

General overview

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

Volumes

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

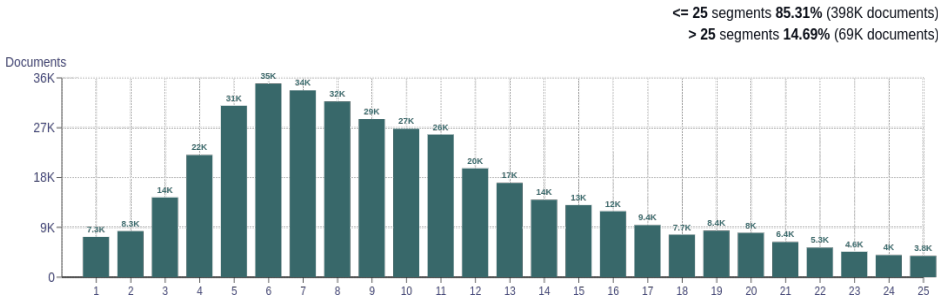
Top 10 domains

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

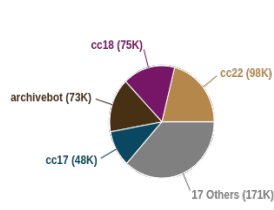
Top 10 TLDs

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

Documents size (in segments)

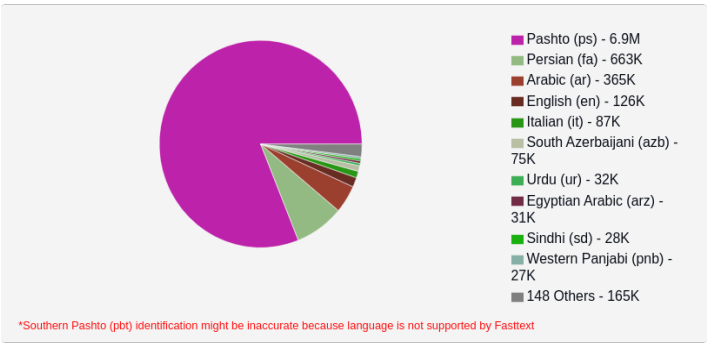


Documents by collection

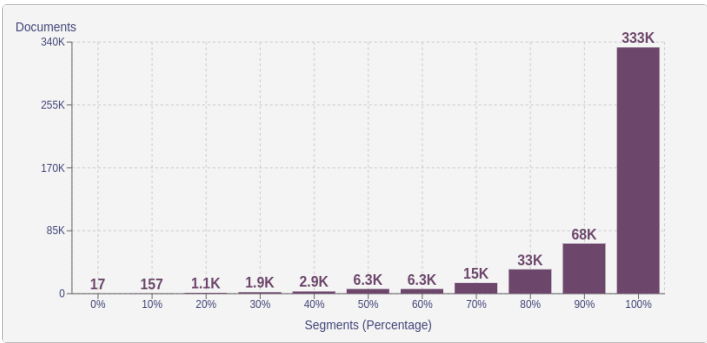


Language Distribution

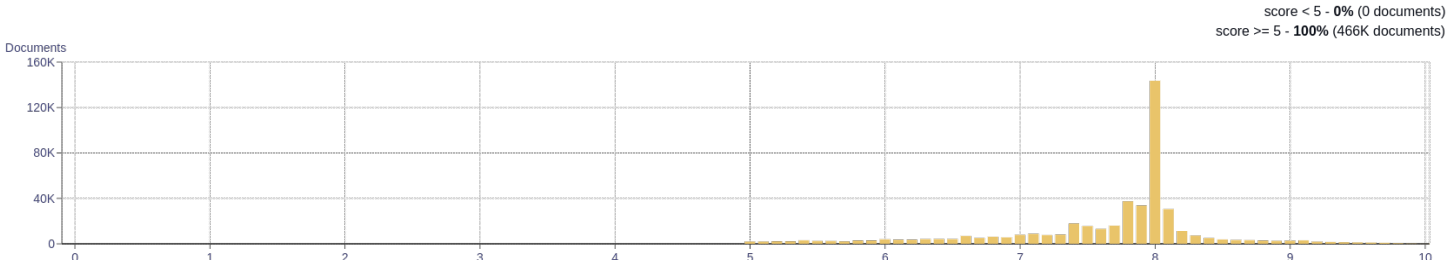
Number of segments



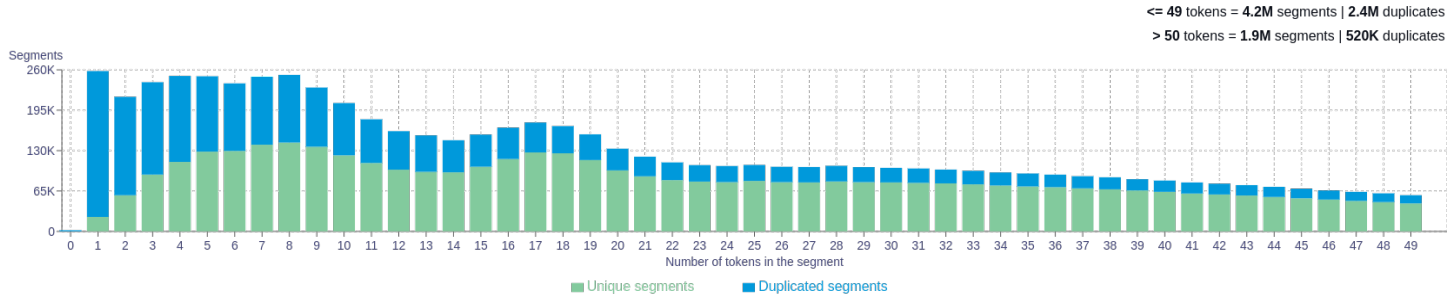
Percentage of segments in Southern Pashto (pbt) inside documents



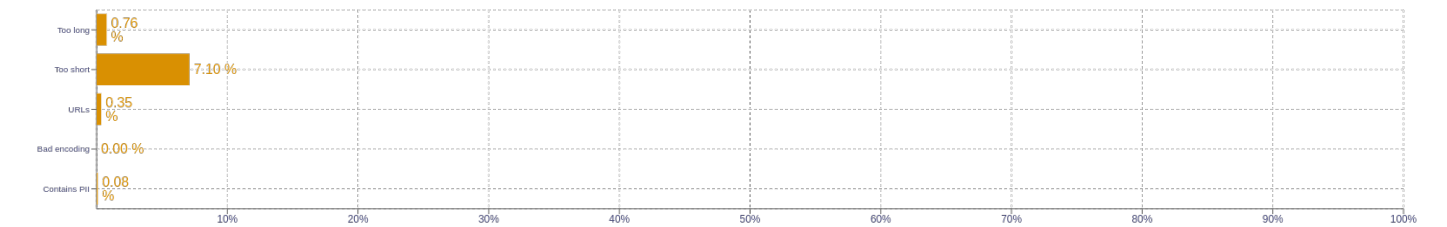
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	5754718   چي   5332453   کي   1706720   دي   1624373   دي   1341005   يي
2	193226   افغا نستان کي   188791   دي چي   130557   برخه کي   116629   کي چي   110752   حال کي
3	56303   چي په دي   51729   چي د افغا نستان   39227   حال کي چي   36106   کي د افغا نستان   31550   چي د دي
4	28146   چي په افغا نستان کي   25413   صلی الله عليه وسلم   16442   صلی الله عليه وسلم   14219   افغا نستان د اسلامي جمهوريت   13877   صلی الله عليه وسلم
5	12440   رسول الله صلی الله عليه   8720   رسول الله صلی الله عليه   7995   رسول الله صلی الله عليه   4799   جمهور رئيس محمد اشرف غني   4178   کي د پښتو او بلوڅو

About HPLT Analytics

Volumes - Segments

Segments correspond to paragraph and list boundaries as defined by HTML elements (<p>, <ul>, <ol>, 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 of 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>, <ul>, <ol>, 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>