

General overview

Corpus	Analytics date	Language
epo_Latn.jsonl.tsv	9/16/2024	Esperanto (eo)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
818,878	20,353,314	7,149,533 (35.13 %)	571M	2.81 GB	2,956,534,128

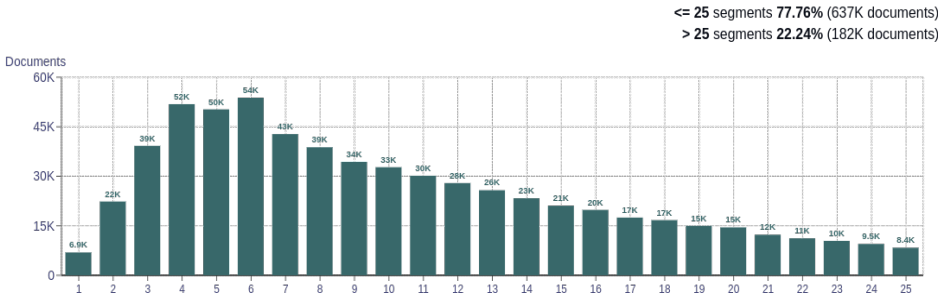
Top 10 domains

Domain	Docs	% of total
wikipedia.org	538K	65.69
blogspot.com	9.9K	1.21
wikitrans.net	8.3K	1.02
esperantio.net	7.2K	0.88
pola-retradio.org	5.7K	0.70
wordpress.com	5.7K	0.70
espero.com.cn	4.6K	0.57
over-blog.com	3.8K	0.46
ikso.net	3.6K	0.44
uea.org	3.6K	0.44

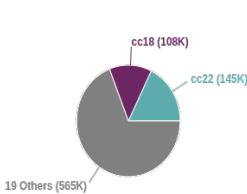
Top 10 TLDs

Domain	Docs	% of total
org	595K	72.69
com	99K	12.15
net	39K	4.74
ru	10K	1.22
cn	6.5K	0.80
info	6.5K	0.79
be	5.2K	0.63
com.cn	4.7K	0.57
de	4.6K	0.56
eu	4.5K	0.55

Documents size (in segments)

