

General overview

Corpus	Date	SL	TL
hplt-v2-en-kn.tsv	1/22/2025	English (en)	Kannada (kn)

Volumes

Segments	SL tokens	SL characters	SL size
720,157	19M	100,021,435	95.78 MB

TL tokens	TL characters	TL size
17M	112,959,684	284.92 MB

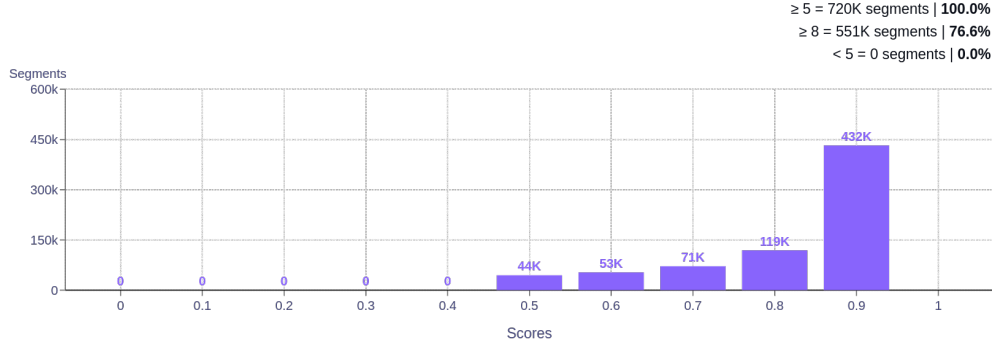
Dataset top 10 domains

SL domain	Segments	TL domain	Segments
wikipedia.org	31.3%	wikipedia.org	22.3%
hebbarskitchen.com	4.9%	hebbarskitchen.com	4.5%
bajajfinserv.in	4.2%	bajajfinserv.in	3.7%
educationbro.com	3.2%	vsaduidoma.com	2.5%
itsmygame.org	2.9%	itsmygame.org	2.4%
vsaduidoma.com	2.6%	educationbro.com	1.6%
vessoft.com	1.5%	wordproject.org	1.6%
schools-wikipedia.org	1.4%	oneindia.com	1.4%
boldsky.com	1.3%	news18.com	1.4%
uber.com	1.2%	boldsky.com	1.3%

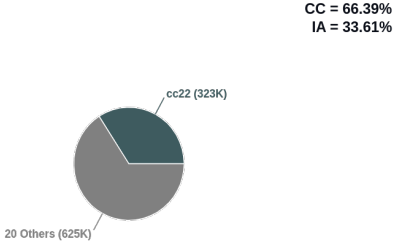
Dataset top 10 TLDs

SL domain	Segments	TL domain	Segments
com	81.8%	com	64.4%
org	44.9%	org	33.6%
in	10.1%	in	10.1%
net	4.6%	net	4.8%
gov.in	0.8%	gov.in	0.8%
cc	0.6%	nic.in	0.5%
nic.in	0.5%	de	0.4%
co.in	0.5%	co.in	0.4%
plus	0.5%	zone	0.4%
zone	0.4%	news	0.4%

