

General overview

Corpus	Analytics date	Language
lij_Latn.jsonl.tsv	9/25/2024	Ligurian (lij)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
8,371	157,715	76,776 (48.68 %)	7.4M	33.67 MB	31,311,759

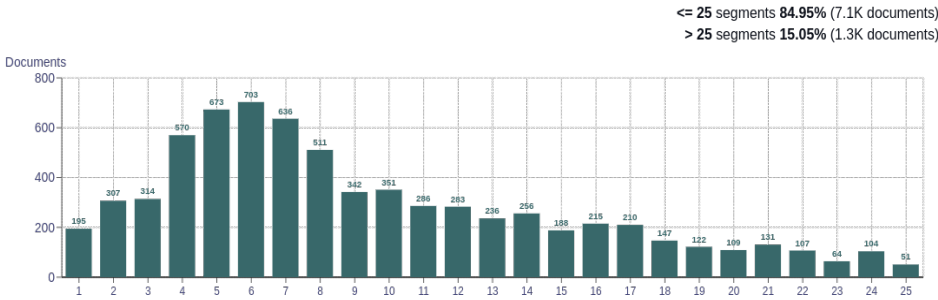
Top 10 domains

Domain	Docs	% of total
wikipedia.org	6.7K	80.48
wikisource.org	126	1.51
ilsecoloxix.it	91	1.09
bible.is	78	0.93
primocanale.it	49	0.59
blogspot.com	36	0.43
eodishasamachar.com	36	0.43
in-links.eu	36	0.43
genovapress.com	34	0.41
bywiki.com	29	0.35

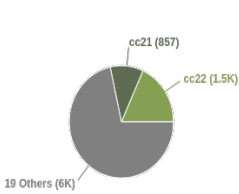
Top 10 TLDs

Domain	Docs	% of total
org	7K	83.97
com	532	6.36
it	356	4.25
net	86	1.03
is	78	0.93
eu	50	0.60
si	20	0.24
com.ar	17	0.20
ru	14	0.17
com.br	14	0.17

Documents size (in segments)

