

General overview

Corpus	Analytics date	Language
te_1.jsonl.tsv	3/21/2024	Telugu (te)

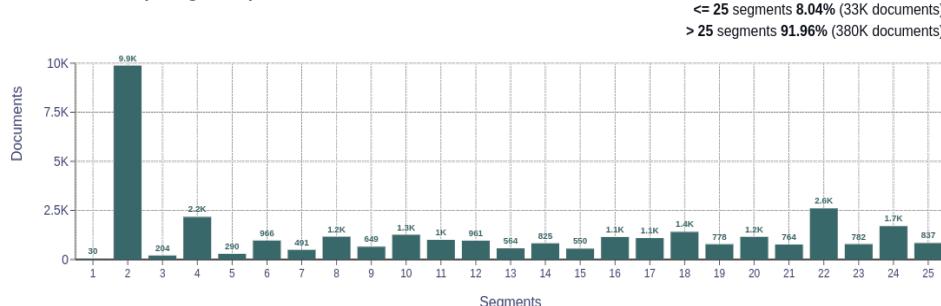
Volumes

Docs	Segments	Unique segments	Tokens	Size
415,598	51,141,409	51,242 (0.10 %)	573M	6.75 GB

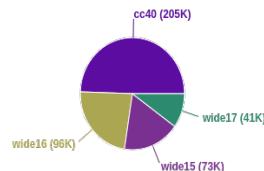
Type-Token Ratio

Telugu (te)
0.02

Documents size (in segments)

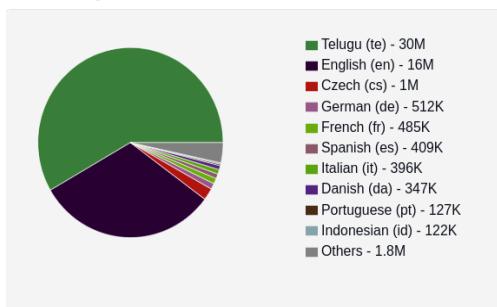


Documents by collection

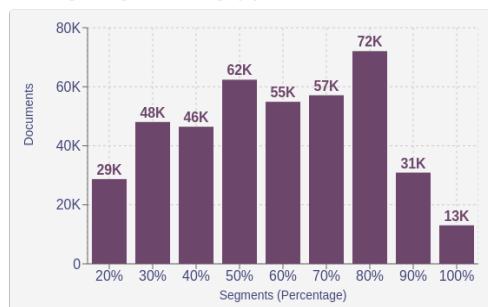


Language Distribution

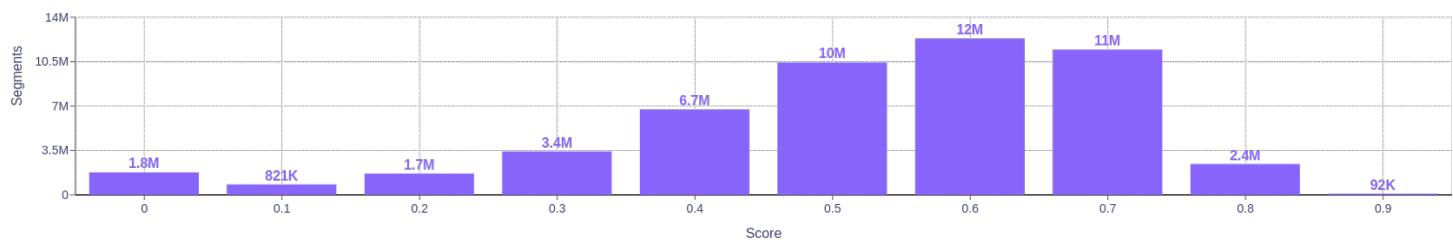
Number of segments



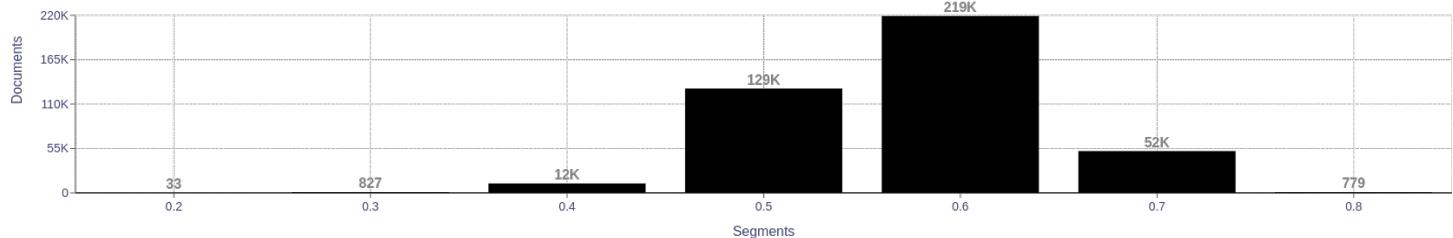
Percentage of segments in Telugu (te) inside documents



Distribution of segments by fluency score



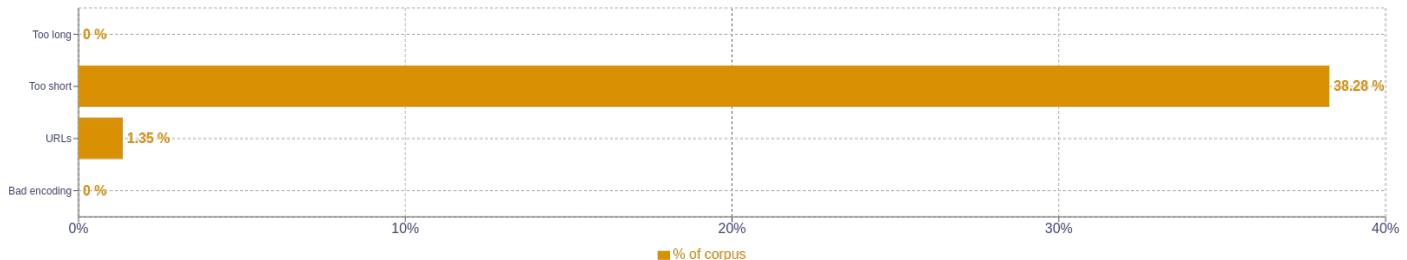
Distribution of documents by average fluency score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	ఏ 2660949 the 1952812 to 1857841 news 1504450 ఉ 1483233
2	span style 216760 telugu news 216156 of the 203200 posted by 198221 read more 198196
3	all rights reserved 147922 to twittershare to 121391 share to twittershare 121391 twittershare to facebookshare 118860 to facebookshare to 118860
4	share to twittershare to 121391 twittershare to facebookshare to 118860 to twittershare to facebookshare 118860 to facebookshare to pinterest 118860
5	తీరుమ తీరుమ తీరుమ తీరుమ 95956
	twittershare to facebookshare to pinterest 118860 to twittershare to facebookshare to 118860 share to twittershare to facebookshare 118860
	తీరుమ తీరుమ తీరుమ తీరుమ 82331 భాగస్వాముల నెయ్యందిfacebookక భాగస్వాముల నెయ్యంpinterestక భాగస్వాము 31802

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 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>, , , 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>).

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.

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>