

General overview

Corpus	Analytics date	Language
ka_1.jsonl.tsv	3/26/2024	Georgian (ka)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
533,070	65,524,284	14,489,425 (22.11 %)	769M	10.03 GB	

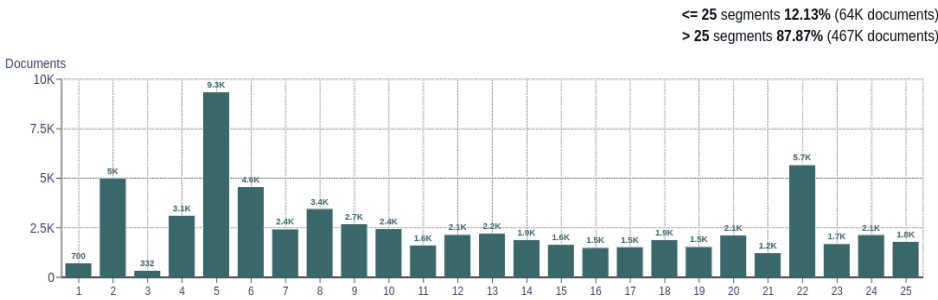
Top 10 domains

Domain	Docs	% of total
extra.ge	30K	5.67
forum.ge	12K	2.33
wikipedia.org	11K	2.14
chinashop.ge	8.4K	1.58
netgazeti.ge	8.2K	1.53
tabula.ge	7.2K	1.35
interpressnews.ge	5K	0.94
news.ge	4.6K	0.86
wordpress.com	4.5K	0.85
on.ge	4.4K	0.83

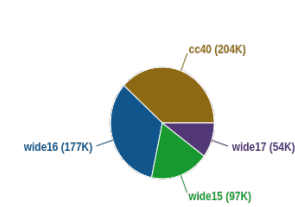
Top 10 TLDs

Domain	Docs	% of total
ge	365K	68.52
com	65K	12.19
org	29K	5.53
net	12K	2.19
gov.ge	7.6K	1.43
com.ge	7.1K	1.32
edu.ge	4.2K	0.79
ch	4K	0.75
org.ge	4K	0.75
info	2.8K	0.52

Documents size (in segments)

