

General overview

Corpus	Analytics date	Language
kn_1.jsonl.tsv	3/19/2024	Kannada (kn)

Volumes

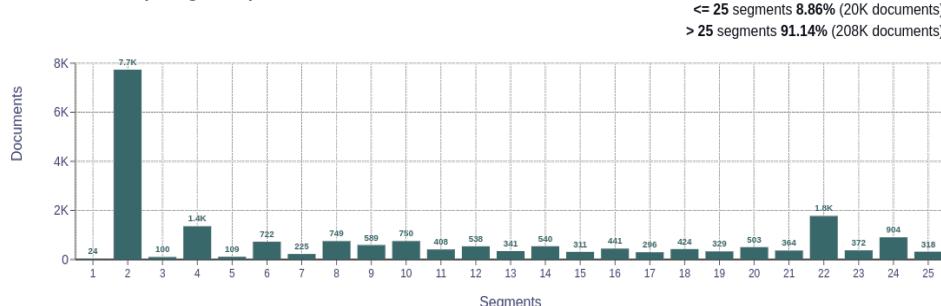
Docs	Segments	Unique segments	Tokens	Size
228,215	29,241,332	29,293 (0.10 %)	301M	3.8 GB

Type-Token Ratio

Kannada (kn)

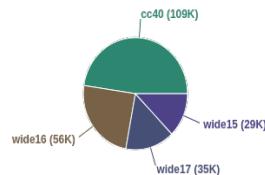
0.03

Documents size (in segments)



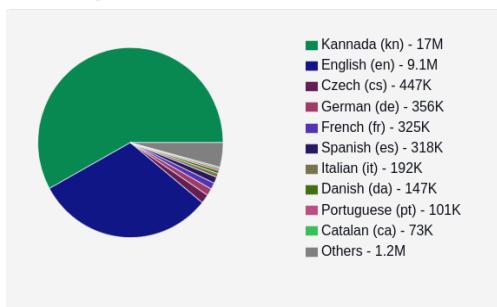
<= 25 segments 8.86% (20K documents)
> 25 segments 91.14% (208K documents)

Documents by collection

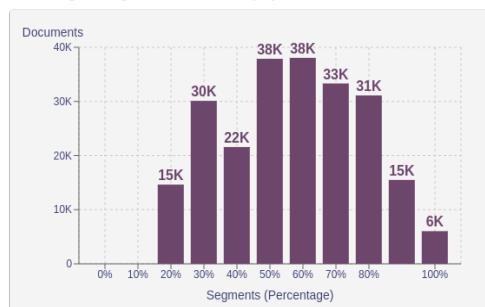


Language Distribution

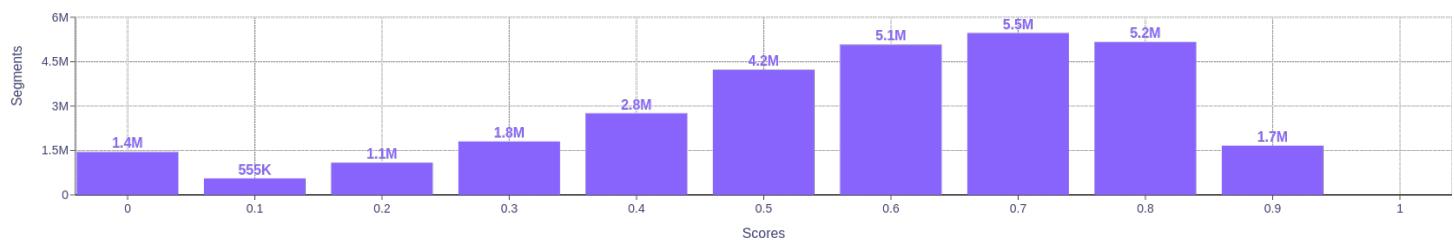
Number of segments



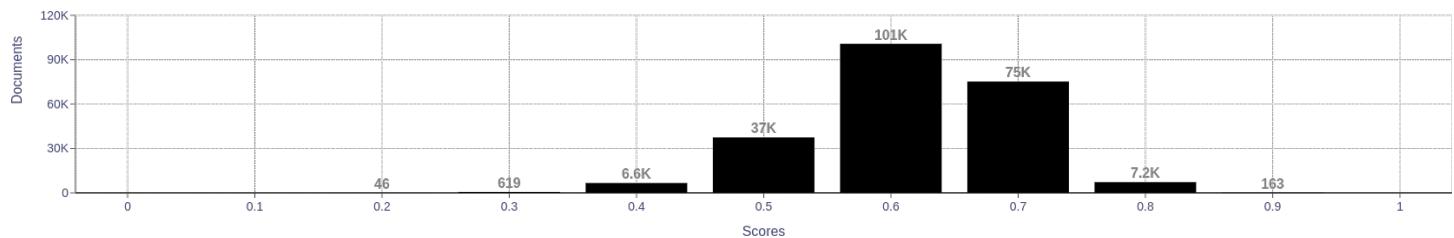
Percentage of segments in Kannada (kn) inside documents



