

General overview

Corpus	Analytics date	Language
pap_Latn.jsonl.tsv	12/3/2024	Papiamento (pap)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
89,812	1,387,382	814,803 (58.73 %)	53M	244.39 MB	252,796,043

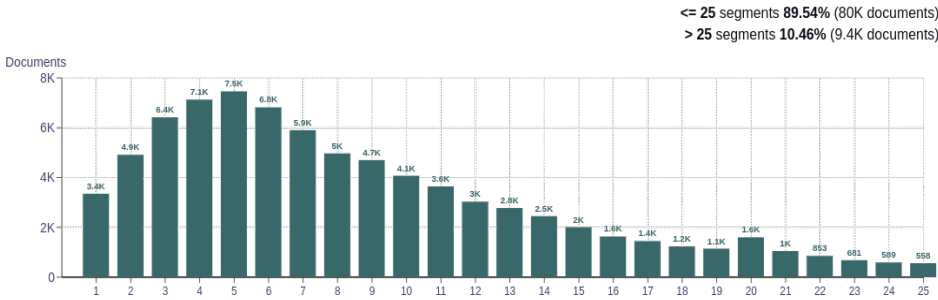
Top 10 domains

Domain	Docs	% of total
diario.aw	4.8K	5.30
kikotapasando.com	4.3K	4.74
wikipedia.org	4.2K	4.67
masnoticia.com	4.2K	4.66
arubanative.com	3.8K	4.26
live99fm.com	3.4K	3.75
awe24.com	3.3K	3.68
noticiacia.com	2.6K	2.94
awemainta.com	2.2K	2.41
1noticia.com	1.6K	1.83

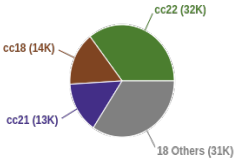
Top 10 TLDs

Domain	Docs	% of total
com	57K	63.50
aw	10K	11.08
org	9.4K	10.43
cw	3.8K	4.22
nl	3.4K	3.81
net	1.6K	1.78
nu	1.4K	1.51
news	1.2K	1.31
today	276	0.31
blog	230	0.26

Documents size (in segments)

