

General overview

Corpus	Analytics date	Language
mar_Deva.jsonl.tsv	9/23/2024	Marathi (mr)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
2,079,680	36,315,277	20,571,649 (56.65 %)	1.2B	16.02 GB	6,587,873,834

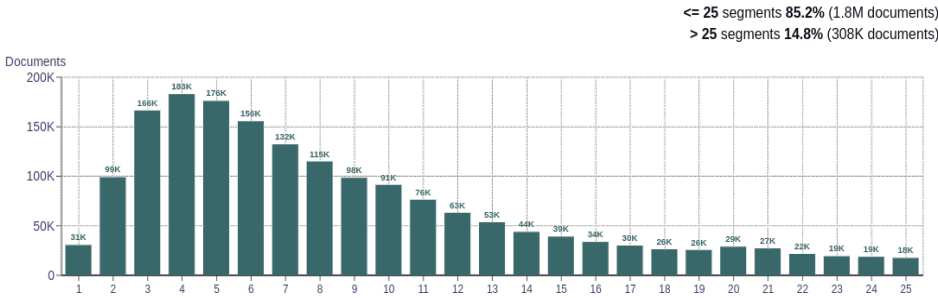
Top 10 domains

Domain	Docs	% of total
blogspot.in	82K	3.95
bhaskar.com	80K	3.85
esakal.com	79K	3.82
loksatta.com	67K	3.23
blogspot.com	66K	3.19
wikipedia.org	57K	2.74
news18.com	56K	2.71
indiatimes.com	37K	1.79
tv9marathi.com	33K	1.60
ibnlokmat.tv	28K	1.36

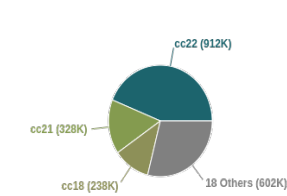
Top 10 TLDs

Domain	Docs	% of total
com	1.4M	67.81
in	330K	15.89
org	122K	5.85
tv	31K	1.49
gov.in	24K	1.16
co.in	24K	1.16
net	24K	1.13
news	21K	1.00
page	19K	0.89
es	12K	0.58

Documents size (in segments)

