

General overview

Corpus	Analytics date	Language
dyu_Latn.jsonl.tsv	11/27/2024	Dyula (dyu)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
1,390	24,558	20,698 (84.28 %)	1.5M	5.7 MB	5,529,102

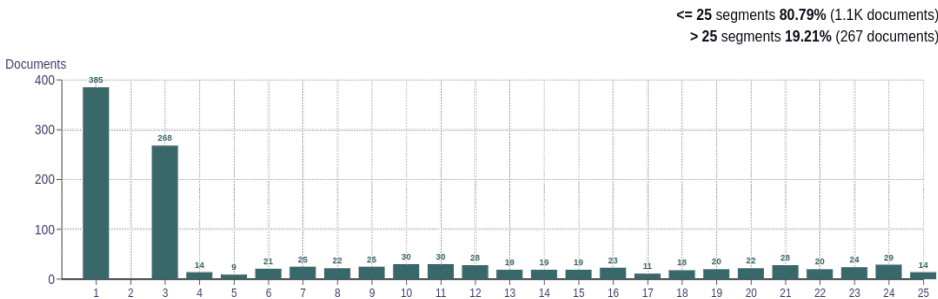
Top 10 domains

Domain	Docs	% of total
bible.is	646	46.47
bibles.org	431	31.01
jw.org	299	21.51
omniglot.com	4	0.29
bible.com	3	0.22
watchtower.org	3	0.22
gospelgo.com	2	0.14
twr360.org	1	0.07
reunion.com	1	0.07

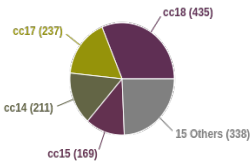
Top 10 TLDs

Domain	Docs	% of total
org	734	52.81
is	646	46.47
com	10	0.72

Documents size (in segments)