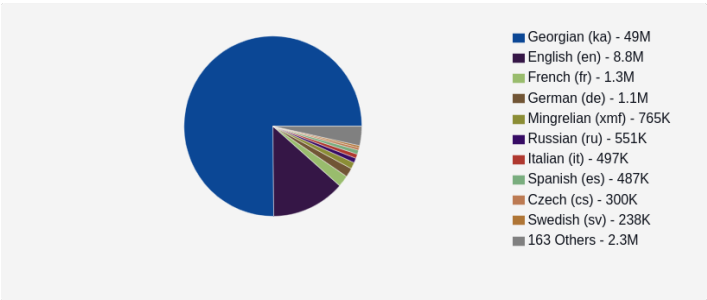


Documents by collection

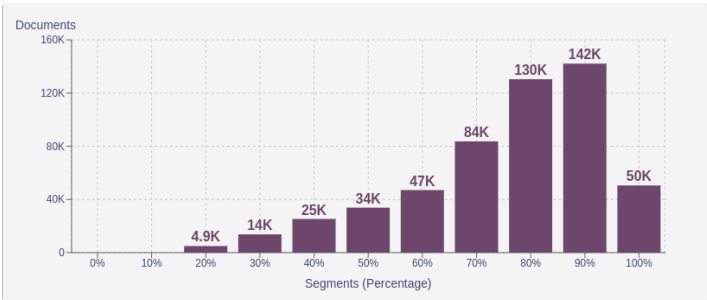


Language Distribution

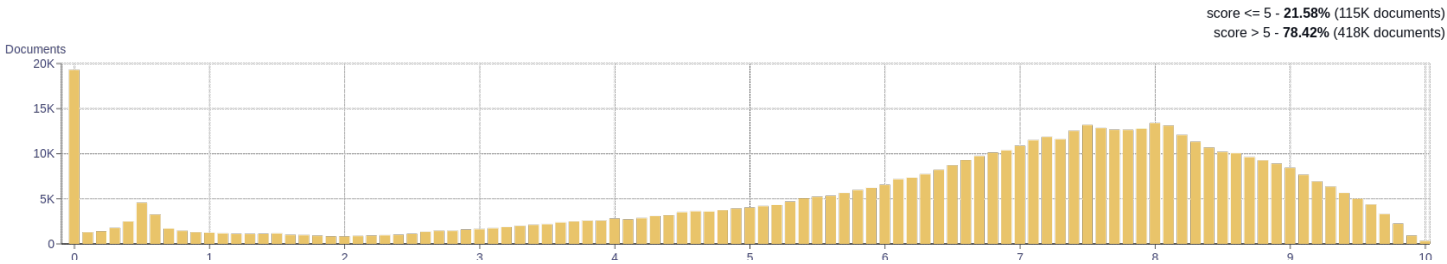
Number of segments



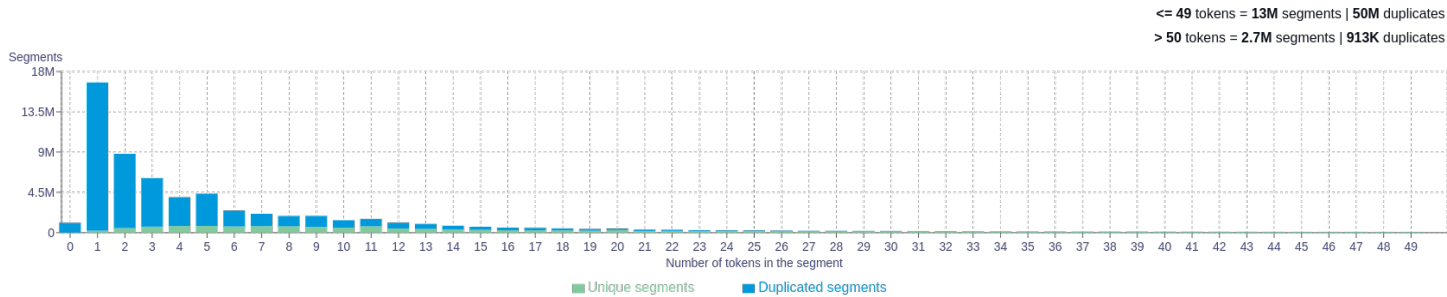
Percentage of segments in Georgian (ka) inside documents



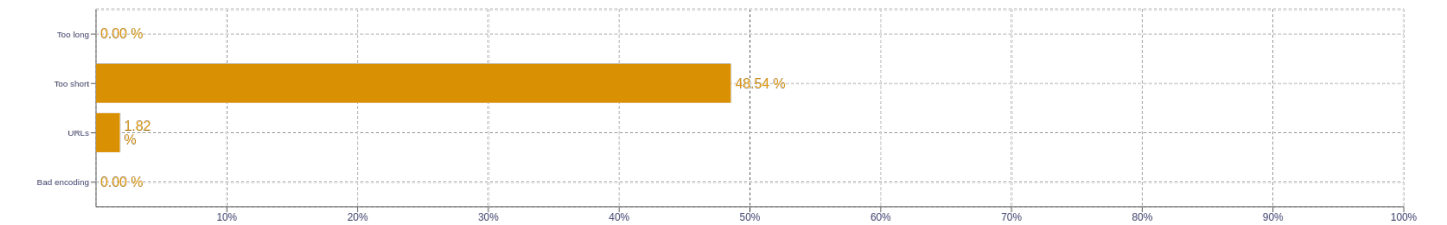
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>ის 1799198</div> <div>ეს 1790348</div> <div>ამ 1727935</div> <div>არის 1460367</div> <div>the 1280215</div>
2	<div>posted by 341878</div> <div>span style 244391</div> <div>ახალი ამბები 211577</div> <div>of the 176848</div> <div>ნაღდი ანგარიშსწორების 158773</div>
3	<div>თბილისის მასშტაბით სრულიად 97511</div> <div>ნაღდი ანგარიშსწორების სურვილის 96448</div> <div>ანგარიშსწორების სურვილის შემთხვევაში 96448</div> <div>შეავსეთ მარტივი ღირმა 96386</div> <div>ლიდაკს და შეავსეთ 96102</div>
4	<div>ნაღდი ანგარიშსწორების სურვილის შემთხვევაში 96448</div> <div>ლიდაკს და შეავსეთ მარტივი 96099</div> <div>ჩვენი კურიერი ადგილზე მოგაწვდით 79409</div> <div>კურიერი ადგილზე მოგაწვდით პროდუქციას 79409</div> <div>მიწოდება თბილისის მასშტაბით სრულიად 69051</div>
5	<div>ლიდაკს და შეავსეთ მარტივი ღირმა 96099</div> <div>ჩვენი კურიერი ადგილზე მოგაწვდით პროდუქციას 79409</div> <div>მიწოდება თბილისის მასშტაბით სრულიად უფასოა 69051</div> <div>თამაში და გართობა არასდროს სრულდება 62522</div> <div>დააკლიკე და გადახაველ სათამაშოების საუკეთესო 62522</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>