

General overview

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

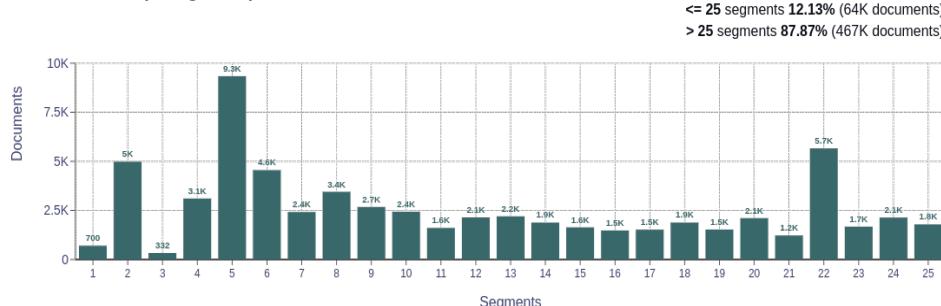
Volumes

Docs	Segments	Unique segments	Tokens	Size
533,070	65,524,284	53,749 (0.08 %)	769M	10.03 GB

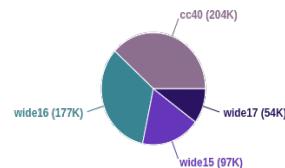
Type-Token Ratio

Georgian (ka)
0.01

Documents size (in segments)

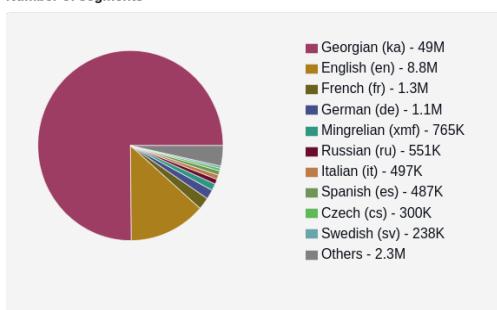


Documents by collection

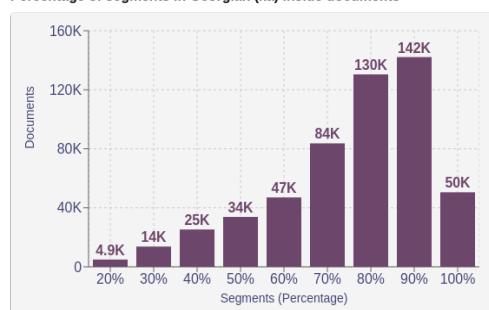


Language Distribution

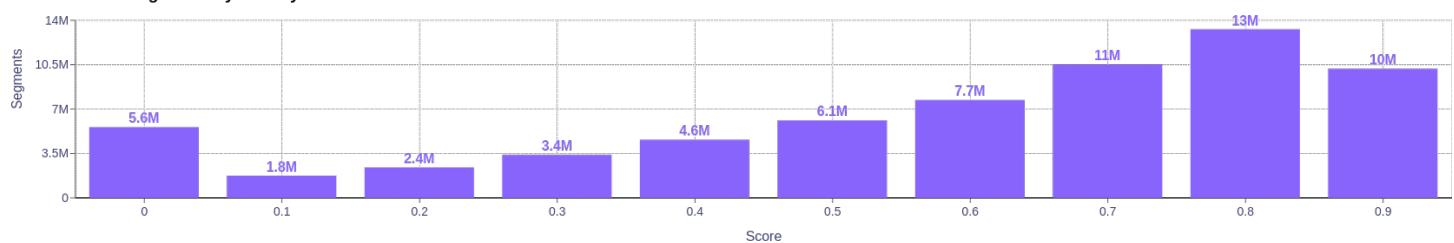
Number of segments



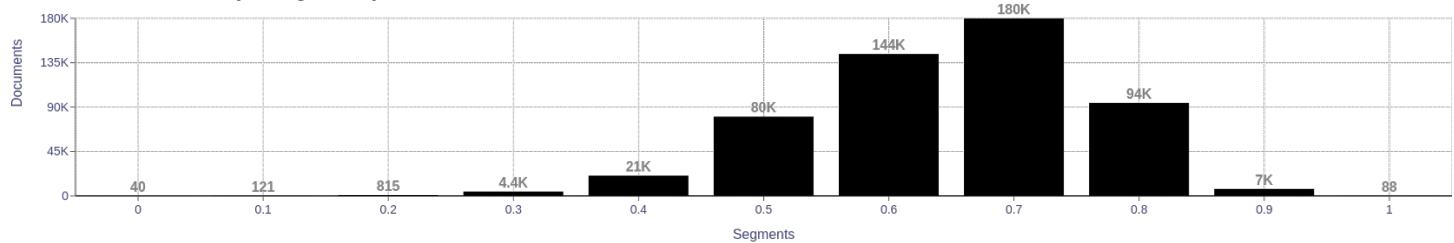
Percentage of segments in Georgian (ka) inside documents



