

General overview

Corpus	Analytics date	Language
taq_Latn.jsonl.tsv	11/28/2024	Tamasheq (taq)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
1,747	13,884	7,842 (56.48 %)	2.3M	9.41 MB	8,833,192

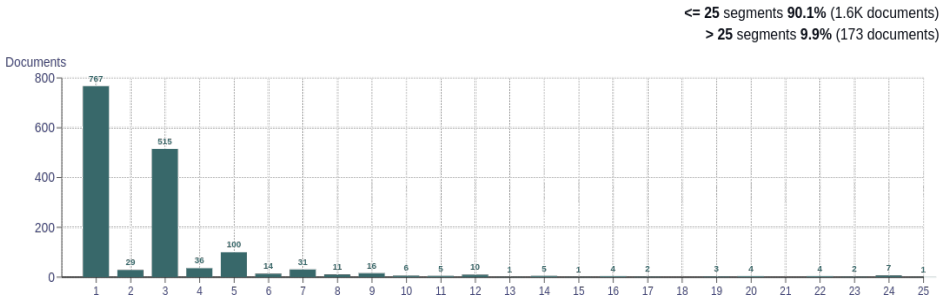
Top 10 domains

Domain	Docs	% of total
bible.is	1.2K	69.83
worddetector.com	79	4.52
newchristianbiblestudy.org	60	3.43
ebible.org	53	3.03
biblehub.com	18	1.03
vanuatubibles.org	16	0.92
tuspalabras.com	14	0.80
blogspot.com	12	0.69
case.edu	9	0.52
omniglot.com	9	0.52

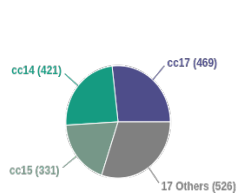
Top 10 TLDs

Domain	Docs	% of total
is	1.2K	69.83
com	250	14.31
org	162	9.27
net	33	1.89
edu	9	0.52
de	8	0.46
co	8	0.46
es	7	0.40
com.pl	4	0.23
me	4	0.23

Documents size (in segments)