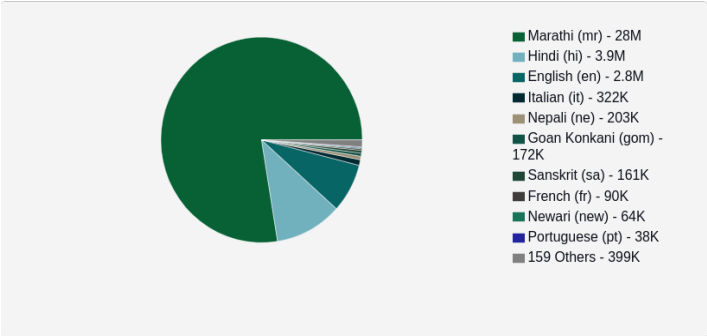


Documents by collection

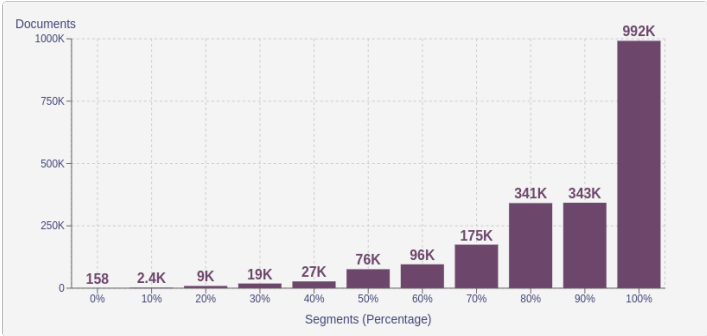


Language Distribution

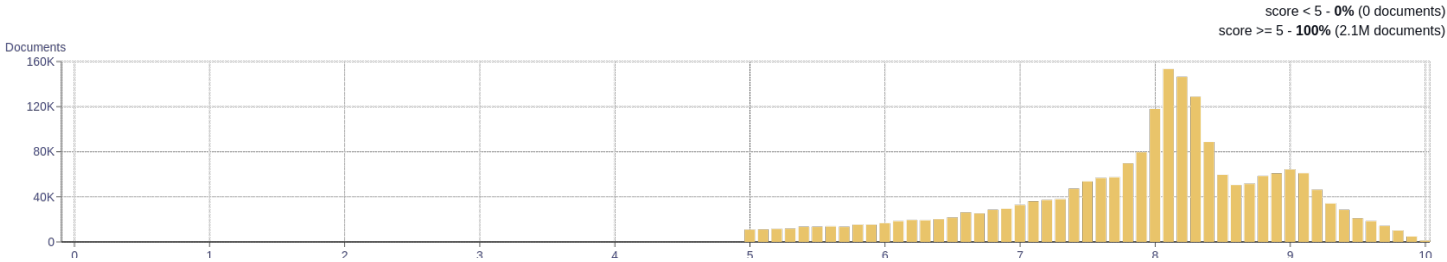
Number of segments



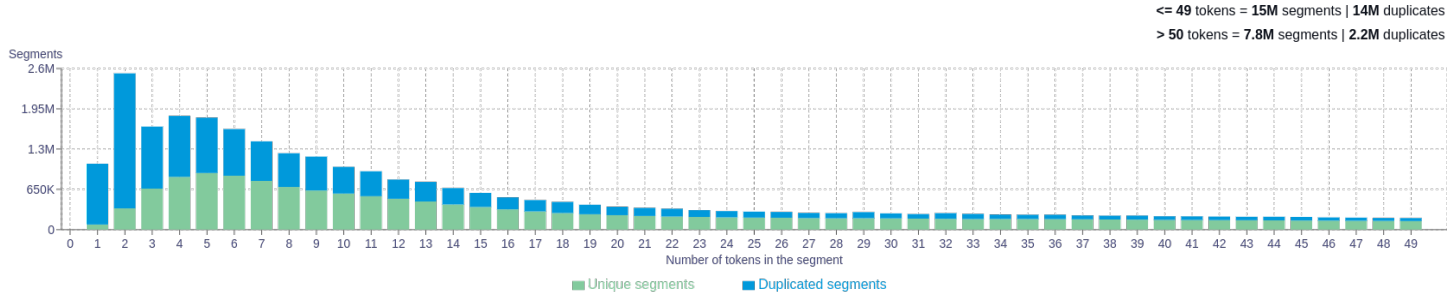
Percentage of segments in Marathi (mr) inside documents



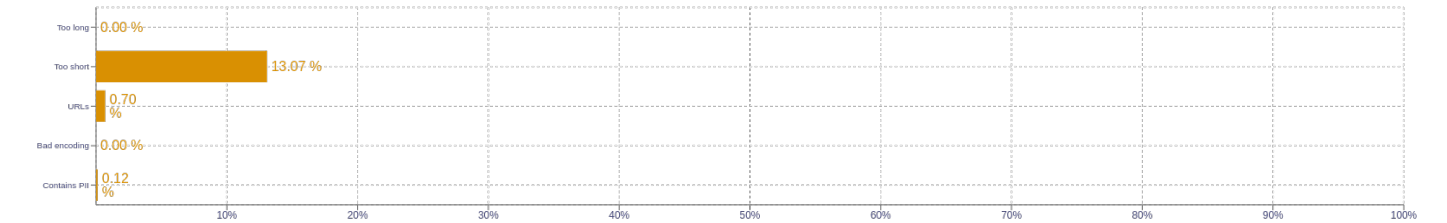
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>यांनी   2721206</div> <div>कारण्यात   1843409</div> <div>रिग्या   1820179</div> <div>त्यांनी   1588940</div> <div>आमण   1551143</div>
2	<div>in marathi   213206</div> <div>मोठ्या प्रमाणात   153558</div> <div>करू शकता   134827</div> <div>पुन्हा एकदा   126837</div> <div>यांनी सांगितले   124606</div>
3	<div>पंतप्रधान नरेंद्र मोदी   64747</div> <div>आम्हाला पोलो करा   50515</div> <div>all rights reserved   48853</div> <div>मुख्यमंत्री उद्धव ठाकरे   42998</div> <div>शहरातील ताज्या बातम्या   42286</div>
4	<div>ताज्या बातम्या आणि ई   42283</div> <div>this website follows the   42117</div> <div>website follows the dnpa   42115</div> <div>the dnpa code of   42113</div> <div>follows the dnpa code   42113</div>
5	<div>शहरातील ताज्या बातम्या आणि ई   42283</div> <div>this website follows the dnpa   42115</div> <div>website follows the dnpa code   42113</div> <div>the dnpa code of ethics   42113</div> <div>follows the dnpa code of   42113</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 or 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>