

General overview

Corpus	Analytics date	Language
gui_gujr.jsonl.tsv	9/16/2024	Gujarati (gu)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
1,134,252	20,639,718	11,424,183 (55.35 %)	667M	7.99 GB	3,366,421,654

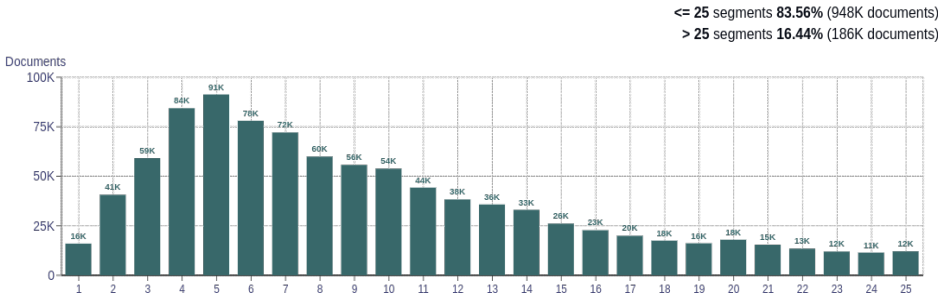
Top 10 domains

Domain	Docs	% of total
divyabhaskar.co.in	83K	7.35
wordpress.com	48K	4.25
news18.com	44K	3.92
oneindia.com	29K	2.58
sandesh.com	27K	2.42
wikipedia.org	24K	2.16
gujurocks.in	23K	2.02
chitralekha.com	18K	1.58
webdunia.com	17K	1.48
vtvgujarati.com	16K	1.38

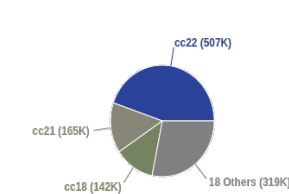
Top 10 TLDs

Domain	Docs	% of total
com	727K	64.10
in	172K	15.18
co.in	101K	8.92
org	76K	6.69
net	19K	1.67
news	4.7K	0.41
app	3K	0.26
online	2.8K	0.24
gov.in	2.4K	0.21
live	1.4K	0.12

Documents size (in segments)