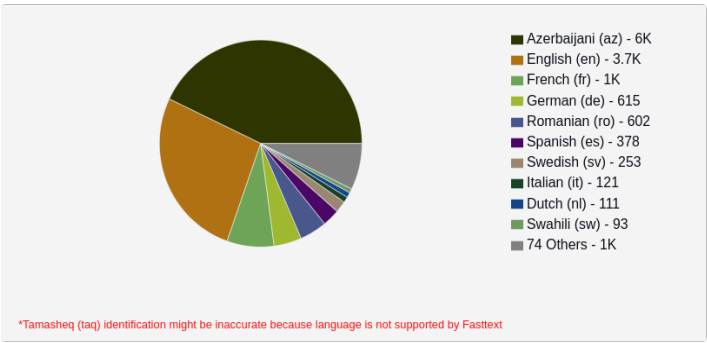


Documents by collection

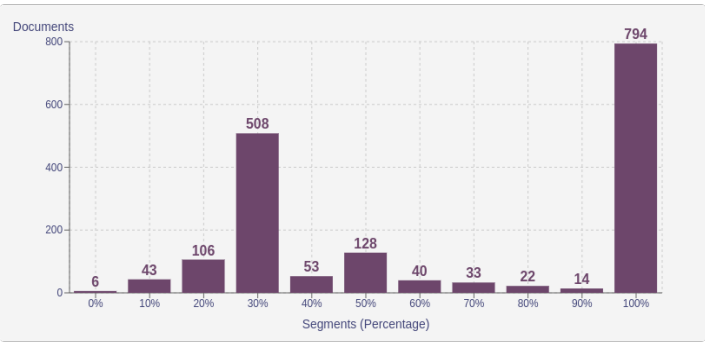


Language Distribution

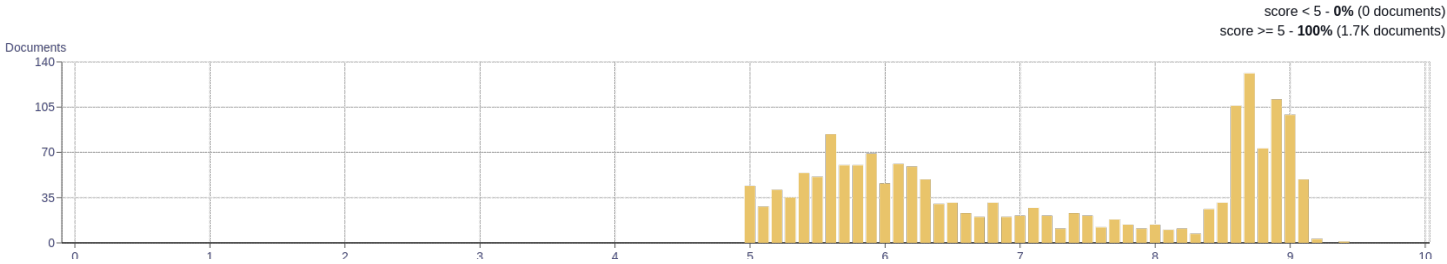
Number of segments



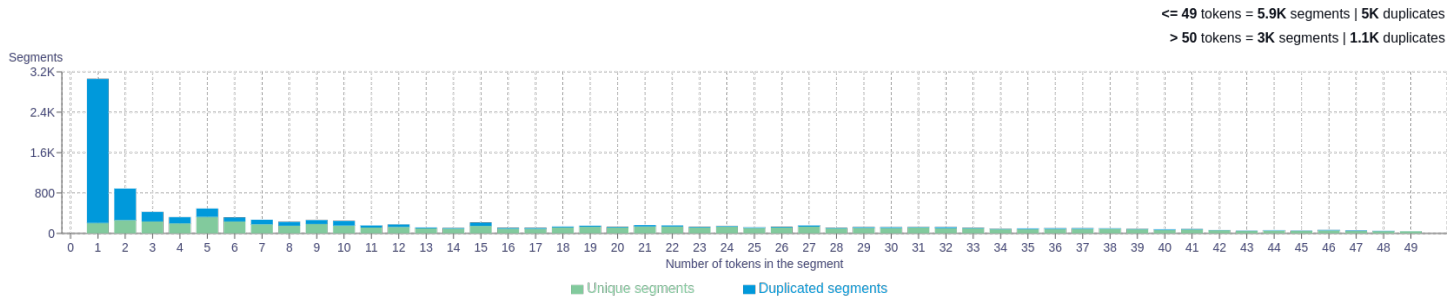
Percentage of segments in Tamasheq (taq) inside documents



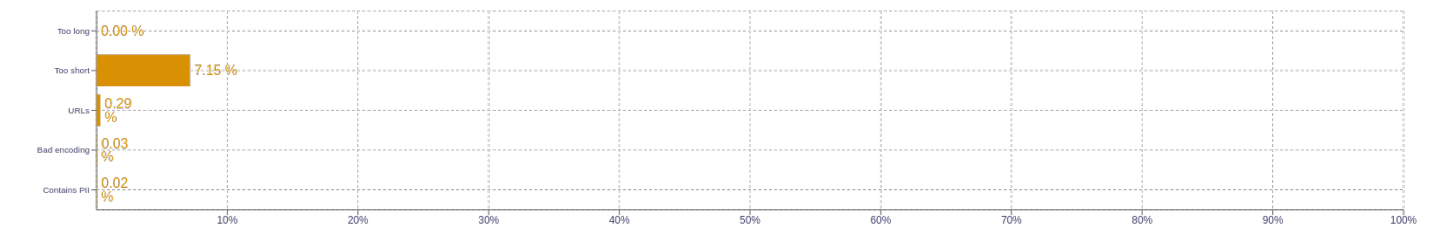
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>n 29099</div> <div>s 23680</div> <div>dåg 22105</div> <div>ahay 21633</div> <div>sə 20455</div>
2	<div>win win 5792</div> <div>ata awan 4379</div> <div>do ahay 3716</div> <div>ata nå 3207</div> <div>do sə 3019</div>
3	<div>win win win 5787</div> <div>mer su way 2029</div> <div>anà do ahay 1261</div> <div>asd asd asd 1200</div> <div>sdjflk asdfkas df 990</div>
4	<div>win win win win 5782</div> <div>laskdj flka sdjflk asdfkas 990</div> <div>flka sdjflk asdfkas df 990</div> <div>d faslkdfj asdlkfj as 990</div> <div>asdf aasd f asdklfj 990</div>
5	<div>win win win win win 5777</div> <div>laskdj flka sdjflk asdfkas df 990</div> <div>asdf aasd f asdklfj as 990</div> <div>flsadf asdf asdf aasd f asdklfj 986</div> <div>aksd flsadf asdf asdf aasd f 986</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.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>