

General overview

Corpus	Date	Language
mal_Mlym.jsonl.tsv	9/21/2024	Malayalam (ml)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
3,104,759	48,003,517	24,484,289 (51.01 %)	1.2B	9,443,613,087	23.7 GB

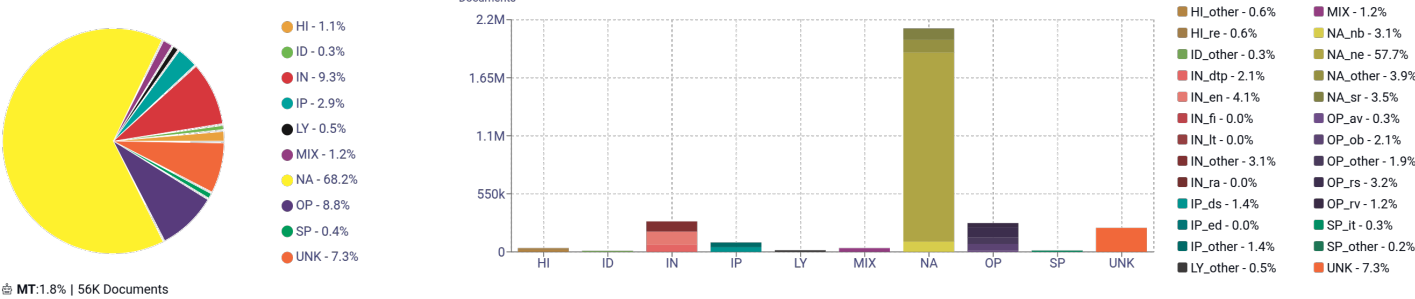
Top 10 domains

Domain	Docs	% of total
blogspot.com	145K	4.68%
wikipedia.org	130K	4.18%
thejasnews.com	118K	3.79%
mathrubhumi.com	88K	2.82%
sirajlive.com	69K	2.23%
blogspot.in	65K	2.08%
news18.com	55K	1.78%
boolokam.com	42K	1.34%
manoramaonline.com	41K	1.32%
indianexpress.com	40K	1.28%

Top 10 TLDs

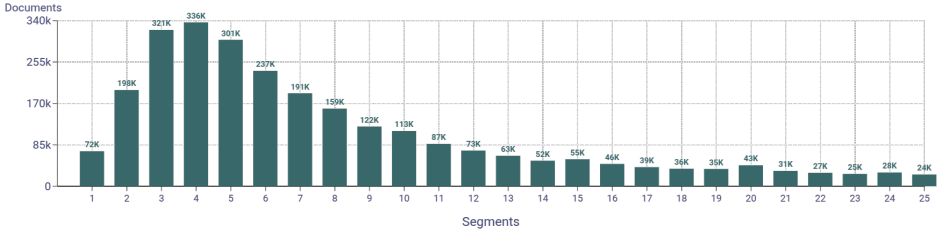
Domain	Docs	% of total
com	2.4M	77.30%
in	347K	11.18%
org	192K	6.17%
net	25K	0.79%
ae	15K	0.47%
tv	9.9K	0.32%
ie	9.5K	0.31%
news	9.5K	0.31%
gov.in	8.3K	0.27%
co.uk	7.2K	0.23%

Register labels

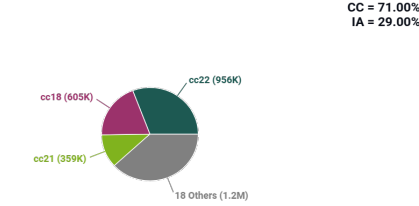


📄 MT:1.8% | 56K Documents

Documents size (in segments)

