General overview

Corpus	us Analytics date	
sat Olck.jsonl.tsv	11/28/2024	Santali (sat)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
2,566	45,801	26,205 (57.21 %)	1.3M	14.52 MB	6,222,595

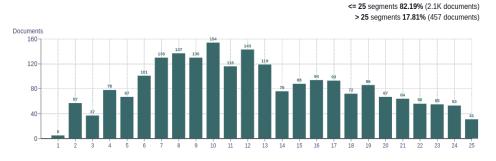
Top 10 domains

Domain	Docs	% of total
wikipedia.org	2.3K	89.44
vikaspedia.in	178	6.94
globalvoices.org	47	1.83
raharahla.com	15	0.58
wikimedia.org	8	0.31
wikiplanet.click	5	0.19
wikiversity.org	4	0.16
know.cf	3	0.12
santalinews.com	3	0.12
mediawiki.org	2	0.08

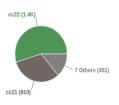
Top 10 TLDs

Domain	Docs	% of total
org	2.4K	91.89
in	178	6.94
com	21	0.82
click	5	0.19
cf	3	0.12
cn	1	0.04

Documents size (in segments)

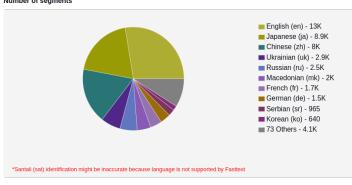


Documents by collection

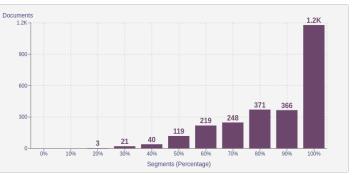


Language Distribution

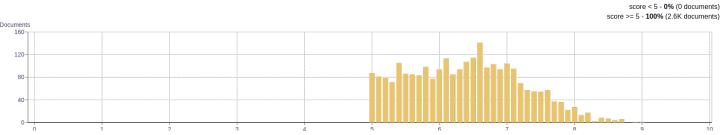
Number of segments



Percentage of segments in Santali (sat) inside documents



Distribution of documents by document score



Segment length distribution by token

<= 49 tokens = 20K segments | 17K duplicates > 50 tokens = 8.3K segments | 2.4K duplicates



Segment noise distribution



Frequent n-grams

Size	n-grams
1	(200200 1799) (302000 8934) (40200 7799) (62 6556)
2	①最近の間の と思り記載で 8849) (いたな P25会所 1134) (U22所属 b5心地 922) (P2® O別いと"® 705) (P2を O別いと"® 665)
3	(wikipedia articles with 394) (from the original 382) (archived from the 378) (שאַס אַרְצָטְאָשׁ אָרָאָשׁ אָרָאָשׁ אָרָאָשׁׁ
4	archived from the original 378) (১৯১৪৯ ১৫৪৯৫ এটর P25৪৯ 322) (from the original on 315) (৪৪৪৯৫ ১৯৫৪০ ১৯ ৫৮৪৭ 218) (১৯৫৪৯ ০৫ ৪৯৪৯৯ ১৮৫৭ 181)
5	(archived from the original on 312) (K2340 32606 W93 P2600 P260 277) (D38600 OZ P9800 OZ PP8800 OZ PP880

About HPLT Analytics

Volumes - Segments

Segments correspond to paragraph and list boundaries as defined by HTML elements (, , , 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 (, , , 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