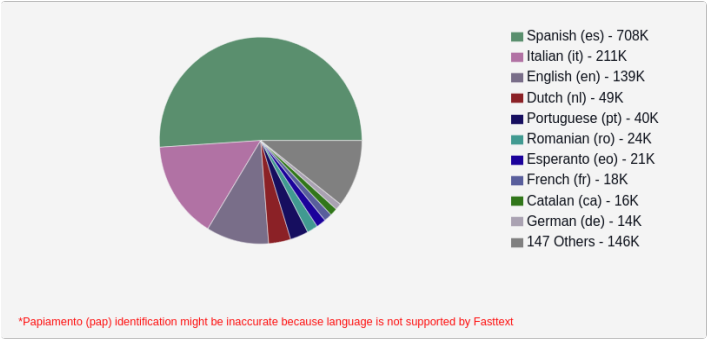


Documents by collection

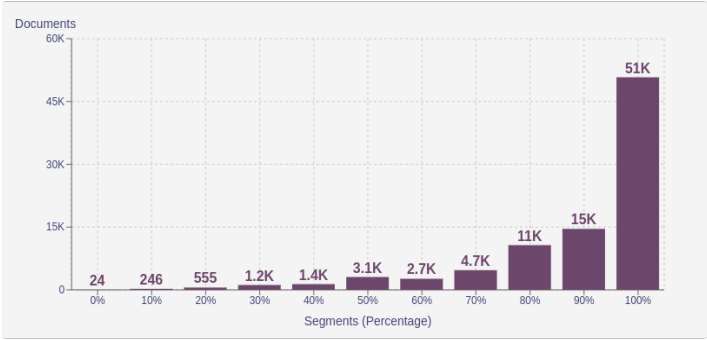


Language Distribution

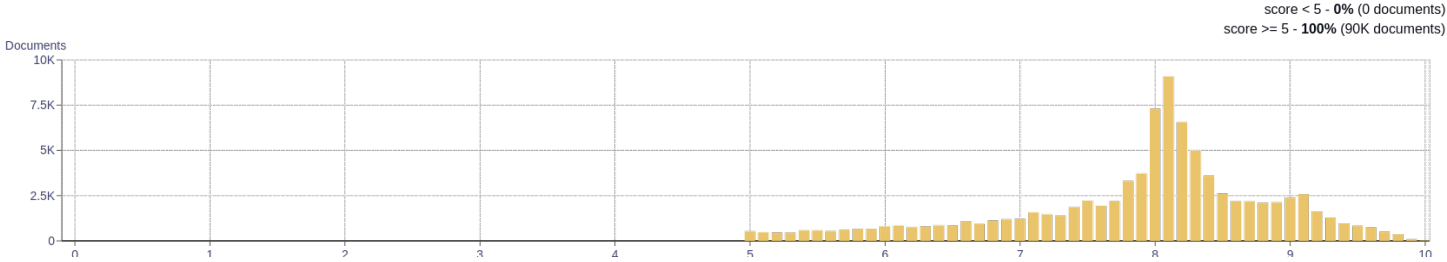
Number of segments



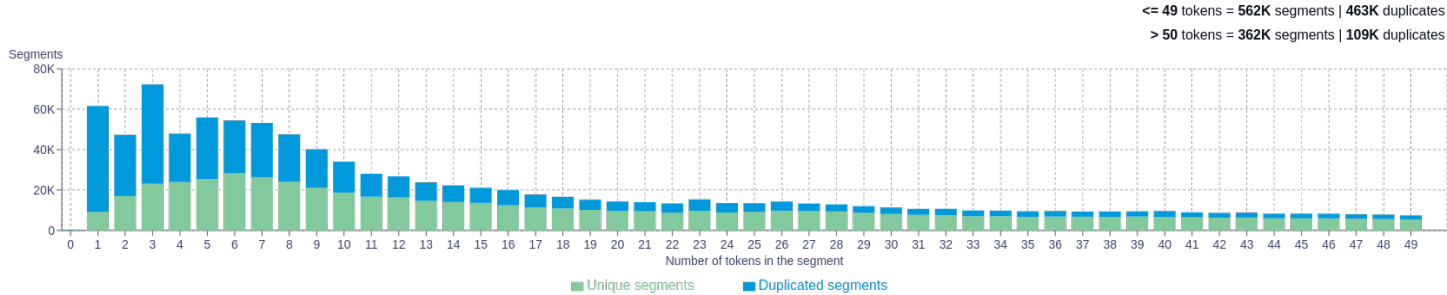
Percentage of segments in Papiamento (pap) inside documents



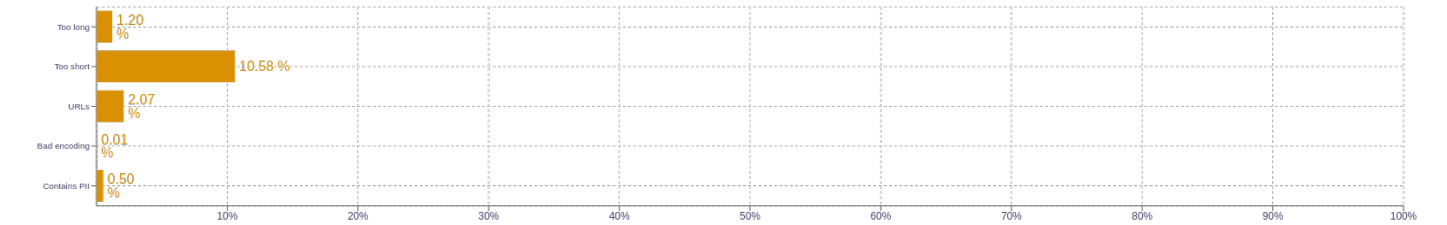
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>pa 1117468</div> <div>na 749639</div> <div>y 733878</div> <div>i 474986</div> <div>nan 381363</div>
2	<div>na aruba 35903</div> <div>pa asina 23408</div> <div>pa nan 21416</div> <div>su mes 19814</div> <div>tur hende 16545</div>
3	<div>na e sitio 7627</div> <div>el a bisa 6354</div> <div>gobierno di aruba 5318</div> <div>loke ta trata 5225</div> <div>na un manera 4495</div>
4	<div>pa loke ta trata 5009</div> <div>prueba prueba prueba prueba 2697</div> <div>na e momentonan aki 2015</div> <div>prome minister evelyn wever 2009</div> <div>yegada di e patruya 1518</div>
5	<div>prueba prueba prueba prueba prueba 2559</div> <div>na yegada di e patruya 1509</div> <div>camara di comercio y industria 1227</div> <div>impuesto di vehiculo di motor 944</div> <div>comercio y industria di aruba 881</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.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>