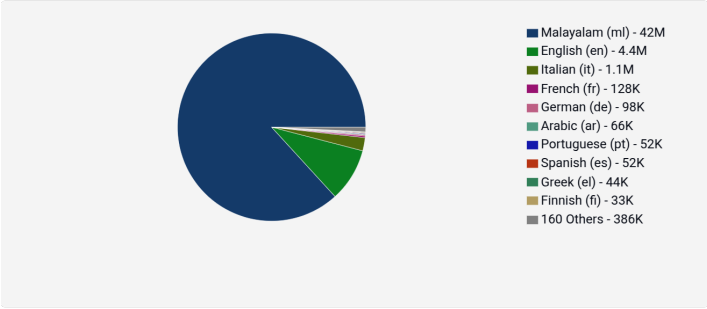


Documents by collection

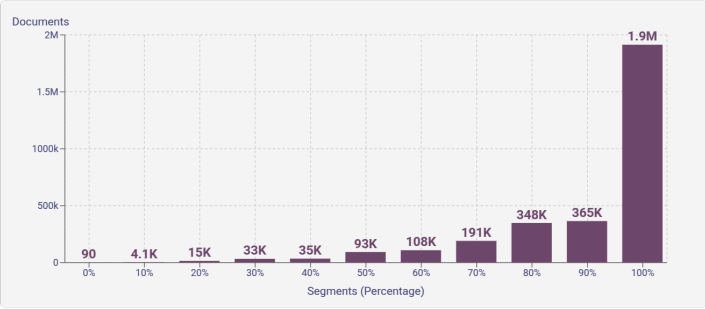


Language Distribution

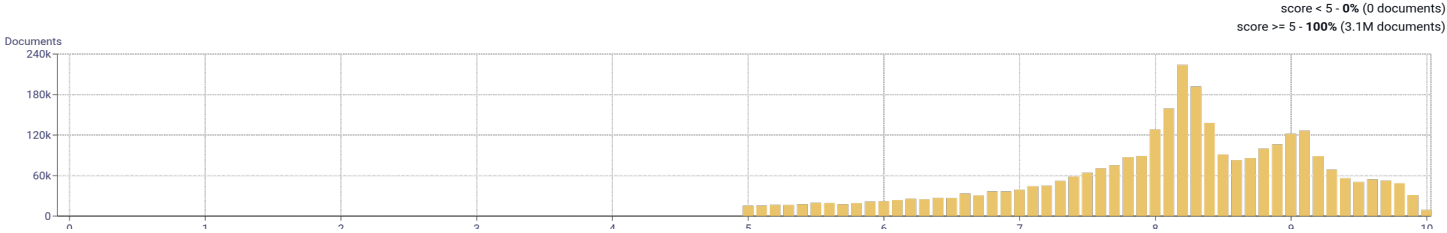
Number of segments in the Malayalam (ml) corpus



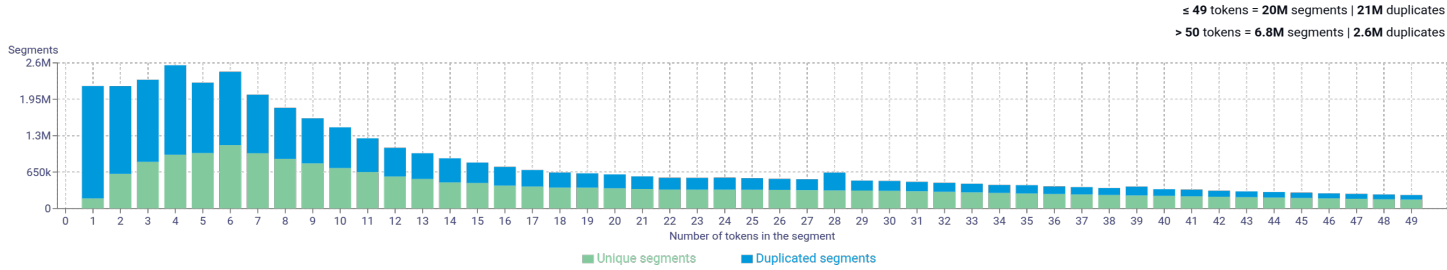
Percentage of segments in Malayalam (ml) inside documents



