

General overview

Corpus	Analytics date	Language
som_Latn.jsonl.tsv	9/23/2024	Somali (so)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
966,507	16,384,689	8,400,746 (51.27 %)	434M	2.39 GB	2,549,211,220

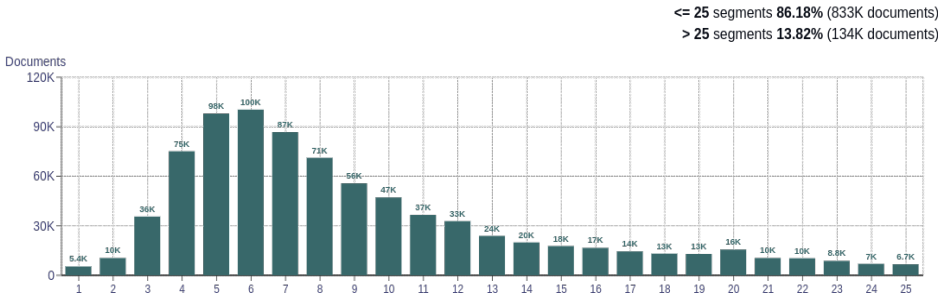
Top 10 domains

Domain	Docs	% of total
somaliitalk.com	36K	3.70
voasomali.com	17K	1.78
caasimada.net	14K	1.41
dunidaonline.com	13K	1.39
somaliand.org	13K	1.34
puntlandi.com	13K	1.31
goobjoog.com	12K	1.27
wikipedia.org	12K	1.21
wordpress.com	11K	1.18
goolfm.net	11K	1.13

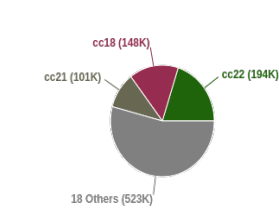
Top 10 TLDs

Domain	Docs	% of total
com	699K	72.28
net	154K	15.93
org	52K	5.34
so	16K	1.70
se	6.2K	0.65
online	5.5K	0.57
ca	4.2K	0.43
info	3.4K	0.35
is	2.8K	0.29
fi	2.2K	0.23

Documents size (in segments)

