

## General overview

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

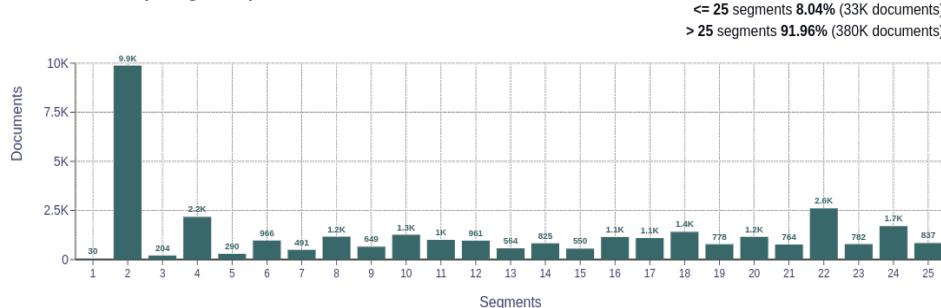
## Volumes

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

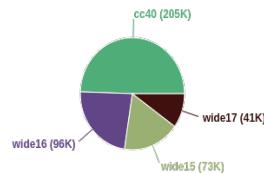
## Type-Token Ratio

Telugu (te)
0.02

## Documents size (in segments)

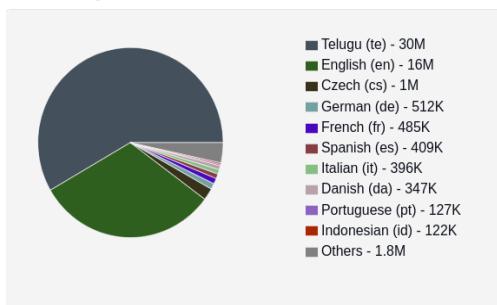


## Documents by collection

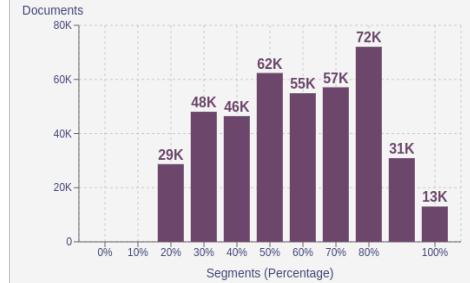


## Language Distribution

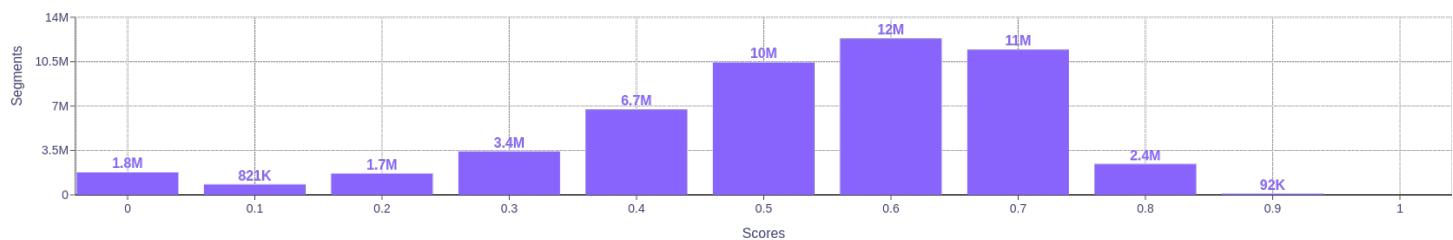
## Number of segments



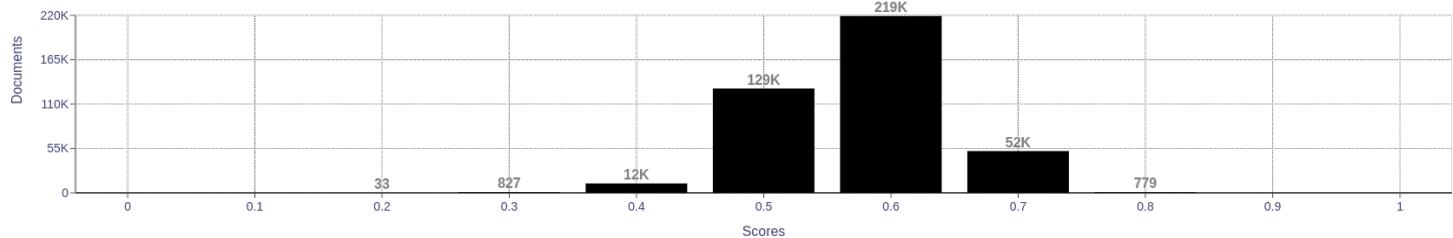
## Percentage of segments in Telugu (te) inside documents



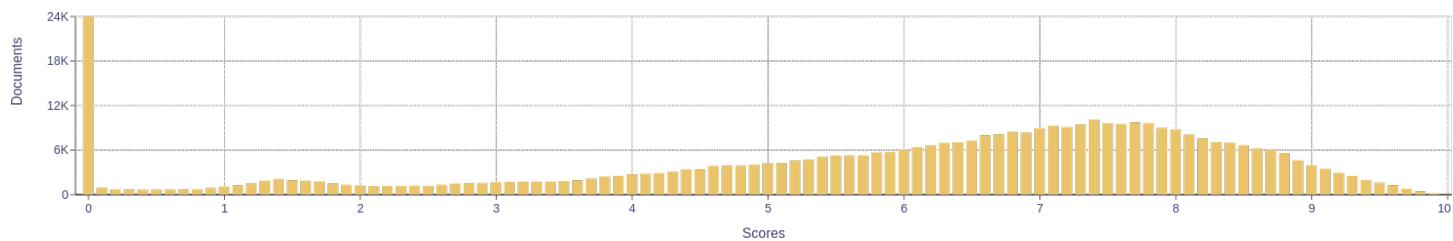
## Distribution of segments by fluency score



## Distribution of documents by average fluency score

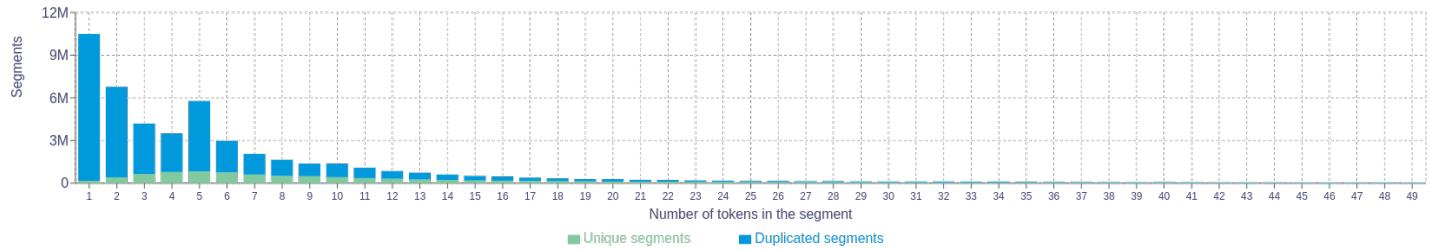


## Distribution of documents by document score



## Segment length distribution by token

<= 49 tokens = 9.5M segments | 39M duplicates  
> 50 tokens = 2.1M segments | 621K duplicates



## Segment noise distribution



## Frequent n-grams

Size	n-grams
1	ఆస   2660949, తిరాము   1952812, తిరాము   1857841, తిరాము   1504450, తిరాము   1483233
2	తిరాము   216760, తిరాము   216156, తిరాము   203200, తిరాము   198221, తిరాము   198196
3	all rights reserved   147922, తిరాము   121391, తిరాము   121391, తిరాము   121391, తిరాము   118860, తిరాము   118860
4	తిరాము   95956
5	తిరాము   82331, తిరాము   82331, తిరాము   31802

## 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.slinfo.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>