Translation likelihood

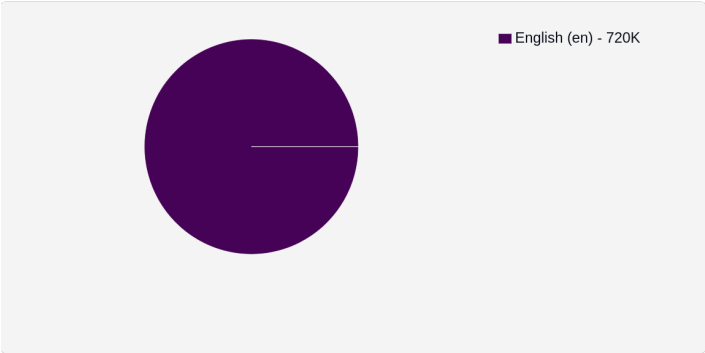


Collections

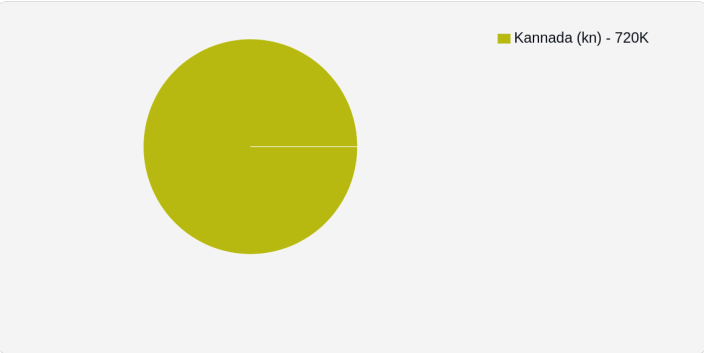


Language Distribution

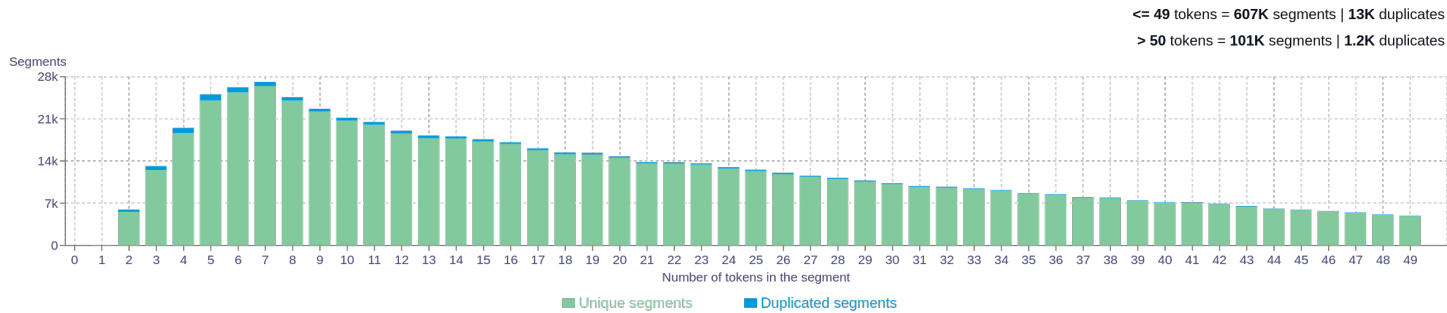
Source



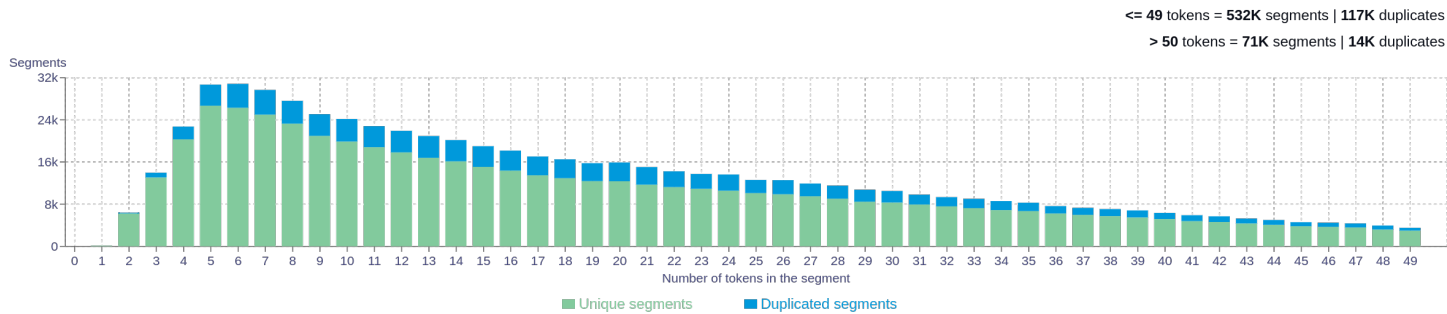
Target



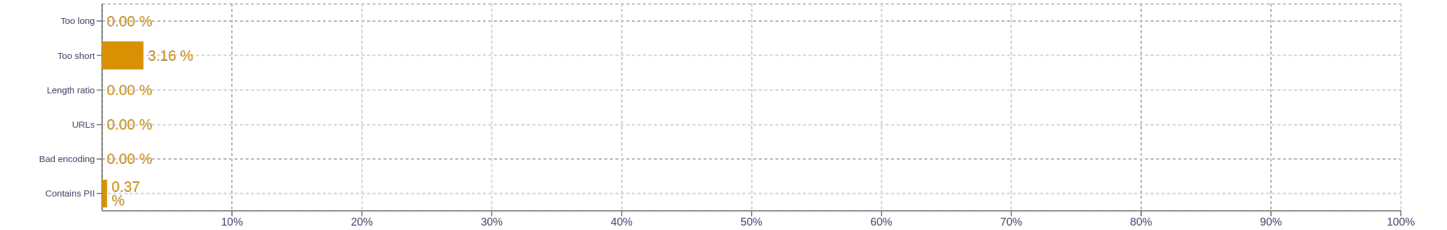
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	also   37380new   31072india   29940one   29857use   23932
2	new delhi   5704prime minister   5495personal loan   5004bajaj finserv   4340united states   4113
3	jammu and kashmir   2270minister narendra modi   2068prime minister narendra   2018play the game   1859embed the game   1621
4	characteristics of the game   2205prime minister narendra modi   2005games like the game   1026step by step photo   959reserve bank of india   823
5	technical characteristics of the game   2205coupon code to your clipboard   680button to save the coupon   680board of control for cricket   652control for cricket in india   646

