

General overview

Corpus	Date	Language
pbt_Arab.jsonl.tsv	9/20/2024	Southern Pashto (pbt)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
466,472	8,454,662	5,506,825 (65.13 %)	306M	1,295,601,474	2.13 GB

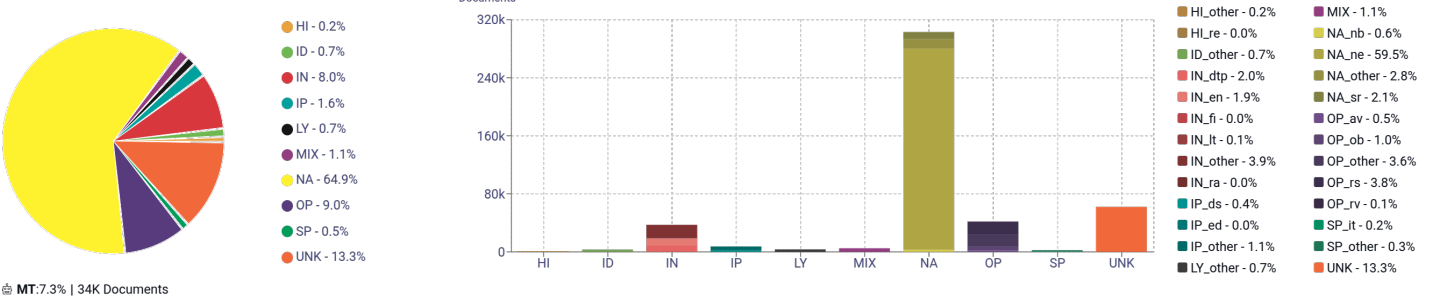
Top 10 domains

Domain	Docs	% of total
pashtovoa.com	38K	8.10%
mashaairadio.com	32K	6.81%
bakhtarnews.af	28K	6.01%
tolonews.com	24K	5.23%
nunn.asia	15K	3.21%
tolafghan.com	13K	2.89%
wikipedia.org	12K	2.63%
larawbar.net	12K	2.48%
taand.com	11K	2.44%
dw.com	10K	2.25%

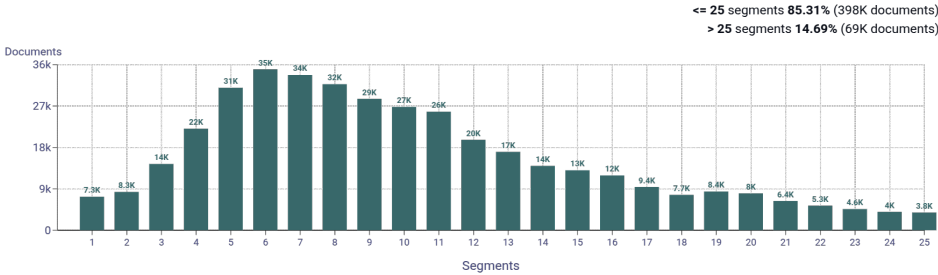
Top 10 TLDs

Domain	Docs	% of total
com	288K	61.68%
af	49K	10.45%
gov.af	25K	5.33%
net	24K	5.11%
org	23K	4.93%
asia	15K	3.23%
cn	6.8K	1.46%
website	4.6K	0.98%
info	3.5K	0.74%
ir	2.7K	0.57%

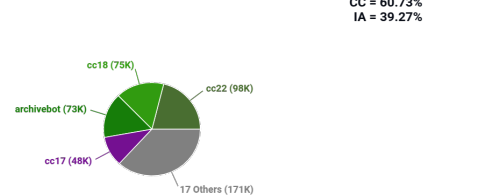
Register labels



Documents size (in segments)

