

General overview

Corpus	Date	Language
kan_Knda.jsonl.tsv	9/18/2024	Kannada (kn)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
1,335,847	24,929,282	12,772,248 (51.23 %)	653M	4,274,156,104	10.46 GB

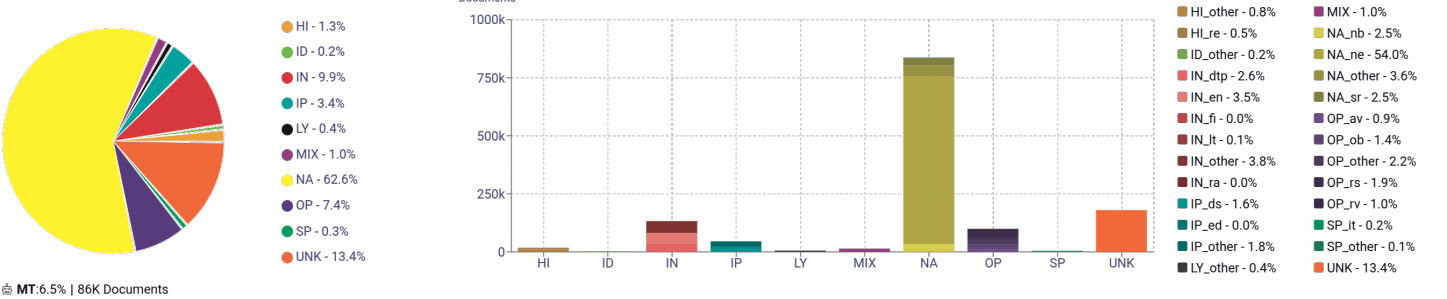
Top 10 domains

Domain	Docs	% of total
<a href="#">prajavani.net</a>	79K	5.92%
<a href="#">udayavani.com</a>	52K	3.90%
<a href="#">wikipedia.org</a>	52K	3.90%
<a href="#">news18.com</a>	47K	3.52%
<a href="#">blogspot.com</a>	45K	3.37%
<a href="#">filmibeat.com</a>	42K	3.13%
<a href="#">oneindia.com</a>	39K	2.92%
<a href="#">asianetnews.com</a>	33K	2.45%
<a href="#">indiatimes.com</a>	31K	2.29%
<a href="#">gulfkannadiga.com</a>	24K	1.78%

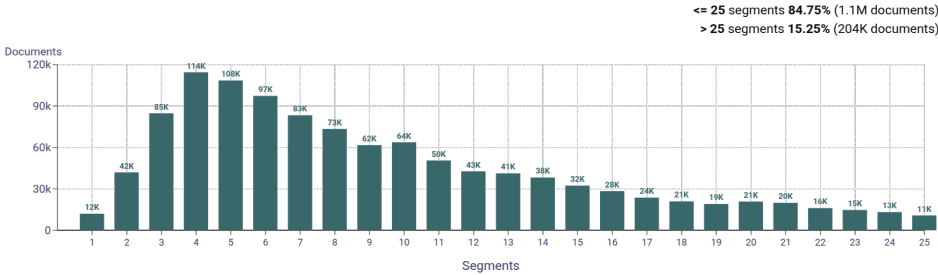
Top 10 TLDs

Domain	Docs	% of total
com	905K	67.75%
in	165K	12.33%
net	116K	8.70%
org	90K	6.75%
news	24K	1.78%
co.in	3.4K	0.25%
live	3.4K	0.25%
gov.in	3.1K	0.23%
today	2.2K	0.17%
online	2K	0.15%

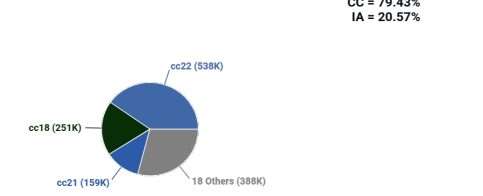
Register labels



Documents size (in segments)