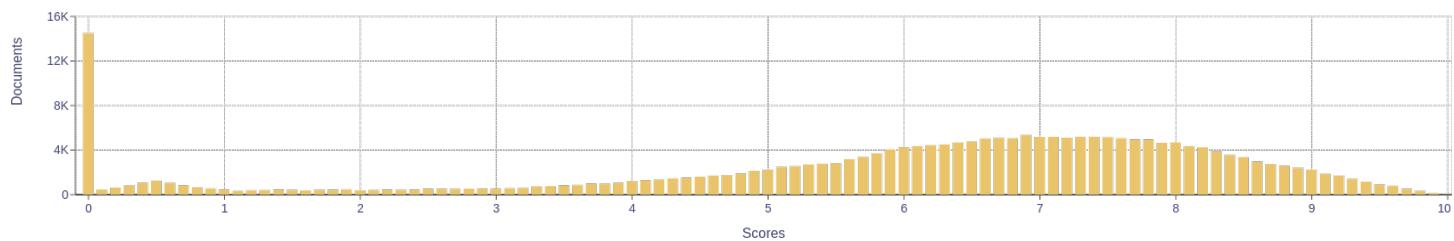
Distribution of segments by fluency score



Distribution of documents by average fluency score

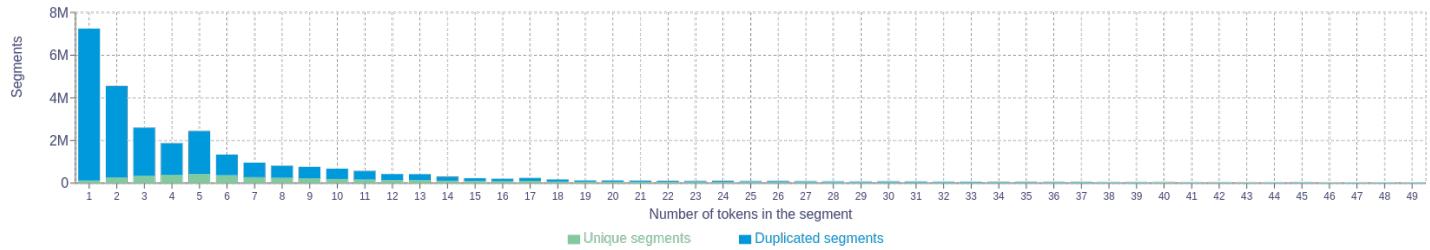


Distribution of documents by document score

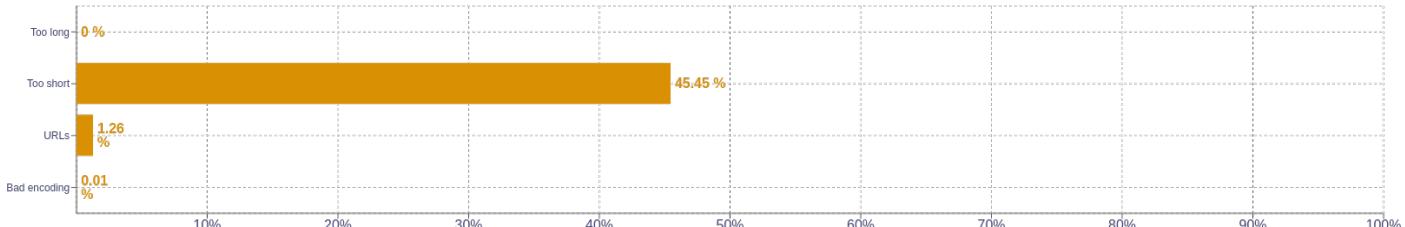


Segment length distribution by token

<= 49 tokens = 5.2M segments | 23M duplicates
 > 50 tokens = 1.1M segments | 315K duplicates



Segment noise distribution



Frequent n-grams

Size	n-grams
1	the 1080755 to 1014534 in 823065 news 801921 of 674903
2	of the 120948 span style 115762 rights reserved 112002 all rights 104013 no comments 96691
3	all rights reserved 103766 opens in new 53064 in new window 53044 to twittershare to 50021 share to twittershare 50021
4	opens in new window 53041 share to twittershare to 50021 twittershare to facebookshare to 45919 to twittershare to facebookshare 45919 to facebookshare to pinterest 45919
5	twittershare to facebookshare to pinterest 45919 to twittershare to facebookshare to 45919 share to twittershare to facebookshare 45919 leave a reply cancel reply 31648 address will not be published 27715

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>