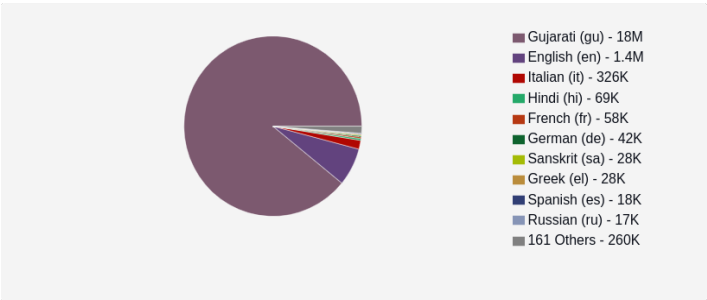


Documents by collection

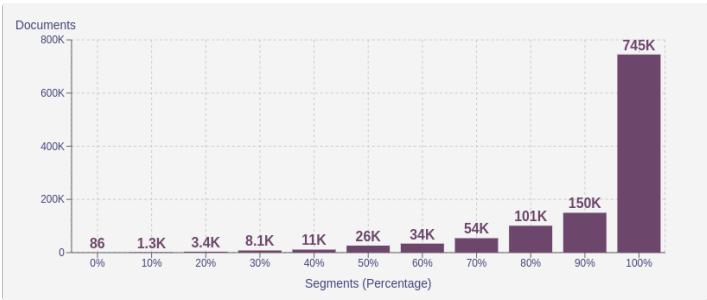


Language Distribution

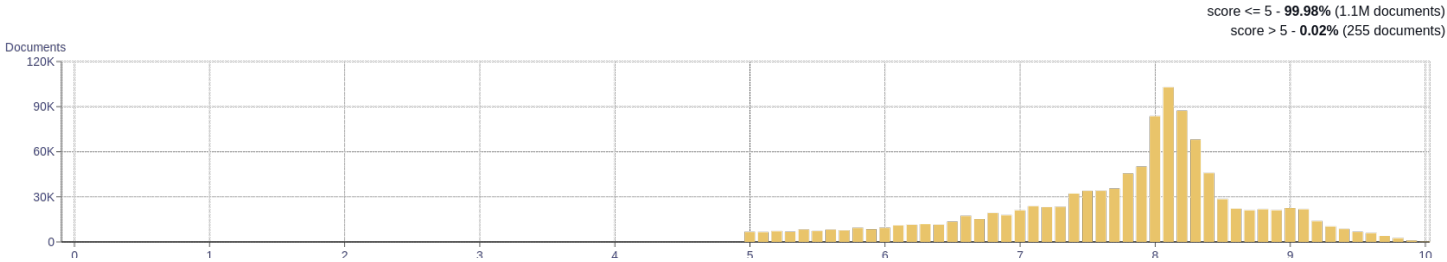
Number of segments



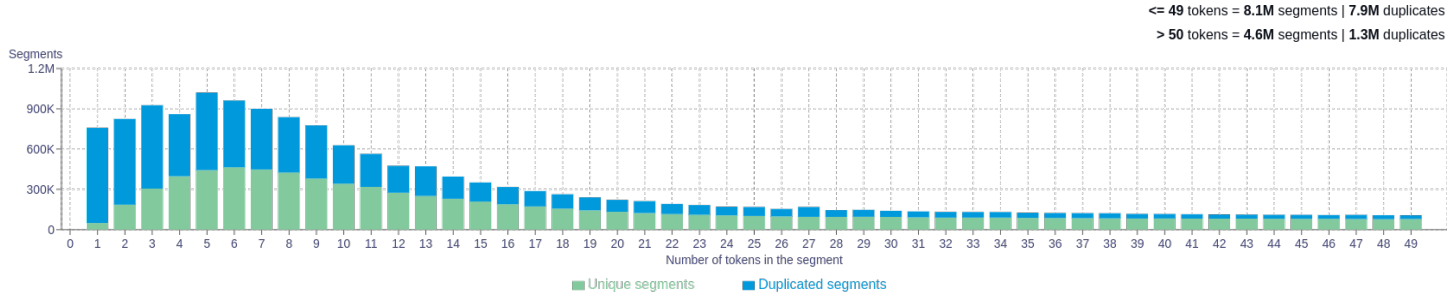
Percentage of segments in Gujarati (gu) inside documents



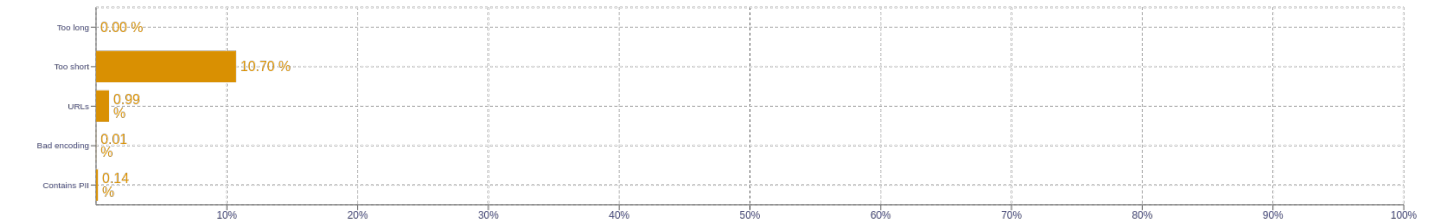
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>સર્વો 2721183</div> <div>દરો 1796456</div> <div>સરવાળો 1275130</div> <div>સરવા 1137108</div> <div>દ્વારા 1131630</div>
2	<div>ફેરફાર કરો 177866</div> <div>સામગ્રી કરો 143408</div> <div>નિંદા મળે 130076</div> <div>સરવાળો સમજાવે 122479</div> <div>સરવાળો સમજાવે 119799</div>
3	<div>all rights reserved 46726</div> <div>db corp ltd 46044</div> <div>code of ethics 46018</div> <div>website follows the 46017</div> <div>this website follows 46017</div>
4	<div>website follows the dnpa 46017</div> <div>this website follows the 46017</div> <div>the dnpa code of 46017</div> <div>follows the dnpa code 46017</div> <div>dnpa code of ethics 46017</div>
5	<div>website follows the dnpa code 46017</div> <div>this website follows the dnpa 46017</div> <div>the dnpa code of ethics 46017</div> <div>follows the dnpa code of 46017</div> <div>સર્વોત્તમ વસ્તુ સમાવેશ કરો -પૃષ્ઠ18 20648</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>