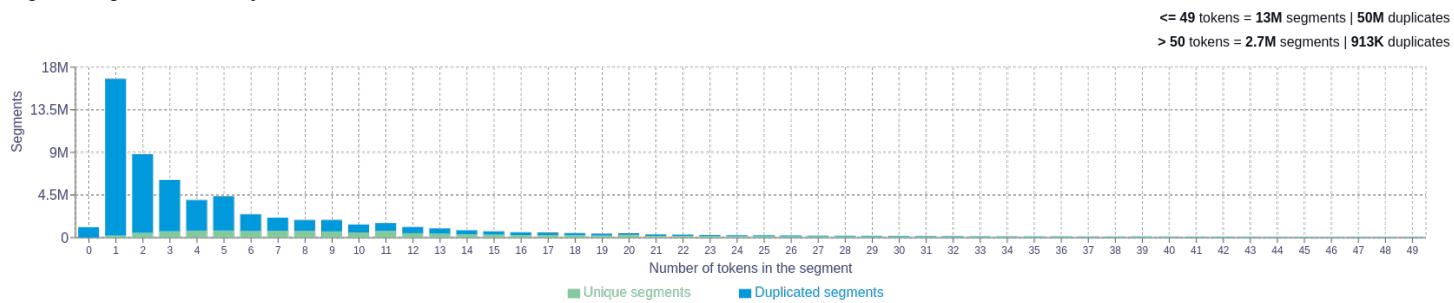
Distribution of segments by fluency score



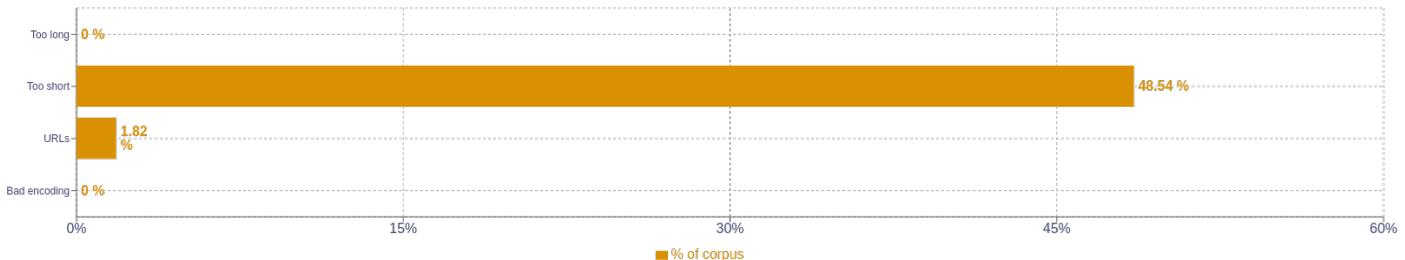
Distribution of documents by average fluency score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	ის 1799198 ქს 1790348 ამ 1727935 არის 1460367 the 1280215
2	posted by 341878 span style 244391 მხდარი ამბეჭით 211577 of the 176848 ნაღმი აწერისშორების 158773
3	თბილისის მასშტაბით სრულიად 97511 ნაღმი აწერისშორების სურვილის 96448 ანგარიშსწორების სურვილის შემთხვევაში 96448 შეაცემ მარტივი ფორმა 96306 ღიაბაჟ და შეავსეთ 96102
4	ნაღმი აწერისშორების სურვილის შემთხვევაში 96448 ღიაბაჟ და შეავსეთ მარტივი 96099 შევნი კურიერი აღვიძე მოგაწვდით 79409 გურიერი აღვიძე მოგაწვდით პროდუქტიას 79409 მინიჭება თბილისი მასშტაბით სრულიად 69051
5	ღიაბაჟ და შეავსეთ მარტივი ფორმა 96099 შევნი კურიერი აღვიძე მოგაწვდით პროდუქტიას 79409 მინიჭება თბილისის მასშტაბით სრულიად უფასოა 69051 თამაში და გარობაზე არასტრის სრულდება 62522 დაბადებულ და გადახავდან სათამაშოების საკუთრების 62522

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.clinfo.com/wp-content/uploads/2014/01/type-token-ratio.pdf>)

Document size (in segments)

Document size (in segments) Segments correspond to paragraph and list boundaries as defined by HTML elements (`<pre>`, ``, ``, etc.) replaced by newlines.

Segments correspond to

Language distribution

Language identified with FastSpell (<https://git>

Distribution of segments by fluency score

Obtained with Monocleaner (<https://github.com/monocleaner/monocleaner>)

Distribution of documents by average fluency score

Obtained with Monocleaner (I)

Segment length distribution by type

Tokenized with http

Frequent n-grams
Taken from https://github.com/bolt-project/data_analytics_toolkit/main/telusko_nlp_info.indd, after removing n-grams starting or ending in a stopword. Stopwords from https://github.com/bolt-project/data_analytics_toolkit/main/telusko_stopwords.txt.