

General overview

Corpus	Analytics date	Language
tgk_Cyrl.jsonl.tsv	9/16/2024	Tajik (tg)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
1,261,259	24,851,003	14,469,071 (58.22 %)	770M	7.75 GB	4,565,968,108

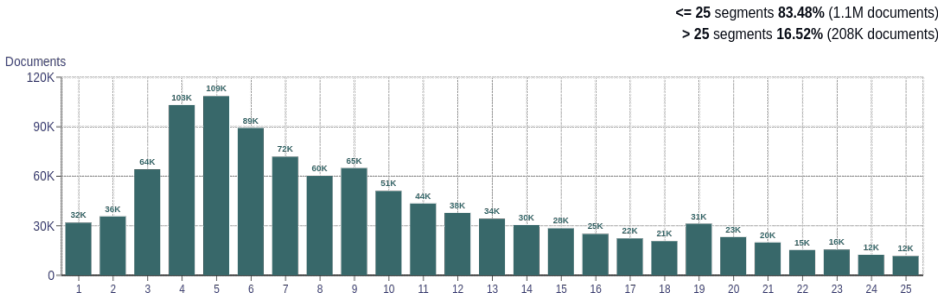
Top 10 domains

Domain	Docs	% of total
ozodi.org	123K	9.77
ozodlik.org	80K	6.36
kun.uz	52K	4.13
wikipedia.org	25K	1.97
ozodagon.com	21K	1.67
khovar.tj	18K	1.40
islom.uz	16K	1.25
fikr.uz	15K	1.18
qalampir.uz	15K	1.17
daryo.uz	13K	0.99

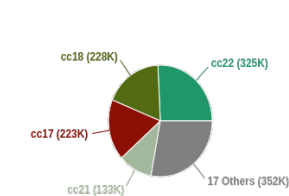
Top 10 TLDs

Domain	Docs	% of total
uz	501K	39.72
org	267K	21.19
com	194K	15.41
tj	188K	14.90
info	21K	1.66
ru	19K	1.52
net	13K	1.06
mobi	12K	0.97
asia	9.9K	0.79
kz	4.8K	0.38

Documents size (in segments)

