

General overview

Corpus	Analytics date	Language
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

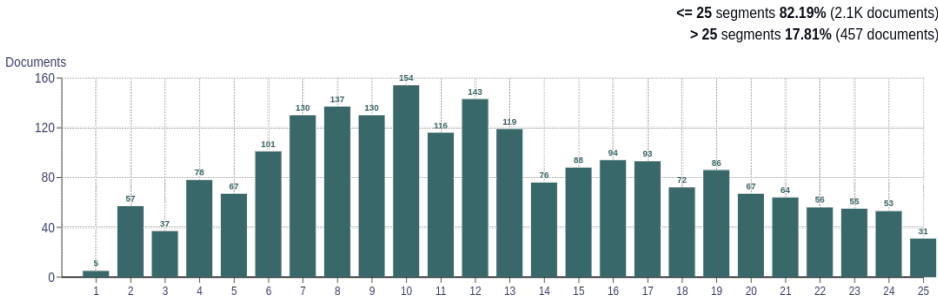
Top 10 domains

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

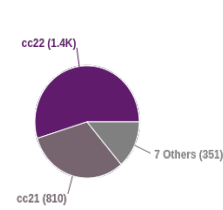
Top 10 TLDs

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)

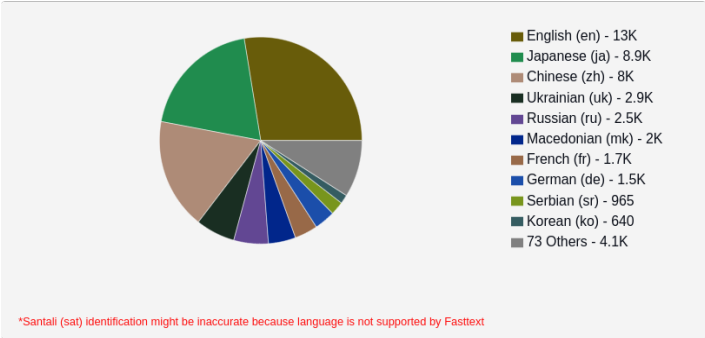


Documents by collection

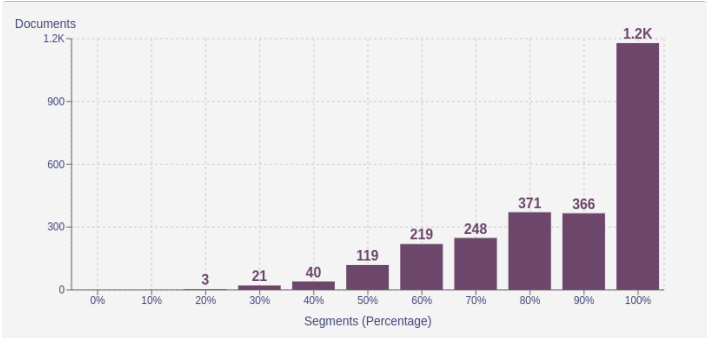


Language Distribution

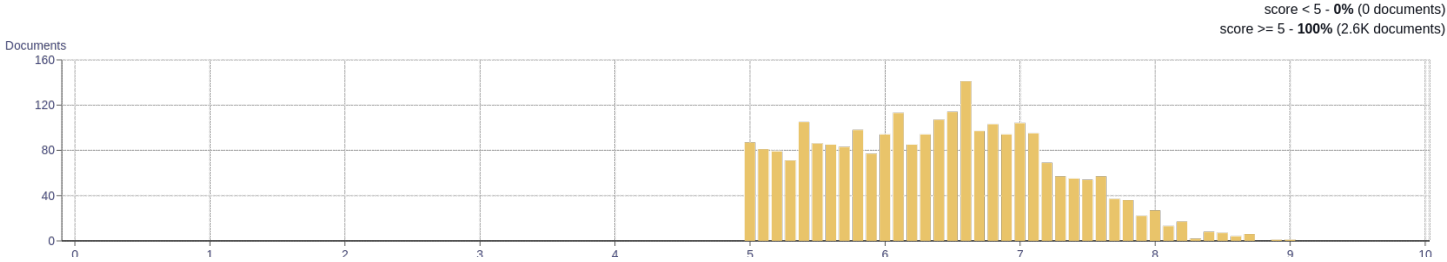
Number of segments



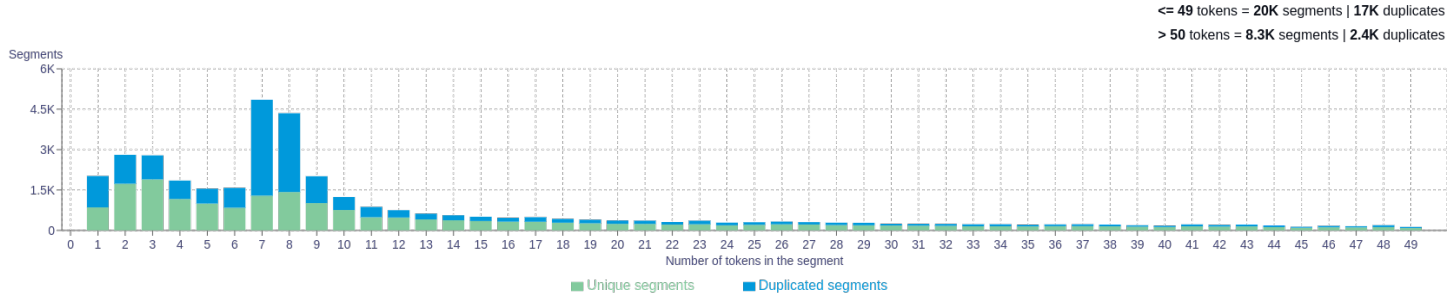
Percentage of segments in Santali (sat) inside documents



Distribution of documents by document score



Segment length distribution by token



Segment noise distribution

