

General overview

Corpus	Analytics date	Language
gl_1.jsonl.tsv	3/21/2024	Galician (gl)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
731,356	92,682,759	19,118,996 (20.63 %)	1B	5.03 GB	

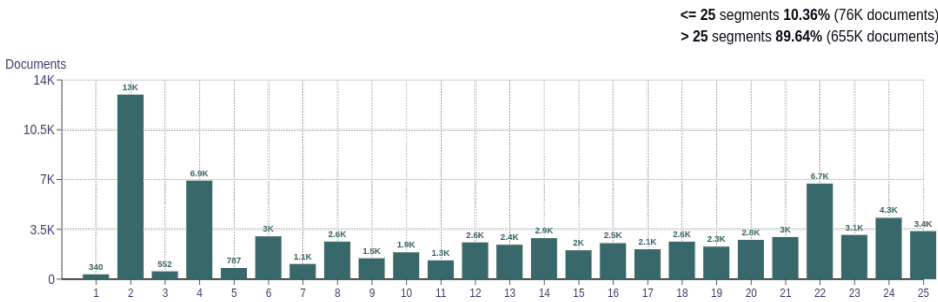
Top 10 domains

Domain	Docs	% of total
blogspot.com.es	76K	10.34
wikipedia.org	30K	4.11
blogspot.com	28K	3.79
xunta.gal	12K	1.67
lugo.gal	11K	1.49
wordpress.com	10K	1.41
pontevedraviva.com	6.5K	0.89
bretemas.gal	6.4K	0.88
vigo.org	5.8K	0.80
galipedia.com	5.4K	0.74

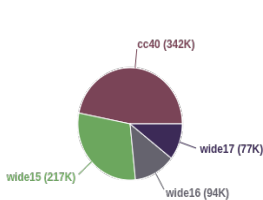
Top 10 TLDs

Domain	Docs	% of total
com	244K	33.34
gal	138K	18.85
org	115K	15.77
es	98K	13.36
com.es	76K	10.36
info	10K	1.38
eu	7.8K	1.06
net	7.6K	1.04
com.ar	5.9K	0.81
pt	3.1K	0.42

Documents size (in segments)

