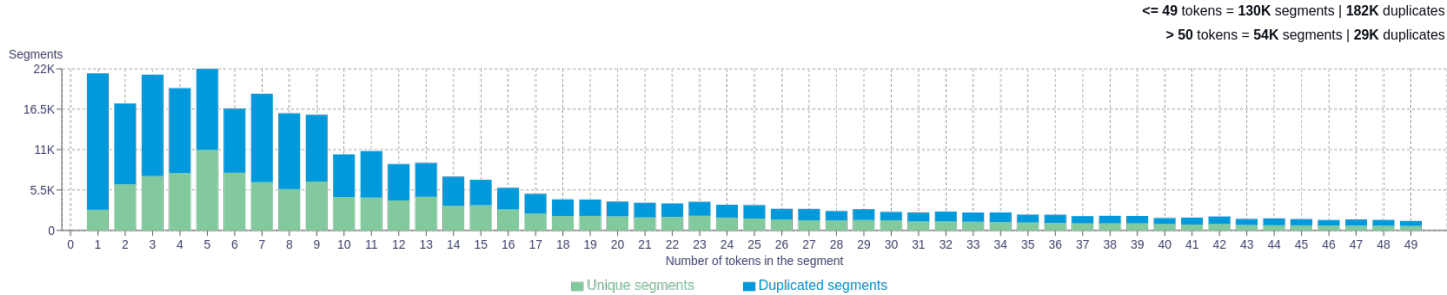
Percentage of segments in Banjar (bjn) inside documents



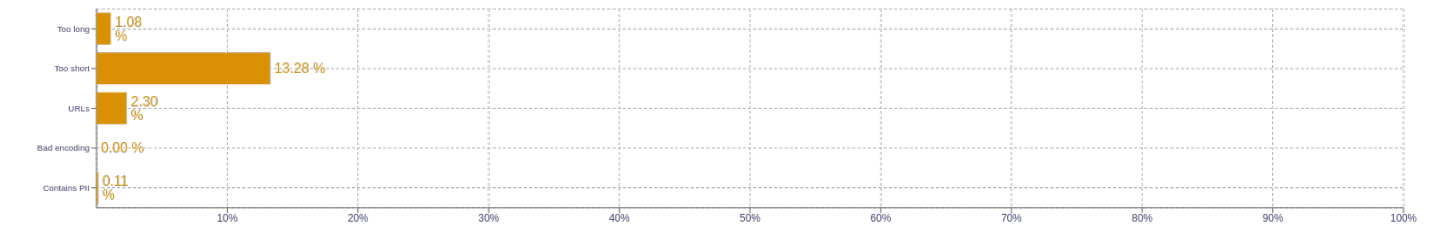
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>arjuna   28581</div> <div>banjarmasin   28345</div> <div>hotel   20187</div> <div>je   19277</div> <div>caps   18504</div>
2	<div>warfarin arjuna   8435</div> <div>caps arjuna   8399</div> <div>sunting sumber   8337</div> <div>caps warfarin   4921</div> <div>warfarin warfarin   4891</div>
3	<div>caps warfarin arjuna   2411</div> <div>warfarin warfarin arjuna   2309</div> <div>jual rok tutu   1714</div> <div>tutu di banjarmasin   1712</div> <div>caps arjuna warfarin   1677</div>
4	<div>rok tutu di banjarmasin   1712</div> <div>grosir jual rok tutu   1284</div> <div>tutu di banjarmasin murah   1070</div> <div>caps warfarin arjuna warfarin   846</div> <div>caps arjuna warfarin arjuna   833</div>
5	<div>jual rok tutu di banjarmasin   1712</div> <div>rok tutu di banjarmasin murah   1070</div> <div>hotel hotel hotel hotel hotel   810</div> <div>caps warfarin arjuna warfarin arjuna   410</div> <div>warfarin warfarin arjuna warfarin arjuna   400</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 "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>, <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>