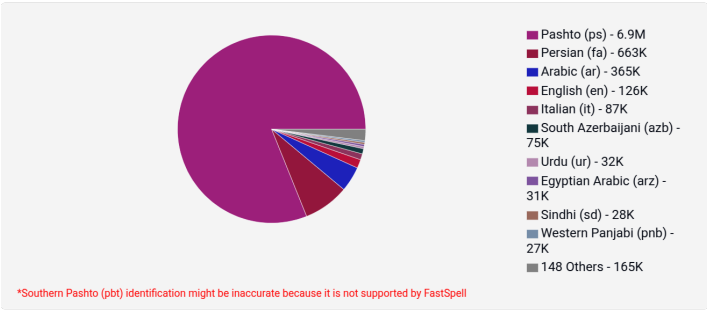


Documents by collection

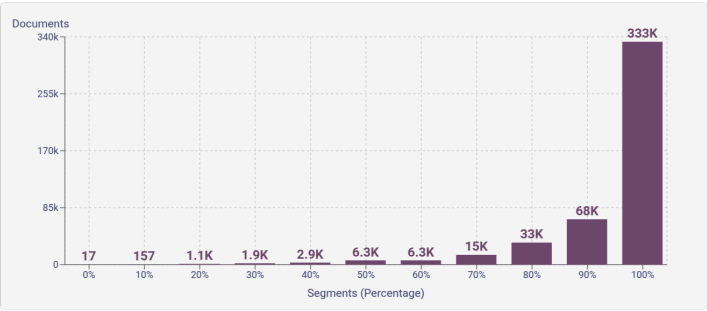


Language Distribution

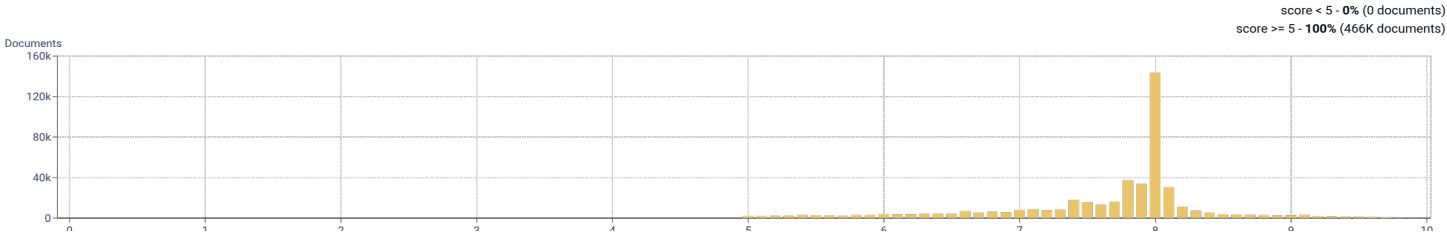
Number of segments in the Southern Pashto (pbt) corpus



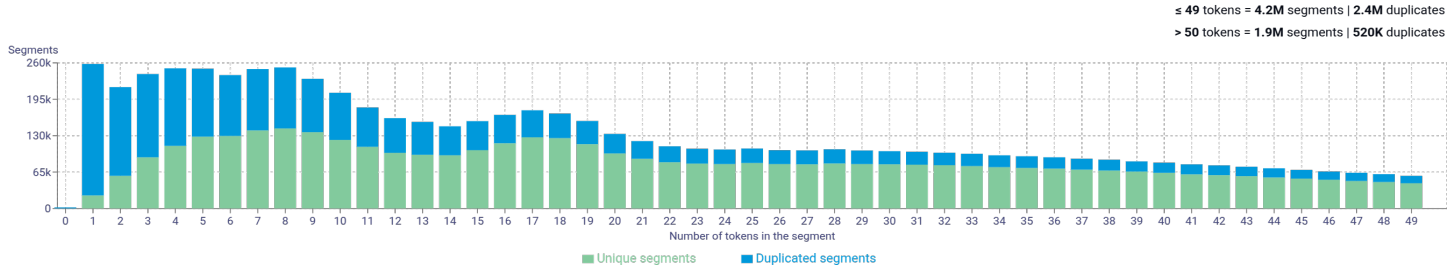
Percentage of segments in Southern Pashto (pbt) inside documents



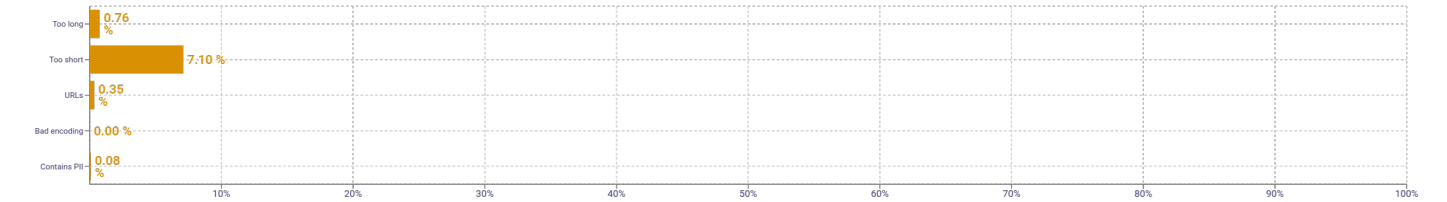
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	چپ 5754718 5332453 کپ 1706720 دی 1624373 1341005 بی
2	حال کی 110752 کی چپ 116629 برخہ کی 130557 دی چپ 188791 افغانستان کی 193226
3	چپ د دی 31550 کی د افغانستان 36106 حال کی چپ 39227 چپ د افغانستان 51729 چپ بہ دی 56303
4	صلی اللہ علیہ وسلم 13877 افغانستان د اسلامی جمہوریہ 14219 صلی اللہ علیہ وسلم 16442 صلی اللہ علیہ وسلم 25413 چپ بہ افغانستان کی 28146
5	کی د پستو او بلوڅو 4178 جمهور رئیس محمد اشرف غنی 4799 رسول اللہ صلی اللہ علیہ 7995 رسول اللہ صلی اللہ علیہ 8720 رسول اللہ صلی اللہ علیہ 12440

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/pablopt16n/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	dtp
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	Enciclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				