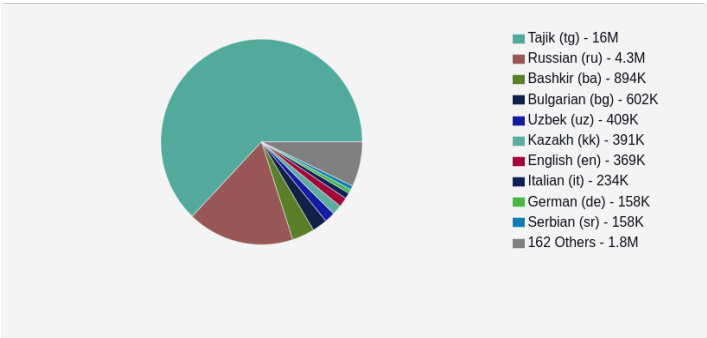


Documents by collection

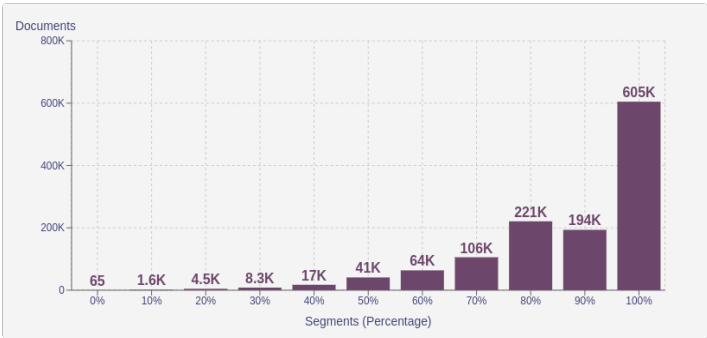


Language Distribution

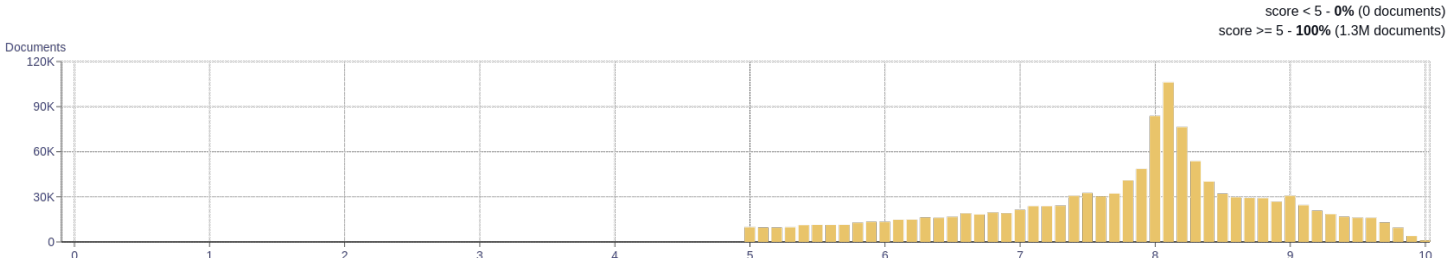
Number of segments



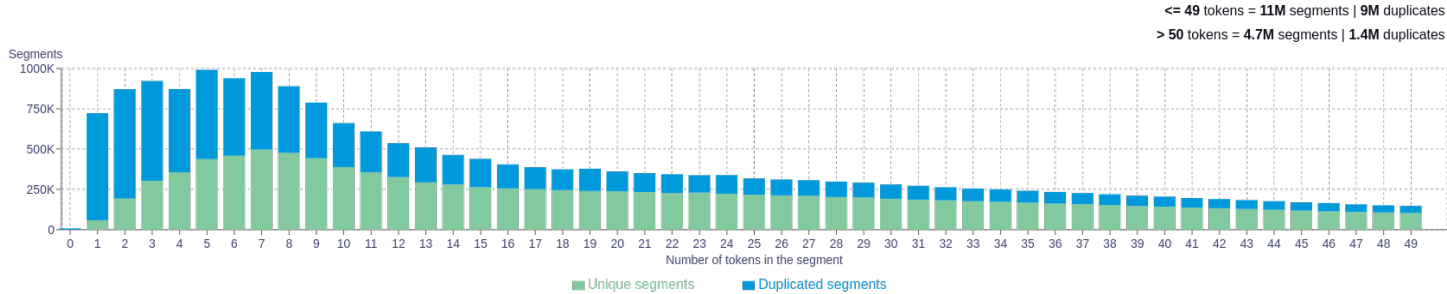
Percentage of segments in Tajik (tg) inside documents



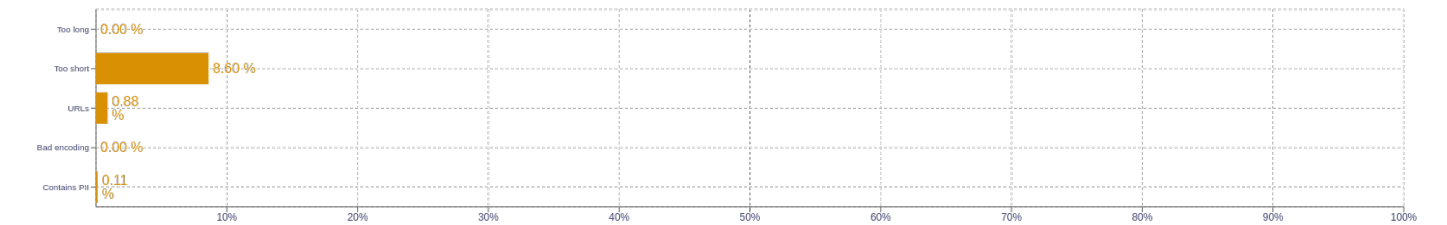
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>билан 3299836</div> <div>бу 2456944</div> <div>бир 2146323</div> <div>учун 2134996</div> <div>тоҷикистон 1655856</div>
2	<div>ҷумҳурии тоҷикистон 661595</div> <div>ўзбекистон республикаси 603614</div> <div>эмомалӣ раҳмон 170101</div> <div>шаҳри душанбе 159888</div> <div>соллаллоҳу алайҳи 147219</div>
3	<div>президенти ҷумҳурии тоҷикистон 124169</div> <div>вайп вайп вайп 82986</div> <div>муҳтарам эмомалӣ раҳмон 75996</div> <div>ўзбекистон республикаси олий 74878</div> <div>ҳукумати ҷумҳурии тоҷикистон 74181</div>
4	<div>вайп вайп вайп вайп 82660</div> <div>асосгузори сулҳу ваҳдати миллий 46294</div> <div>ўзбекистон республикаси вазирлар маҳкамасининг 46155</div> <div>ўзбекистон республикаси олий мажлиси 41447</div> <div>соллаллоҳу алайҳи ва саллам 41224</div>
5	<div>вайп вайп вайп вайп вайп 82394</div> <div>сағи пидорасу сағи пидорасу сағи 39003</div> <div>пидорасу сағи пидорасу сағи пидорасу 38904</div> <div>президенти ҷумҳурии тоҷикистон муҳтарам эмомалӣ 36225</div> <div>ҷумҳурии тоҷикистон муҳтарам эмомалӣ раҳмон 35809</div>

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 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>, , , 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>