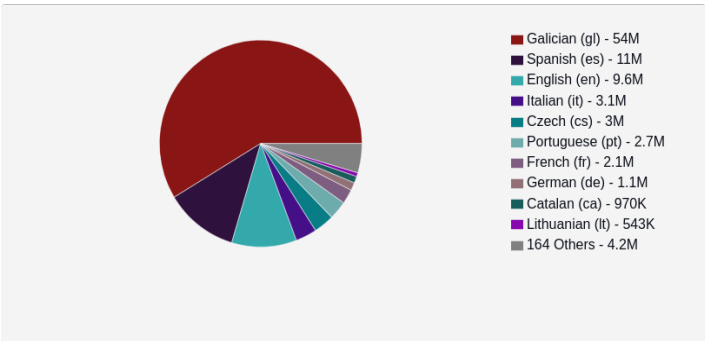


Documents by collection

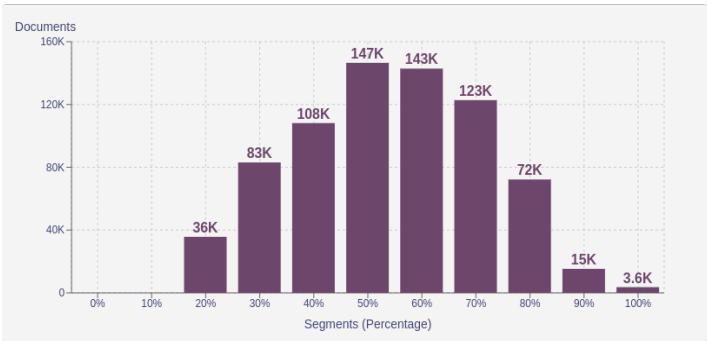


Language Distribution

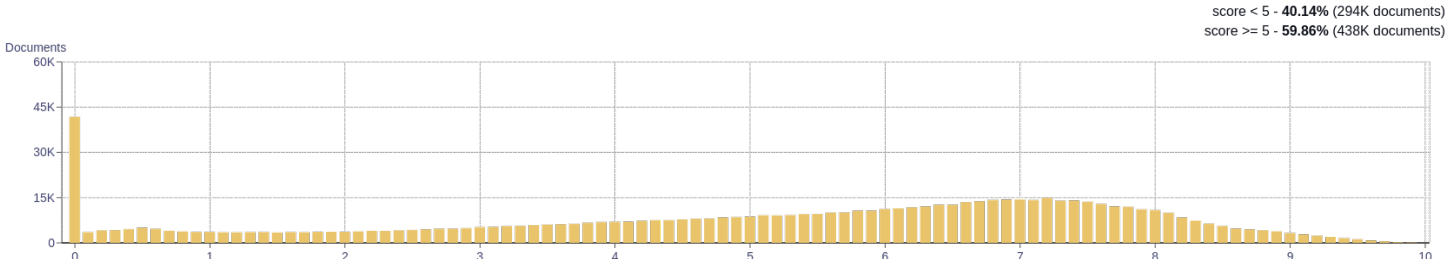
Number of segments



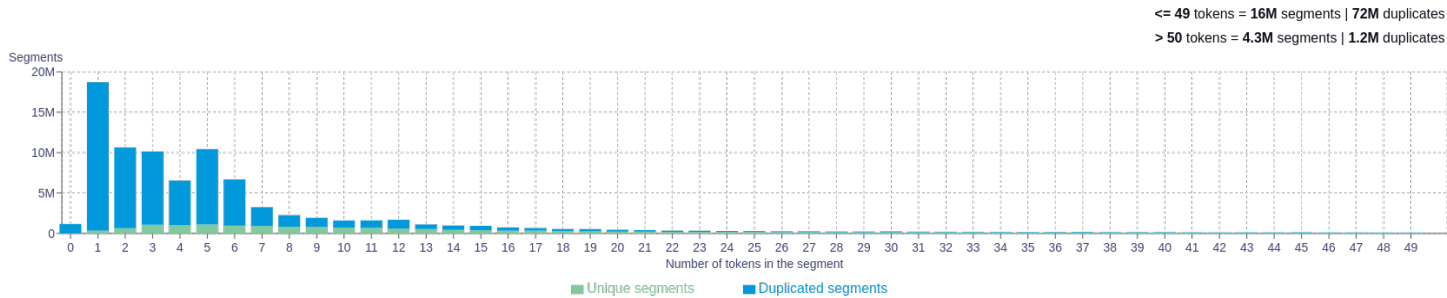
Percentage of segments in Galician (gl) inside documents



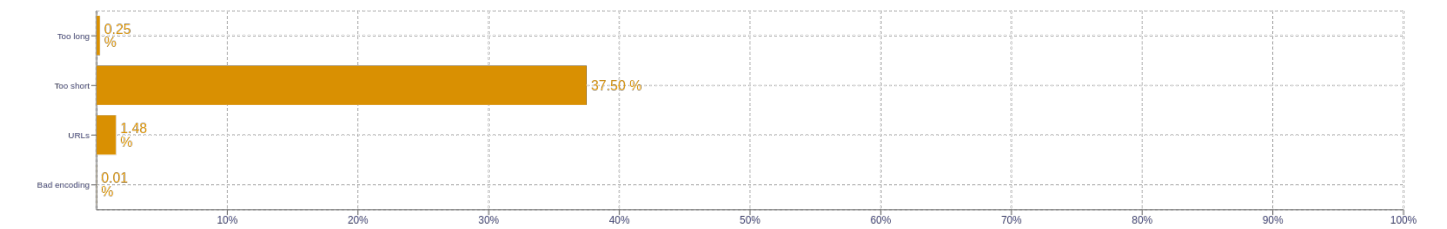
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	y 2371096galicia 1898093día 1110181abril 1085104marzo 1051215
2	hay comentarios 300026correo electrónicoescribe 292932correo electrónico 231870aviso legal 220674sitio web 215834
3	enviar por correo 450385facebookcompartir en pinterest 433535blogcompartir con twittercompartir 292942electrónicoescribe un blogcompartir 292931twittercompartir con facebookcompartir 289851
4	enviar por correo electrónicoescribe 292932correo electrónicoescribe un blogcompartir 292931enviar por correo electrónicoblogthis 143634enlaces a esta entrada 68538entradas antiguas página principal 53522
5	electrónicoescribe un blogcompartir con twittercompartir 292931blogcompartir con twittercompartir con facebookcompartir 289851twittercompartir con facebookcompartir en pinterest 289850compartir en twittercompartir en facebookcompartir 143696twittercompartir en facebookcompartir en pinterest 143685

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.stinfo.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>