

General overview

Corpus	Analytics date	Language
npi_Deva.jsonl.tsv	9/25/2024	Nepali (npi)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
2,777,574	37,138,196	22,365,347 (60.22 %)	1.2B	17.88 GB	7,221,503,597

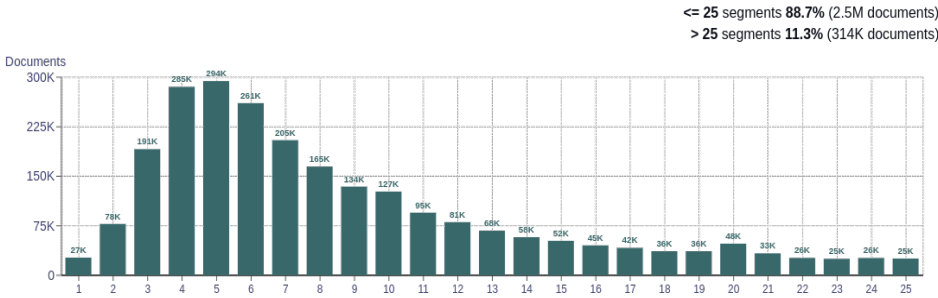
Top 10 domains

Domain	Docs	% of total
onlinekhabar.com	36K	1.31
abhiyan.com.np	29K	1.05
ekantipur.com	28K	1.00
ujyaaloonline.com	25K	0.91
ktmkhabar.com	25K	0.89
wikipedia.org	22K	0.81
setopati.com	22K	0.80
blogspot.com	21K	0.76
ratopati.com	19K	0.67
sajhasabal.com	18K	0.64

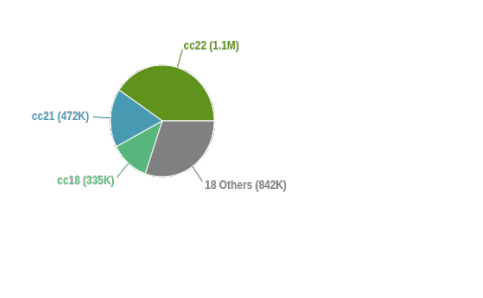
Top 10 TLDs

Domain	Docs	% of total
com	2.4M	87.77
com.np	107K	3.83
org	70K	2.52
net	31K	1.12
gov.np	31K	1.10
tv	26K	0.93
org.np	14K	0.49
com.au	6.7K	0.24
co.il	6K	0.22
news	5K	0.18

Documents size (in segments)