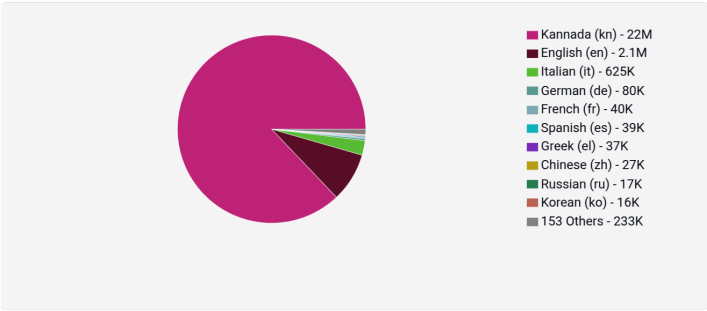


Documents by collection

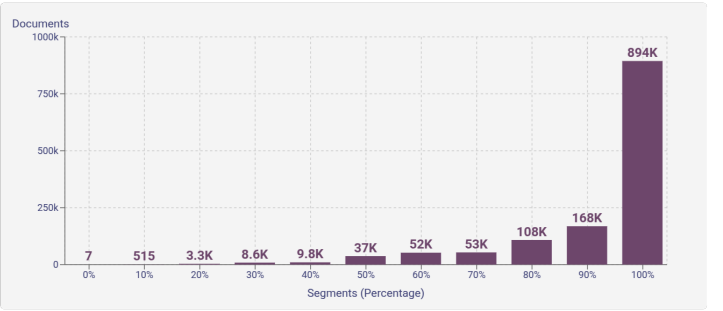


Language Distribution

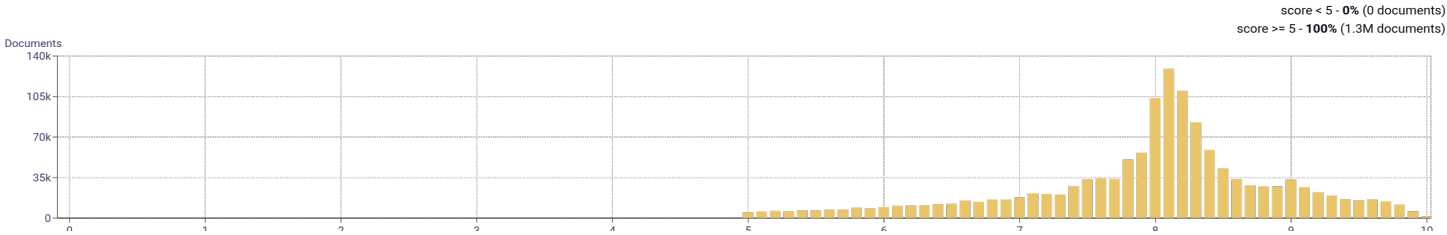
Number of segments in the Kannada (kn) corpus



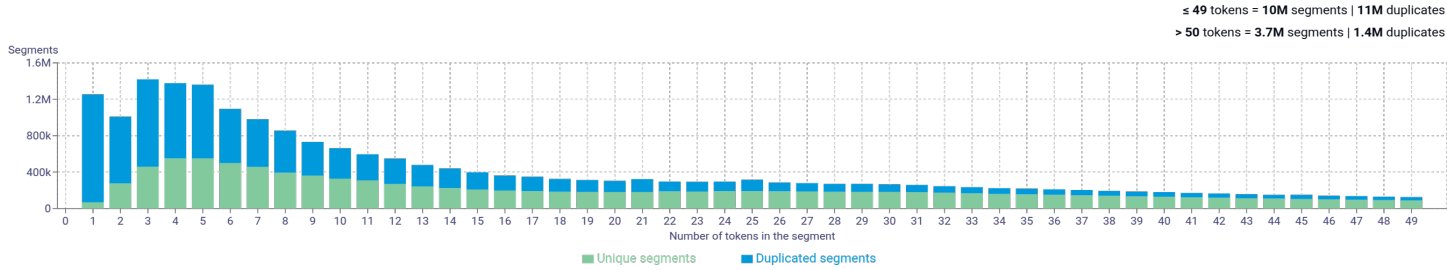
Percentage of segments in Kannada (kn) inside documents



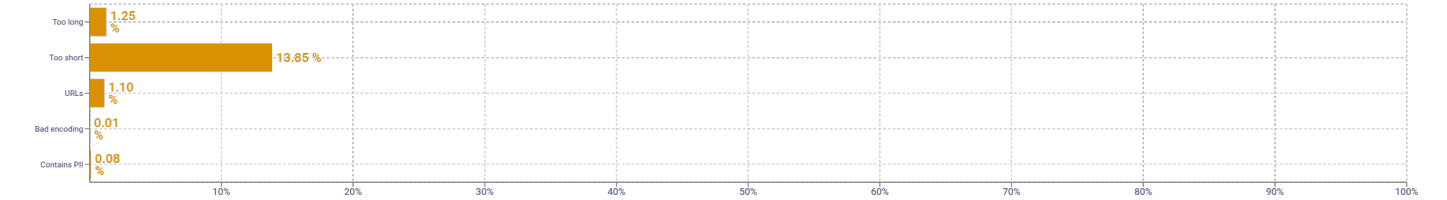
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	ಸವು   1068630, ನಿವು   1018629, ನೀವು   922077, ನಾವು   912901, ನನ್ನ   878025
2	ಇದನ್ನೂ ಓದಿ   89081, of the   78259, ಕೋಟಿ ರೂ   64493, ಸರೇಯ ಮೋದಿ   59053, pm ist   52952
3	lua error in   33784, ಪ್ರಧಾನಿ ಸರೇಯ ಮೋದಿ   33331, from the original   32458, archived from the   32371, error in ಮ್ಯಾಕೋ   31277
4	archived from the original   32364, lua error in ಮ್ಯಾಕೋ   31250, from the original on   30807, to compare number with   29362, compare number with nil   29362
5	archived from the original on   29966, to compare number with nil   29362, attempt to compare number with   29362, ಸುದ್ದಿಗಾಗಿ ಪ್ರಜಾವಾಣಿ ಆನ್ ಲೈನ್   24587

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 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>, <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>

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				