

General overview

Corpus	Date	Language
ory_Orya.jsonl.tsv	9/22/2024	Odia (ory)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
412,895	3,595,619	2,426,565 (67.49 %)	121M	778,355,779	1.91 GB

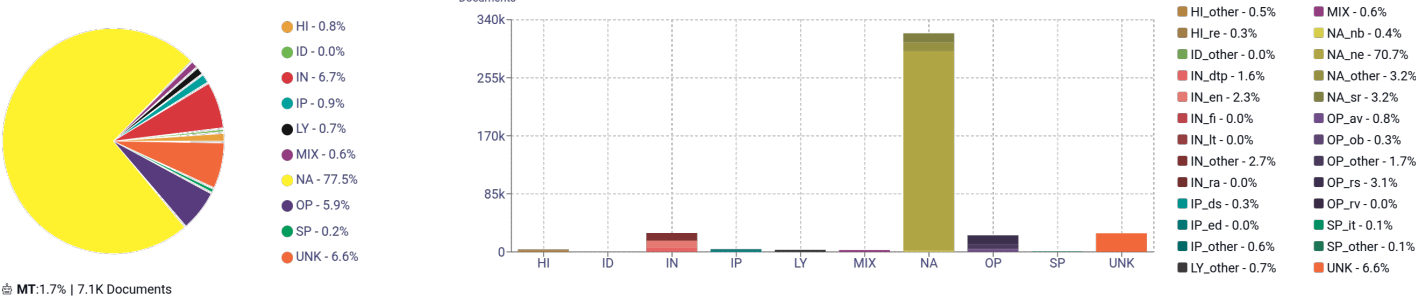
Top 10 domains

Domain	Docs	% of total
prameyanews7.com	39K	9.48%
sambad.in	39K	9.39%
kanaknews.com	30K	7.18%
dharitri.com	20K	4.95%
samajajive.in	19K	4.55%
odishareporter.in	17K	4.02%
news18.com	13K	3.14%
eodishasamecher...	10K	2.54%
wikipedia.org	10K	2.49%
prameya.com	5.7K	1.38%

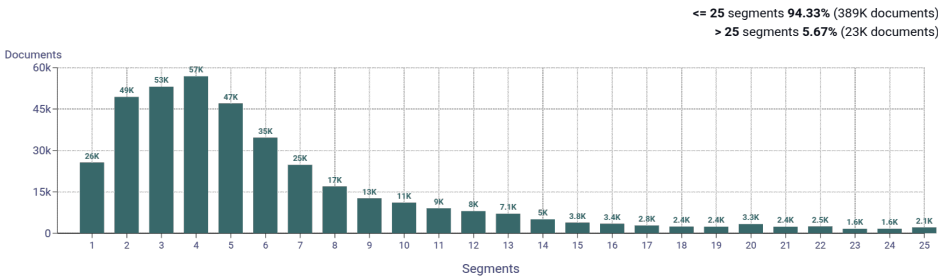
Top 10 TLDs

Domain	Docs	% of total
com	254K	61.61%
in	123K	29.77%
org	24K	5.74%
me	4.3K	1.04%
live	1.8K	0.44%
net	1.2K	0.28%
nic.in	948	0.23%
gov.in	910	0.22%
is	562	0.14%
co.in	479	0.12%

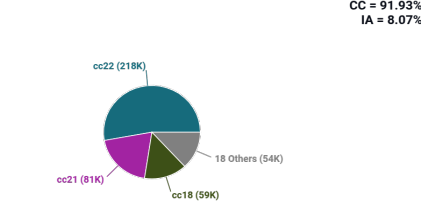
Register labels



Documents size (in segments)