Target n-grams

Size	n-grams
1	ನಿವೃ   96932ನಿವೃ   87887ನಾವು   44718ನವು   36413ನಾನು   28695
2	ಸಹಾಯ ಮಾಡುತ್ತದೆ   4365ಬಜಾಜ್ ಫಿನ್ಸರ್ವ್   3091ವಿರಾಟ್ ಕೊಹ್ಲಿ   3071ನರೇಂದ್ರ ಮೋದಿ   3062ನಿವೃ ವೈಯಕ್ತಿಕ   2964
3	ಪ್ರಧಾನಿ ನರೇಂದ್ರ ಮೋದಿ   1866ನಿವೃ ವೆಬ್ಸೈಟ್ ಅಟದ   1630ನಿವೃ ವೈಯಕ್ತಿಕ ಮಾಹಿತಿಯನ್ನು   1152ಪೆಟ್ರೋಲ್ ಮತ್ತು ಡೀಸೆಲ್   1127ಚಿನ್ನ ಸೂಪರ್ ಕಿಂಗ್ಸ್   1118
4	ಕ್ರಿವ್ಲೋರ್ಡ್ ಕೂಪನ್ ಕೋಡ್ ಉಳಿಸಲು   751ಕೋಡ್ ಉಳಿಸಲು ಗುಂಡಿಯನ್ನು ಕ್ಲಿಕ್   751ಕೂಪನ್ ಕೋಡ್ ಉಳಿಸಲು ಗುಂಡಿಯನ್ನು   751ಪಾಕವಿಧಾನದ ಈ ಪೋಸ್ಟ್‌ನಿಂದಿಗೆ ನನ್ನ   371ಬಜಾಜ್ ಫಿನ್ಸರ್ವ್ ಪರ್ಫನಲ್ ಲೋನ್   370
5	ಕ್ರಿವ್ಲೋರ್ಡ್ ಕೂಪನ್ ಕೋಡ್ ಉಳಿಸಲು ಗುಂಡಿಯನ್ನು   751ಕೂಪನ್ ಕೋಡ್ ಉಳಿಸಲು ಗುಂಡಿಯನ್ನು ಕ್ಲಿಕ್   751ವೈಶಿಷ್ಟ್ಯಗಳನ್ನು ನೋಡಲು ವಿವರಗಳು ಪುಟಕ್ಕೆ ಹೋಗಿ   246ಉತ್ಕೃಷ್ಟ ಕನ್ನಡಿಪಟ ಅಯ್ಯಗಳಲ್ಲಿ ಇವುಗಳು ಸೇರಿವೆ   211comನಿವೃ ಕ್ರಿವ್ಲೋರ್ಡ್ ಕೂಪನ್ ಕೋಡ್ ಉಳಿಸಲು   199

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>