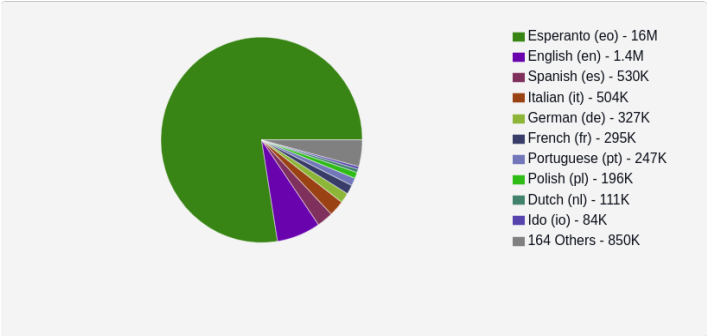


Documents by collection

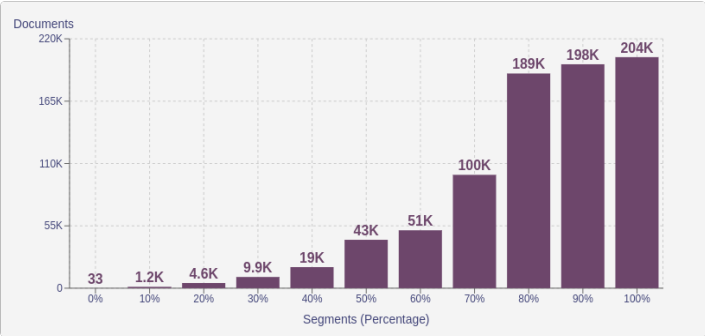


Language Distribution

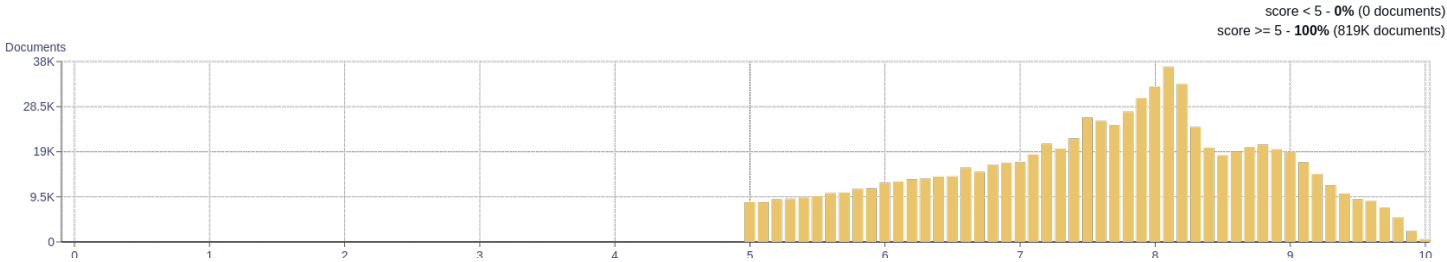
Number of segments



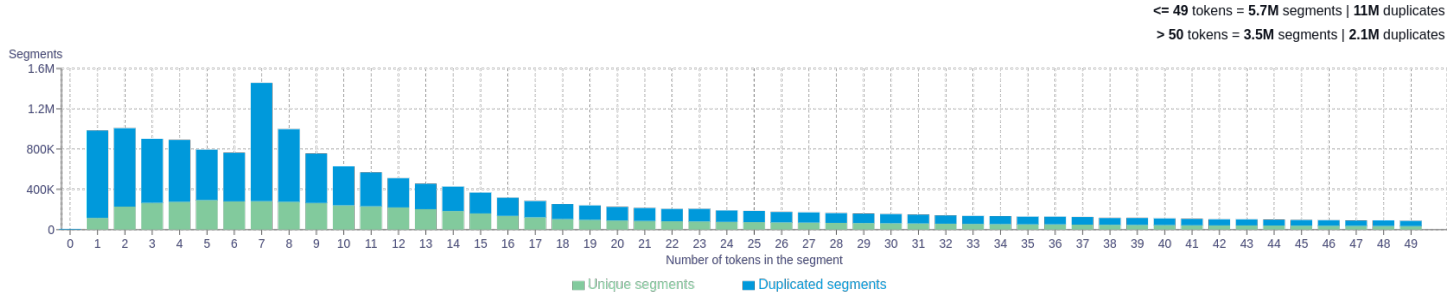
Percentage of segments in Esperanto (eo) inside documents



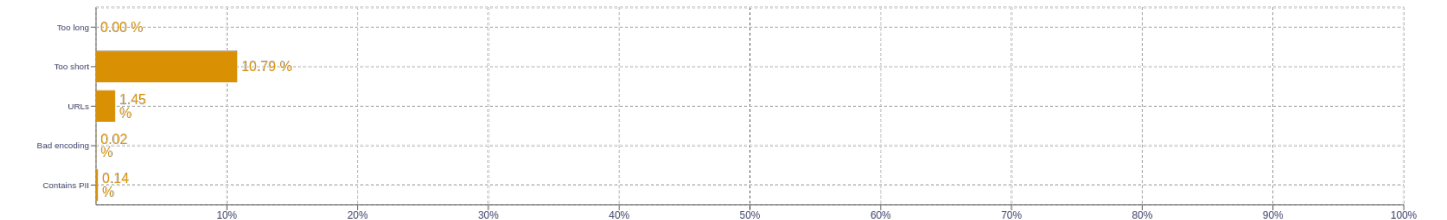
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>redakti 3013130</div> <div>kun 2377814</div> <div>pri 2275296</div> <div>kiel 2189232</div> <div>el 2157729</div>
2	<div>redakti fonton 1405770</div> <div>povas esti 202460</div> <div>of the 118446</div> <div>eksteraj ligiloj 113827</div> <div>temas pri 108717</div>
3	<div>iom post iom 30695</div> <div>per la retarkivo 26652</div> <div>retarkivo wayback machine 26472</div> <div>el la jaro 25792</div> <div>ekde la jaro 17807</div>
4	<div>per la retarkivo wayback 26569</div> <div>arkivita el la originalo 14507</div> <div>archived from the original 11946</div> <div>from the original on 11183</div> <div>el la plej gravaj 10845</div>
5	<div>per la retarkivo wayback machine 26472</div> <div>archived from the original on 11114</div> <div>ankaŭ en la vikimedia komunejo 8898</div> <div>vidu ankaŭ en la vikimedia 8880</div> <div>kolekto de bildoj kaj plurmediaj 8852</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>