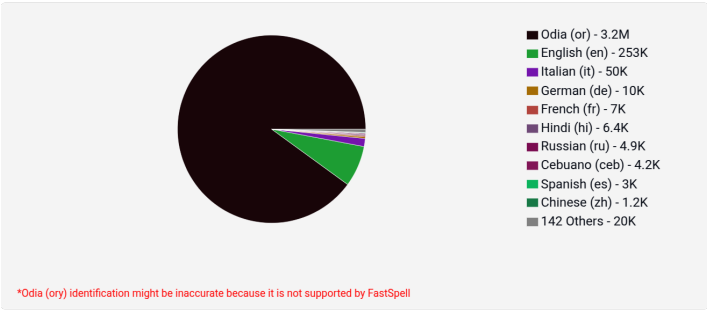


Documents by collection

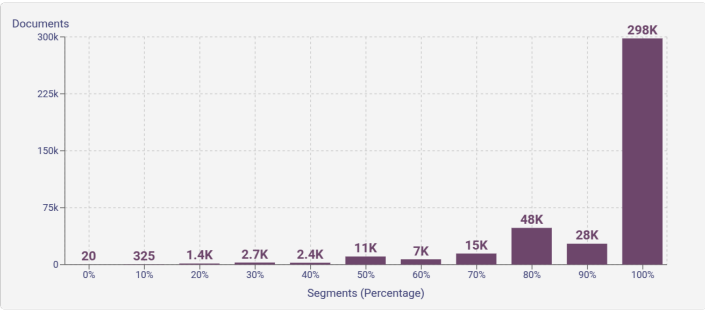


Language Distribution

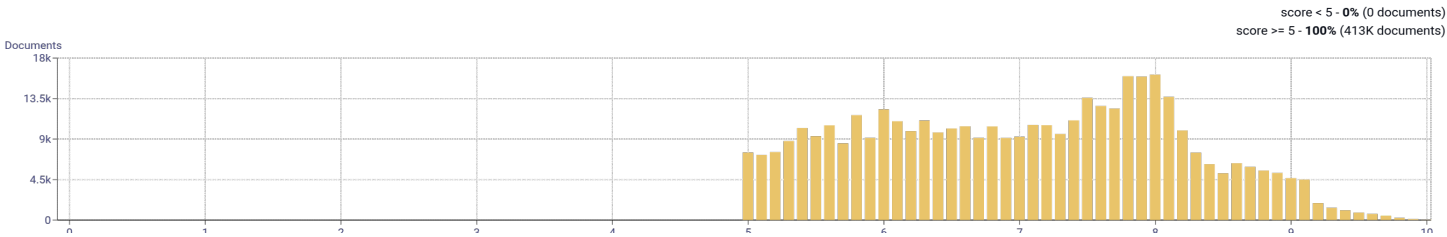
Number of segments in the Odia (ory) corpus



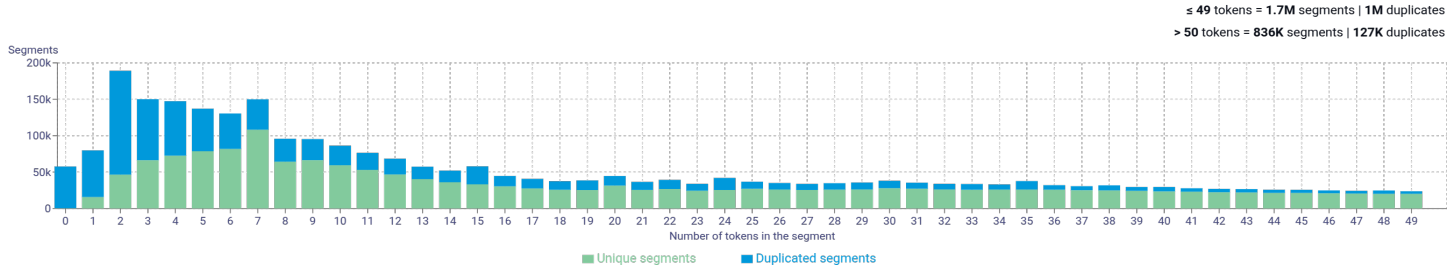
Percentage of segments in Odia (ory) inside documents



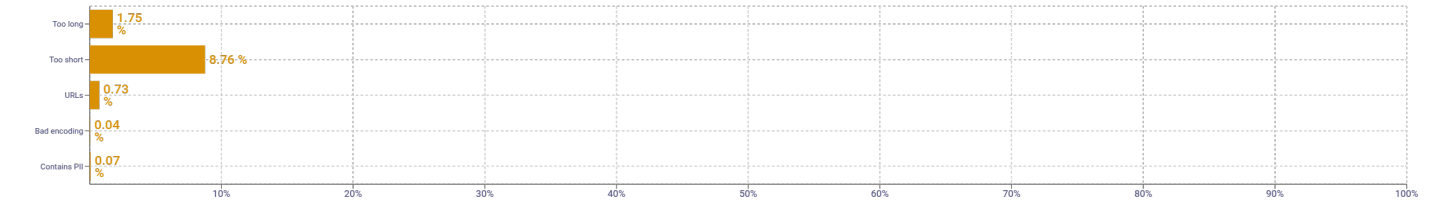
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	ସେ 659801 ଏକ 600284 କରବା 421415 ବୋଲି 371775 ଶେ 312114
2	କହିଛନ୍ତି ଶେ 53427 ଗଭିରା ଗଭିରା 31309 ଶେଉଁଠି ଶେଉଁଠି 26936 ସେ କହିଛନ୍ତି 26410 ଶେଉଁ ଦୃଷ୍ଟି 25708
3	ସୁସ୍ଥମାନବୀ ନବୀନ ପଦ୍ମନାଭ 12409 ଶେଉଁଠି କହିଲୁ କହିଲୁ 10140 ପାଉବା ପାଇଁ ଶେଉଁଠି 10137 ଗଭିରା ବଡ଼ ଶେଉଁଠି 10125 ଶେ ଗଭିରା ବଡ଼ ଶେଉଁଠି 10125
4	ଗଭିରା ବଡ଼ ଶେଉଁଠି ପାଉବା 10125 ଶେ ଗଭିରା ବଡ଼ ଶେଉଁଠି 10125 ଶେଉଁଠି ଶେଉଁଠି ଶେଉଁଠି କହିଲୁ 10125 ପାଉବା ପାଇଁ ଶେଉଁଠି କହିଲୁ 10125 ପଦ୍ମନାଭ ପଦ୍ମନାଭ ଶେଉଁଠି ଶେଉଁଠି 10125
5	ଶେ ଗଭିରା ବଡ଼ ଶେଉଁଠି ପାଉବା 10125 ଶେଉଁଠି ଶେଉଁଠି ଶେଉଁଠି କହିଲୁ ଶେ 10125 ବଡ଼ ଶେଉଁଠି ପାଉବା ପାଇଁ ଶେଉଁଠି 10125 ପାଉବା ପାଇଁ ଶେଉଁଠି କହିଲୁ କହିଲୁ 10125 କହିଲୁ ଶେ ଗଭିରା ବଡ଼ ଶେଉଁଠି 10125

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>

Register labels

Name	Abbr.	Name	Abbr.	Name	Abbr.
Machine-translated	MT	How-to or instructions	HI	Description of a thing or person	dtip
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				