

General overview

Corpus	Date	Language
kat_Georgjsonl.tsv	9/19/2024	Georgian (ka)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
3,335,164	63,722,098	29,708,949 (46.62 %)	1.6B	10,095,225,458	24.7 GB

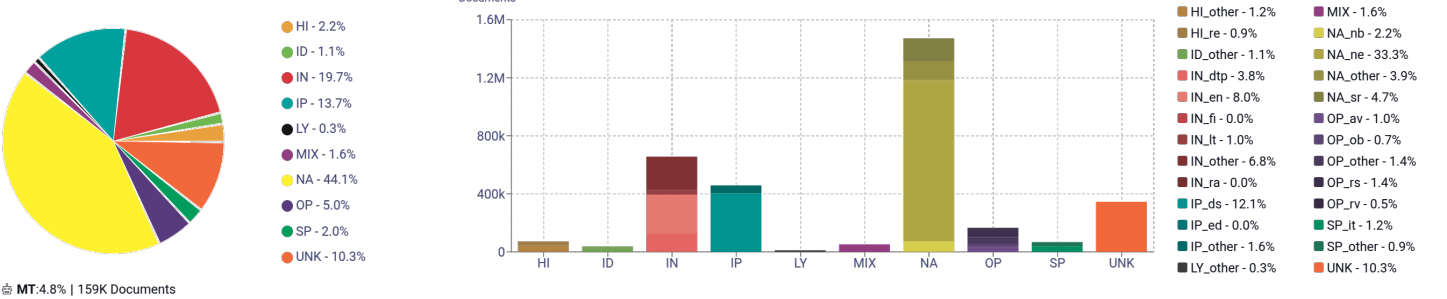
Top 10 domains

Domain	Docs	% of total
wikipedia.org	248K	7.44%
radiotavisupleb...	141K	4.24%
gancxadebebl.ge	96K	2.89%
wordpress.com	69K	2.05%
netgazeti.ge	34K	1.03%
1tv.ge	28K	0.83%
on.ge	27K	0.82%
sana.ge	26K	0.78%
blogspot.com	26K	0.78%
ambebi.ge	26K	0.77%

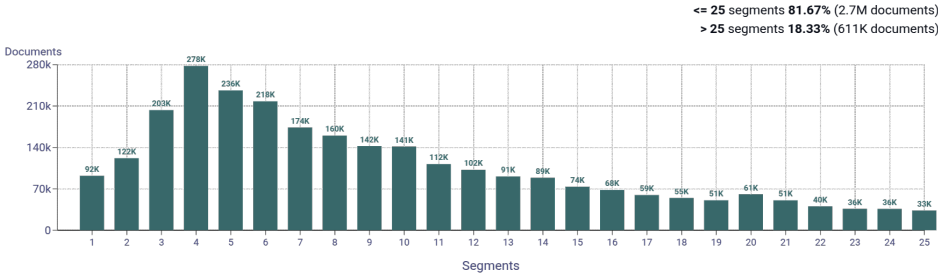
Top 10 TLDs

Domain	Docs	% of total
ge	2.1M	63.33%
com	495K	14.85%
org	313K	9.38%
gov.ge	111K	3.32%
net	68K	2.03%
edu.ge	51K	1.53%
com.ge	34K	1.03%
org.ge	21K	0.63%
am	16K	0.49%
info	15K	0.46%

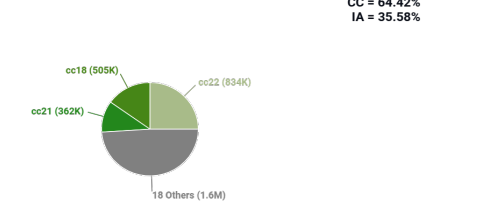
Register labels



Documents size (in segments)