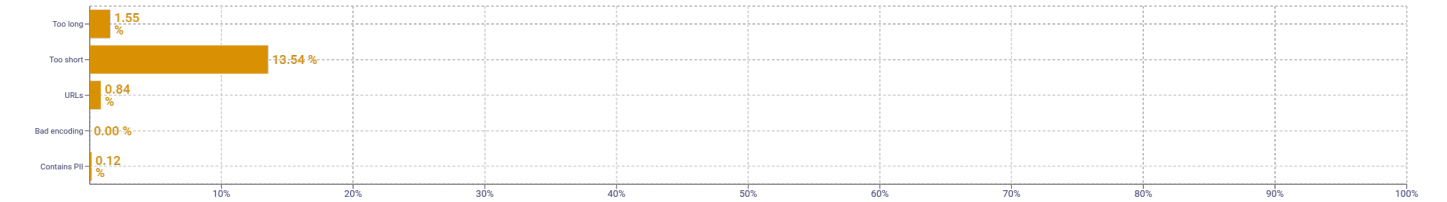
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	പറഞ്ഞു 2103251 പി 1874554 അന്ന് 1614324 കെ 1614032 സി 1539377
2	read more 244935 posted by 177507 കലയാളത്തിലോ ഇംഗ്ലീഷിലോ 167197 വായനക്കാരന്മാരെ അഭിപ്രായങ്ങളു് 141905 കഴിഞ്ഞ ദിവസം 138736
3	കലയാളത്തിലോ ഇംഗ്ലീഷിലോ എഴുതുക 128200 അവസംഗനം വൃക്കിപരമായ അധിക്ഷേപങ്ങളും 119801 അധിക്ഷേപങ്ങളും അശ്ലീല പാശ്ചാത്യങ്ങളും 111305 വൃക്കിപരമായ അധിക്ഷേപങ്ങളും അശ്ലീല 103018 വായനക്കാരന്മാരെ അഭിപ്രായങ്ങളു് കാഴ്ച 102810
4	വൃക്കിപരമായ അധിക്ഷേപങ്ങളും അശ്ലീല പാശ്ചാത്യങ്ങളും 103018 അവസംഗനം വൃക്കിപരമായ അധിക്ഷേപങ്ങളും അശ്ലീല 103018 വായനക്കാരന്മാരെ അഭിപ്രായങ്ങളു് കാഴ്ച എഴുതാനുവന്നതാണ് 102810 നവനവയി അവസംഗനം വൃക്കിപരമായ അധിക്ഷേപങ്ങളും 102604 വായനക്കാരന്മാരെ അഭിപ്രായ പ്രകാശങ്ങളിലോ അധിക്ഷേപങ്ങളിലോ 102560
5	അവസംഗനം വൃക്കിപരമായ അധിക്ഷേപങ്ങളും അശ്ലീല പാശ്ചാത്യങ്ങളും 103018 വൃക്കിപരമായ അധിക്ഷേപങ്ങളും അശ്ലീല പാശ്ചാത്യങ്ങളും ബിനാശക 102364 വായനക്കാരന്മാരെ അഭിപ്രായ പ്രകാശങ്ങളിലോ അധിക്ഷേപങ്ങളിലോ അശ്ലീല 102364 നവനവയി അവസംഗനം വൃക്കിപരമായ അധിക്ഷേപങ്ങളും അശ്ലീല 102364 അഭിപ്രായ പ്രകാശങ്ങളിലോ അധിക്ഷേപങ്ങളിലോ അശ്ലീല പാശ്ചാത്യങ്ങളിലോ 102364

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>

Register labels

Name	Abbr.	Name	Abbr.	Name	Abbr.
Machine-translated	MT	How-to or instructions	HI	Description of a thing or person	dtp
Lyrical	LY	Recipe	re	FAQ	fi
Spoken	SP	Informational persuasion	IP	Legal terms & conditions	lt
Interview	it	Description with intent to sell	ds	Opinion	OP
Interactive discussion	ID	News & opinion blog or editorial	ed	Review	rv
Narrative	NA	Informational description	IN	Opinion blog	ob
News report	ne	Encyclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				