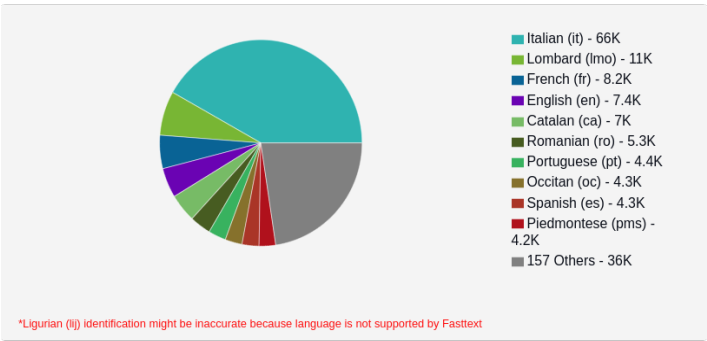


Documents by collection

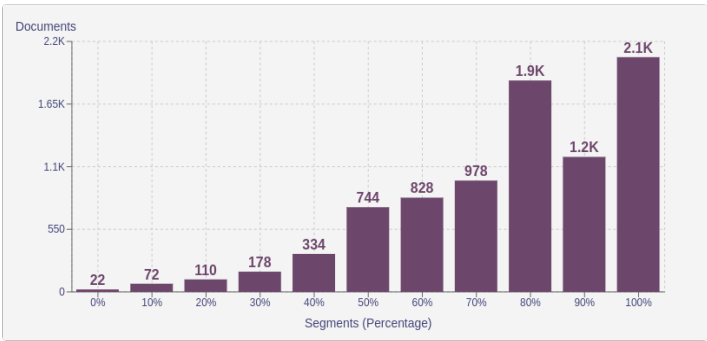


Language Distribution

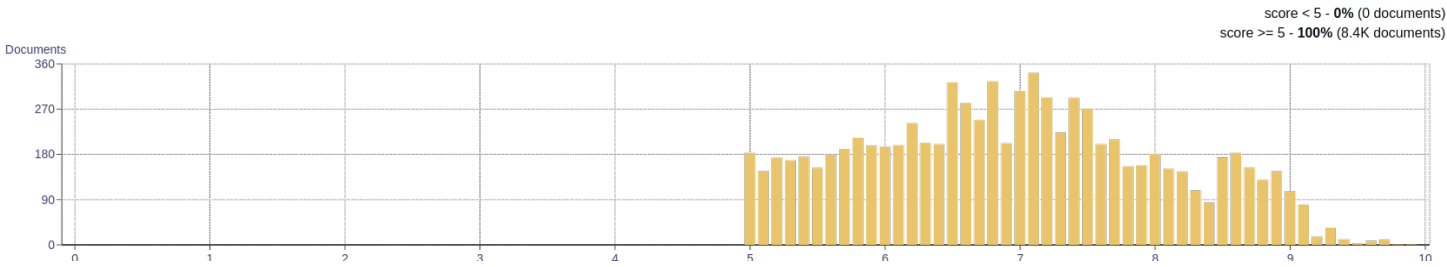
Number of segments



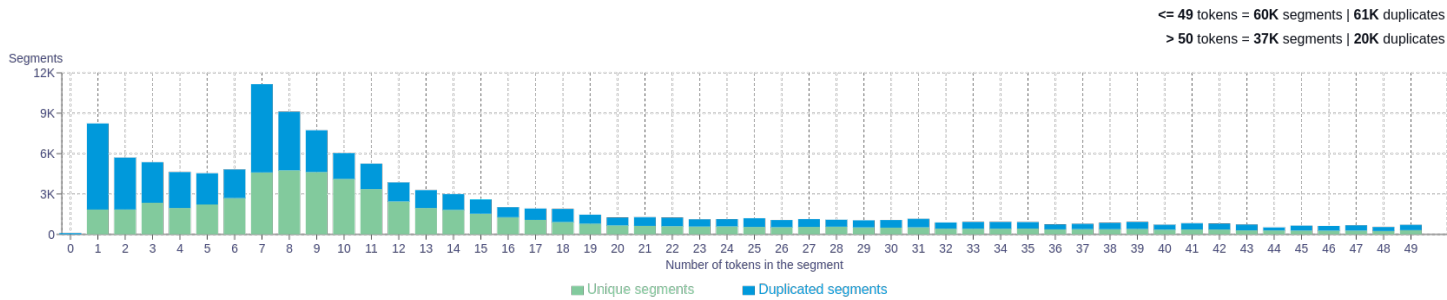
Percentage of segments in Ligurian (lij) inside documents



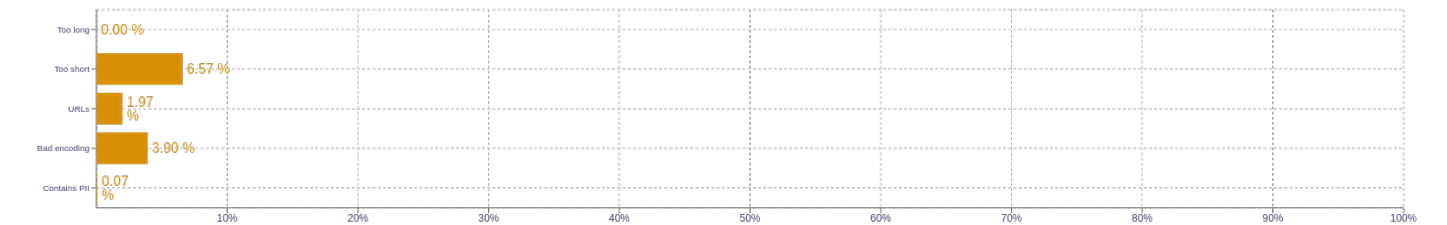
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>l 167155</div> <div>d 45088</div> <div>n 37951</div> <div>è 34758</div> <div>u 32327</div>
2	<div>modifica wikipèdista 8239</div> <div>modifica wikipèdista 4938</div> <div>u l 3034</div> <div>url consultà 2327</div> <div>u n 2272</div>
3	<div>ò ò ò 1852</div> <div>commons a contègne 1733</div> <div>p e r 1148</div> <div>c o n 1092</div> <div>z io n 966</div>
4	<div>ò ò ò ò 1848</div> <div>wikimedia commons a contègne 1733</div> <div>innàgine ò di àtri 760</div> <div>zeauriniàle zeauriniàle zeauriniàle zeauriniàle 602</div> <div>miou miou miou miou 434</div>
5	<div>ò ò ò ò ò 1844</div> <div>commons a contègne di files 934</div> <div>innàgine ò di àtri files 760</div> <div>commons a contègne di innàgine 740</div> <div>zeauriniàle zeauriniàle zeauriniàle zeauriniàle zeauriniàle 592</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>