

General overview

Corpus	Analytics date	Language	
urd Arab.isonl.tsv	9/22/2024	Urdu (ur)	

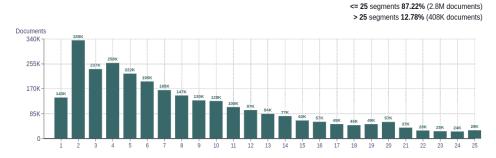
Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
2 102 002	E0 620 040	29,400,119	2.20	16 4 CP	0.050.062.100

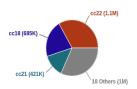
Top 10 domains

op 10 domains				Top 10 TLDs			
Domain	Docs	% of total		Domain	Docs	% of total	
urduvoa.com	141K	4.40		com	1.8M	55.04	
dailypakistan.com.pk	110K	3.44		com.pk	348K	10.89	
urdupoint.com	105K	3.28		org	225K	7.04	
wikipedia.org	93K	2.91		tv	221K	6.92	
arynews.tv	85K	2.65		pk	217K	6.78	
siasat.com	72K	2.24		net	127K	3.96	
nawaiwaqt.com.pk	58K	1.80		info	31K	0.97	
geourdu.com	49K	1.54		in	25K	0.78	
express.pk	39K	1.21		xyz	23K	0.73	
news18.com	34K	1.08		ir	22K	0.70	

Documents size (in segments)

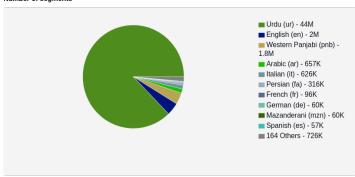


Documents by collection

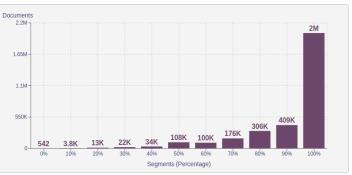


Language Distribution

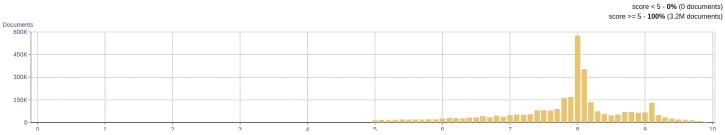
Number of seaments



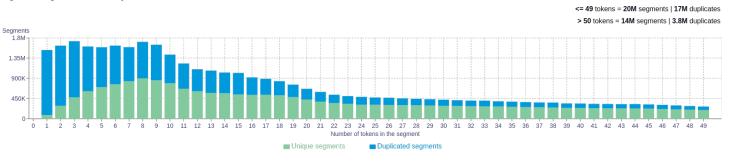
Percentage of segments in Urdu (ur) inside documents



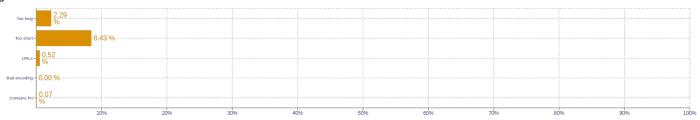
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	(ت 31036762 کی 31036762 (کی 30463430 (کی 31036762 (کی ا
2	(اس کا 1445823 (سب سب ا 1599731) (آپ کو 1599731) (انہوں نہ 2442696 (نہ کہا)
3	(كرت كا ليا 479172 (كا كينا تقا 491060 (الله عليه وسلم 547055 (انبون تاكيا 729302 (صلن الله عليه عالم
4	(الله عليه وآل وسلم 157488 (الله عليه الله عليه الله عليه الله عليه وسلم الله عليه وسلم الله عليه وسلم الله عليه وسلم ا (157532 (الله عليه وسلم الله عليه وسلم ا (15752)
5	آب صلى الله عليه وسلم 84318 (الله صلى الله عليه وسلم 116307 (سول الله صلى الله عليه واله وسلم الله عليه وسلم الله وسلم الل

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)

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

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/).

Fraguent n grame

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