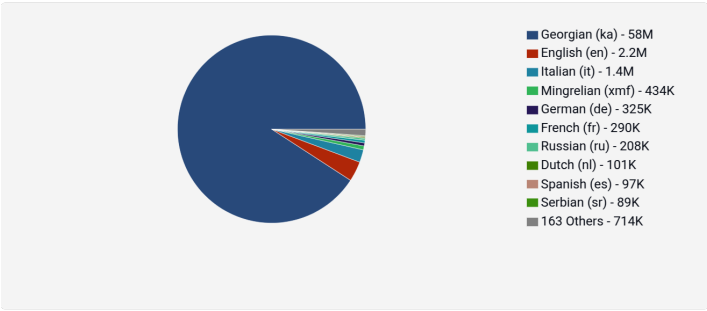


Documents by collection

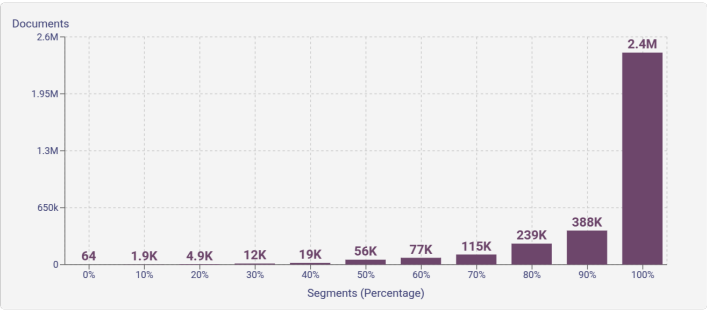


Language Distribution

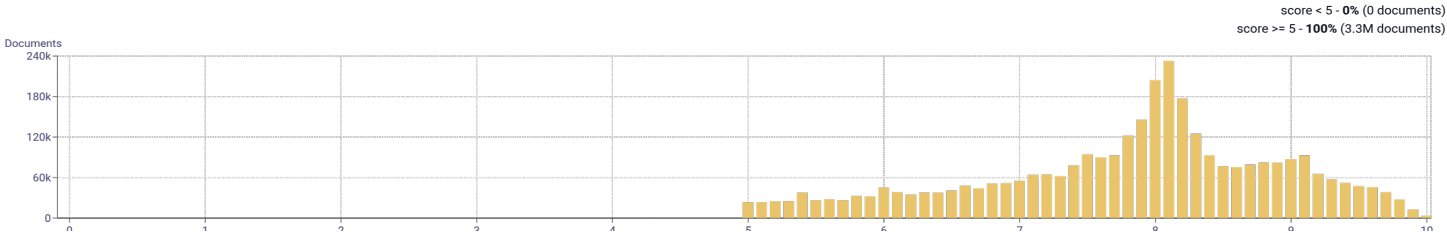
Number of segments in the Georgian (ka) corpus



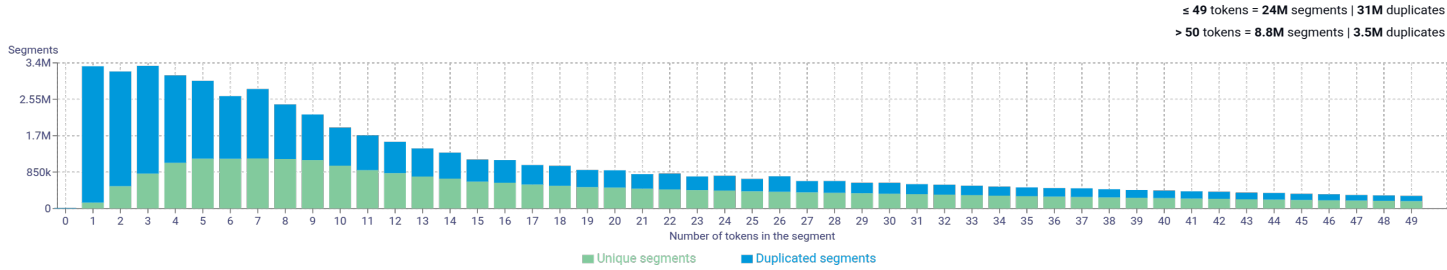
Percentage of segments in Georgian (ka) inside documents



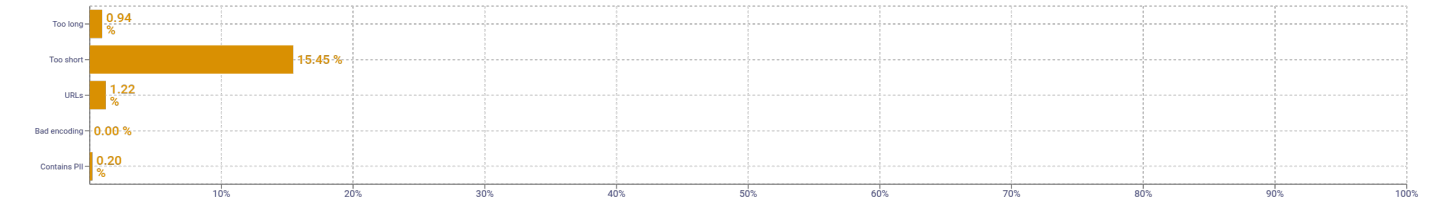
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	ეს 5971888 ამ 5686353 ის 5099032 არის 4570049 იყო 3701223
2	ეს არის 543370 წყაროს რედაქტირება 355269 რა ლქმა 306579 მიუხედავად იმისა 242247 წლის განმავლობაში 208765
3	განათლებისა და მეცნიერების 61610 ამა თუ იმ 61200 ეს არ არის 60572 ჟრანგი და ქართველი 56155 ბაზარზე იყო ორიენტირებული 56068
4	ჰინტის განვითარებისთვის საქართველოში რამდენიმე 56058 ძირითადად ჟრანგულ ბაზარზე იყო 56058 ჟრანგულ ბაზარზე იყო ორიენტირებული 56058 საქართველოში რამდენიმე არაკომერციული პროექტის 56058 კომპანია ძირითადად ჟრანგულ ბაზარზე 56058
5	ჰინტის განვითარებისთვის საქართველოში რამდენიმე არაკომერციული 56058 ძირითადად ჟრანგულ ბაზარზე იყო ორიენტირებული 56058 კომპანია ძირითადად ჟრანგულ ბაზარზე იყო 56058 განვითარებისთვის საქართველოში რამდენიმე არაკომერციული პროექტის 56058 წესს ჟრანგი და ქართველი დამფუძნებლების 56052

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 $\frac{\text{number of types (uniques)}}{\text{number of tokens}}$, 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/pablo16n/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	Encyclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				