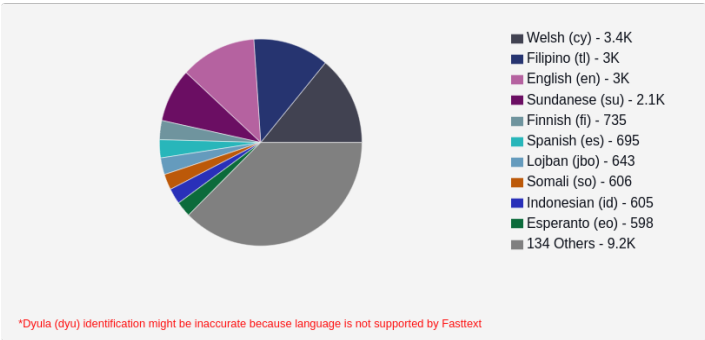


Documents by collection

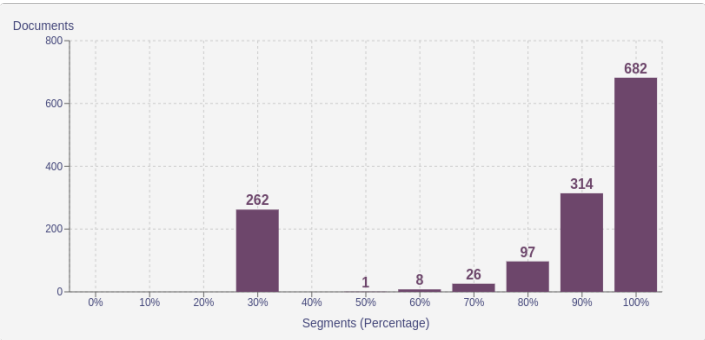


Language Distribution

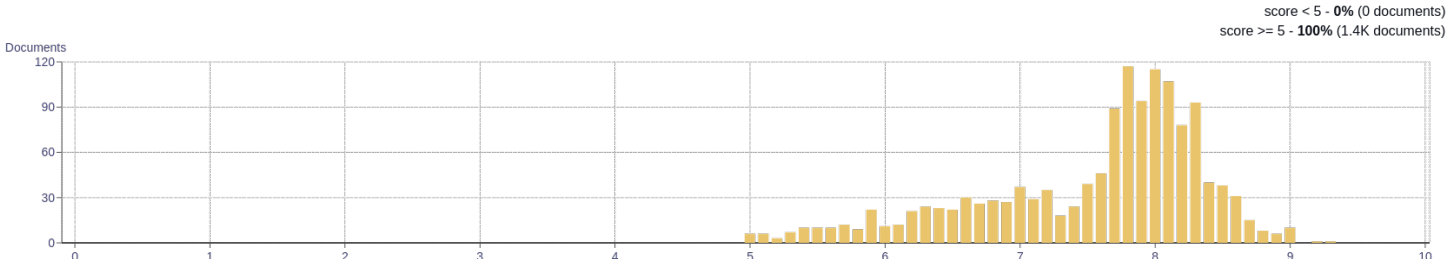
Number of segments



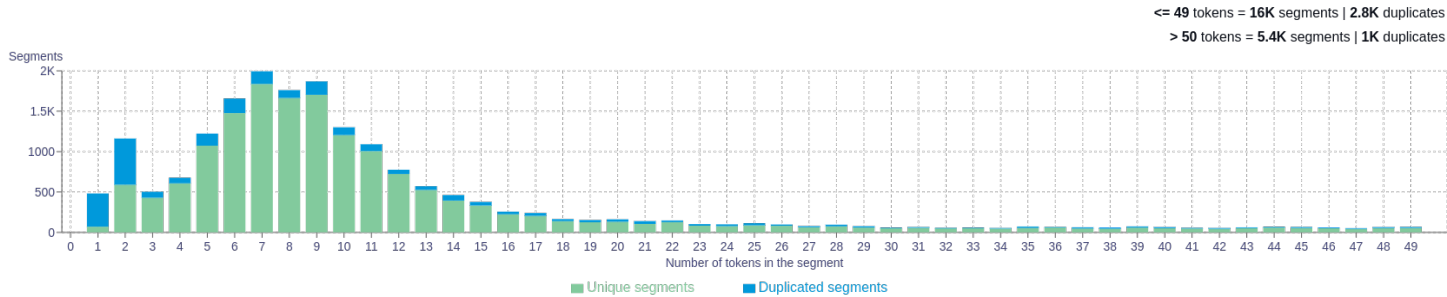
Percentage of segments in Dyula (dyu) inside documents



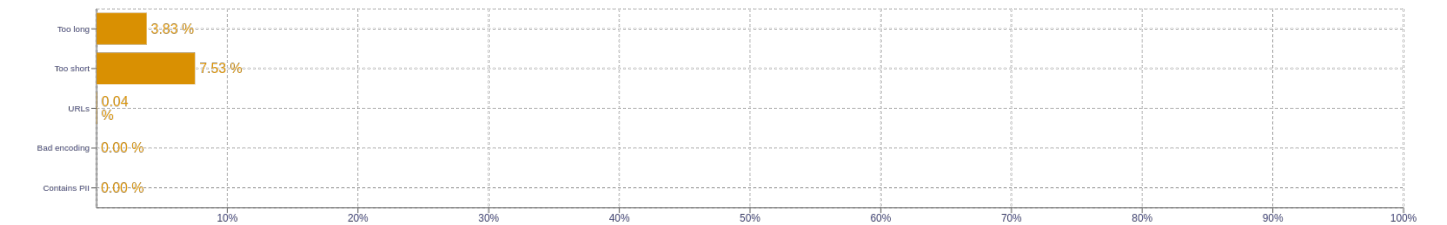
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div><div>o   29654</div><div>u   26214</div><div>be   24276</div><div>ko   22739</div><div>k   19305</div></div>
2	<div><div>tun be   4211</div><div>tuma min   2018</div><div>o tigi   1861</div><div>cogo min   1771</div><div>min be   1713</div></div>
3	<div><div>be se k   824</div><div>u ye ko   802</div><div>ala ka kuma   789</div><div>ala ka masaya   482</div><div>aw ye ko   477</div></div>
4	<div><div>fɔra ala ka kuma   267</div><div>ala ka mɔɔ woɓomanin   151</div><div>masaba aw ka ala   148</div><div>masaba le ko ten   112</div><div>mine ka taga n   100</div></div>
5	<div><div>ala ka kuma na ko   183</div><div>ala ka mɔɔ woɓomanin nin   99</div><div>ne masaba le ko ten   55</div><div>aw ye ko ni mɔɔ   52</div><div>o yɔɔ la ka taga   51</div></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>