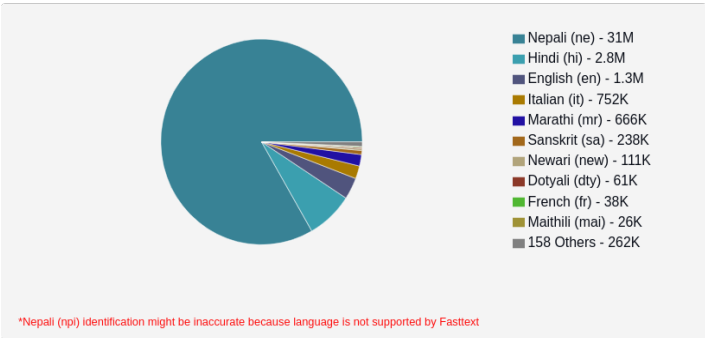


Documents by collection

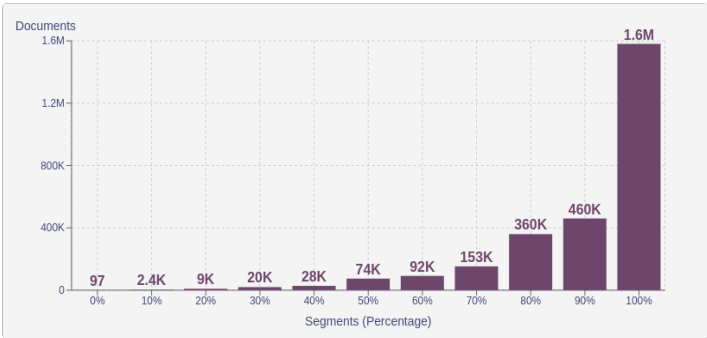


Language Distribution

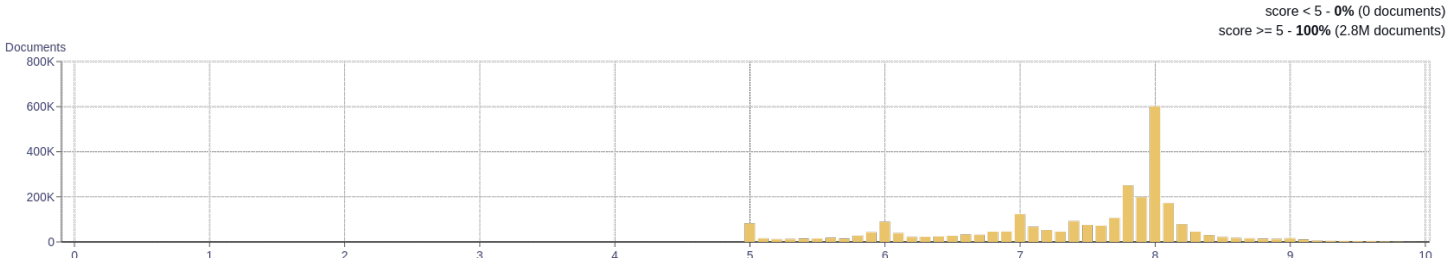
Number of segments



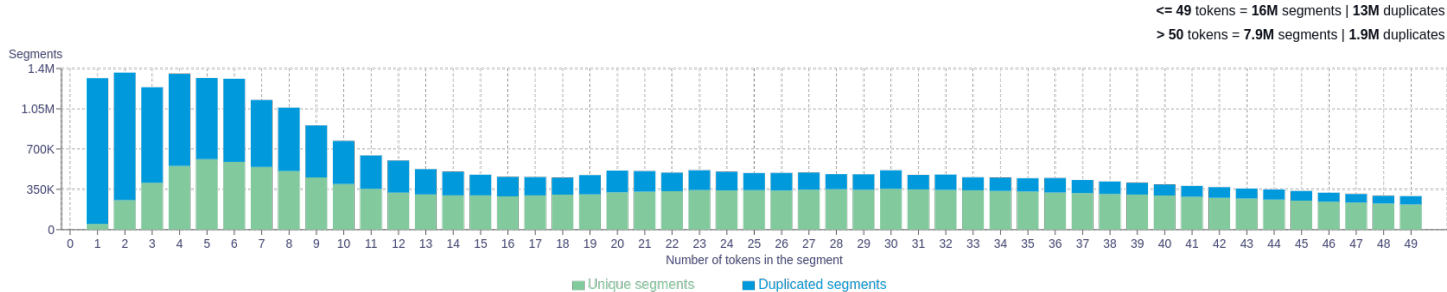
Percentage of segments in Nepali (npi) inside documents



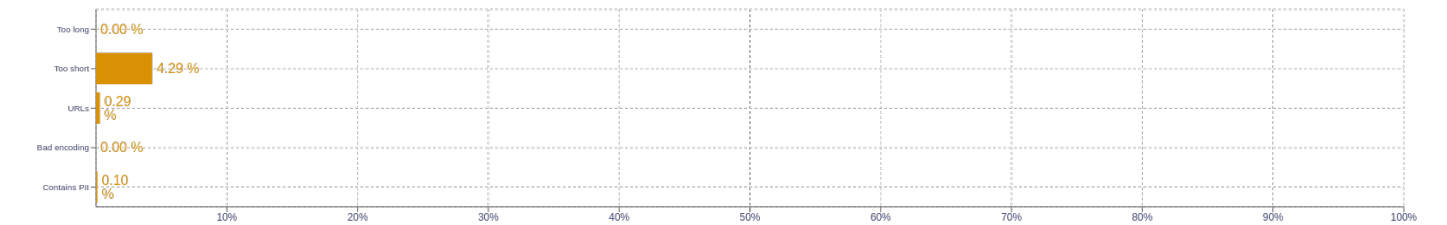
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>नेपाल 2495926</div> <div>नेपाली 2231383</div> <div>गरिएको 2186314</div> <div>काम 2120619</div> <div>भएका 1710736</div>
2	<div>जानकारी दिनुभयो 188639</div> <div>प्रहरी कार्यालय 181360</div> <div>यहा नं 176503</div> <div>read more 176500</div> <div>केपी शर्मा 155024</div>
3	<div>जिल्ला प्रहरी कार्यालय 100852</div> <div>प्रधानमन्त्री केपी शर्मा 85087</div> <div>प्रमुख जिल्ला अधिकारी 76609</div> <div>केपी शर्मा ओलीले 68217</div> <div>अध्यक्ष पुष्पकमल दाहाल 51563</div>
4	<div>प्रधानमन्त्री केपी शर्मा ओलीले 41964</div> <div>आएको सोही अवधिको तुलनामा 25032</div> <div>गत आएको सोही अवधिको 24899</div> <div>घालू आएको तेस्रो वैशाखमा 24806</div> <div>लाख खुद मुनाफा आर्जन 24724</div>
5	<div>गत आएको सोही अवधिको तुलनामा 24853</div> <div>घालू आएको तेस्रो वैशाखमा बैङ्कको 24579</div> <div>बैङ्कको खराब कर्जा शून्य दशमलव 24499</div> <div>वैशाखमा बैङ्कको खराब कर्जा शून्य 24481</div> <div>आएको तेस्रो वैशाखसम्ममा बैङ्कले रु 24478</div>

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>).

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>