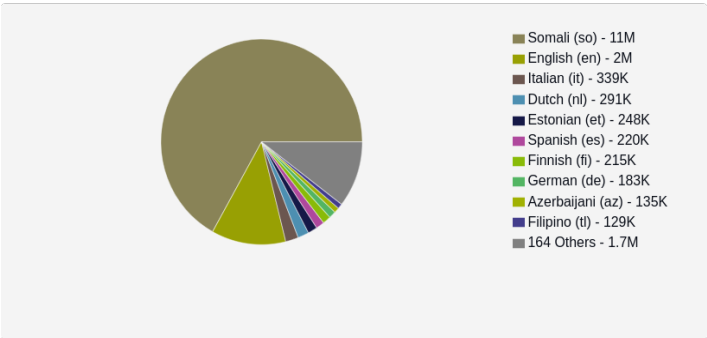


Documents by collection

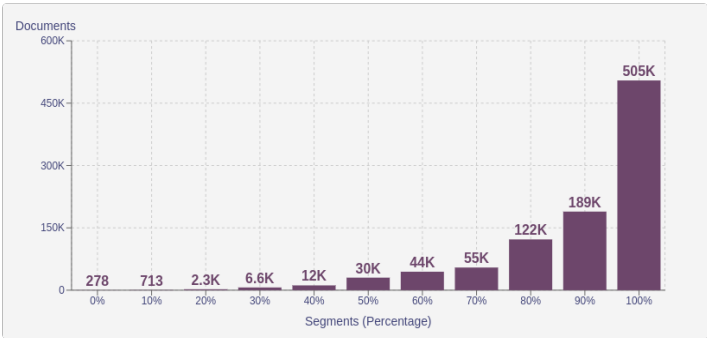


Language Distribution

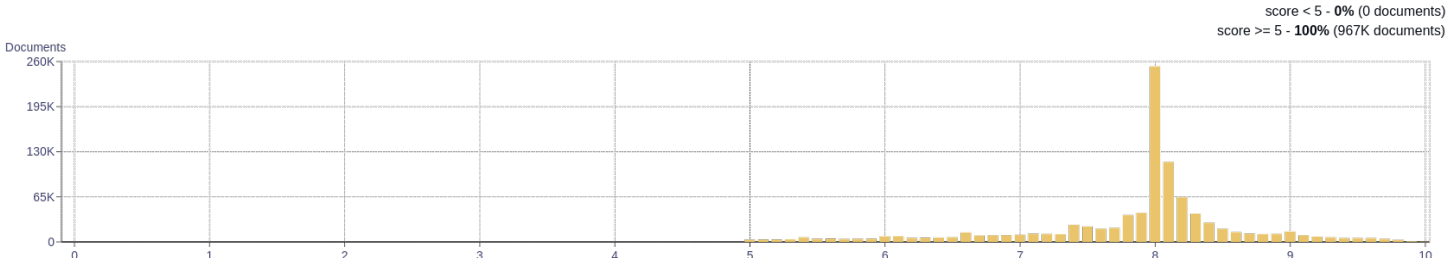
Number of segments



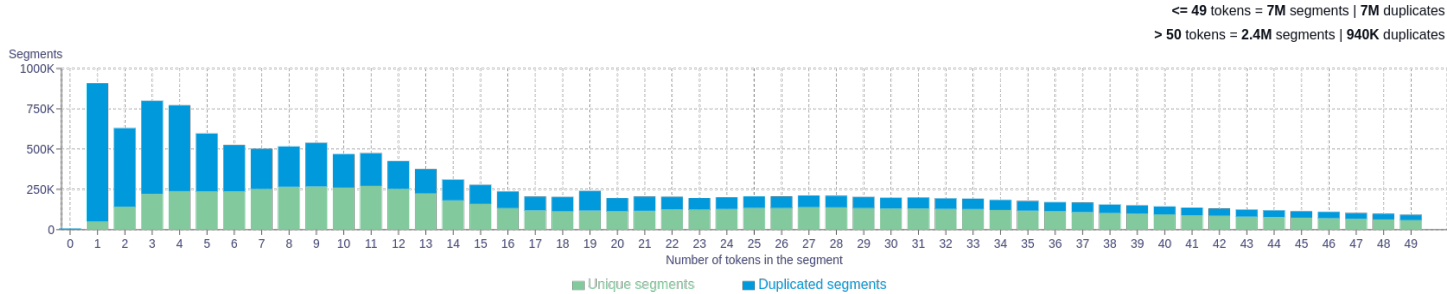
Percentage of segments in Somali (so) inside documents



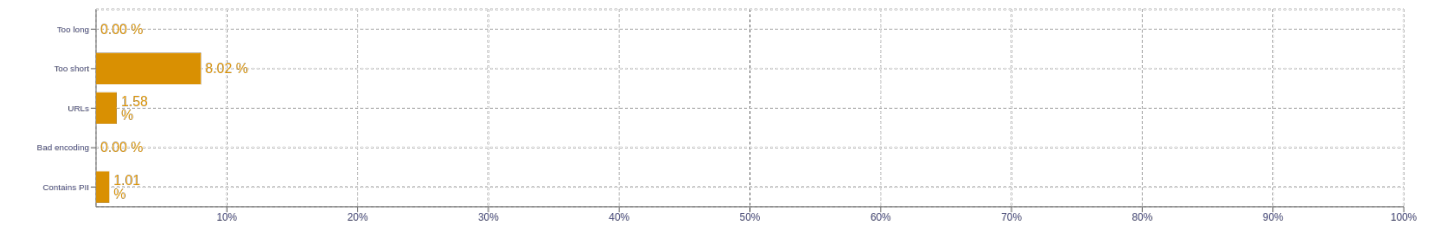
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>iyo 7052255</div> <div>ah 6346633</div> <div>ee 6345824</div> <div>u 5569995</div> <div>la 4116900</div>
2	<div>mid ah 577476</div> <div>ah ee 546674</div> <div>kala duwan 254500</div> <div>magaalada muqdisho 206322</div> <div>ee dalka 197778</div>
3	<div>qaar ka mid 93434</div> <div>mid ka mid 88474</div> <div>waxaa ka mid 63976</div> <div>kala duwan ee 62056</div> <div>ee magaalada muqdisho 46958</div>
4	<div>qaar ka mid ah 88647</div> <div>mid ka mid ah 82206</div> <div>waxaa ka mid ah 45806</div> <div>reer binu israa 'iil 24524</div> <div>ah oo ku saabsan 17948</div>
5	<div>badan oo ka mid ah 11966</div> <div>madaxweynaha jamhuuriyadda federaalka soomaaliya mudane 8169</div> <div>soomaaliya mudane xasan sheekh maxamuud 7434</div> <div>madaxweynaha soomaaliya xasan sheekh maxamuud 6727</div> <div>sida uu hadalka u dhigay 6500</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>