

General overview

Corpus	Date	Language
hplt-v3-spa_Latn	9/24/2025	Spanish

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
725,580,145	16,311,481,939	9,110,045,395 (55.85 %)	512B	2,737,390,990,249	2.54 TB

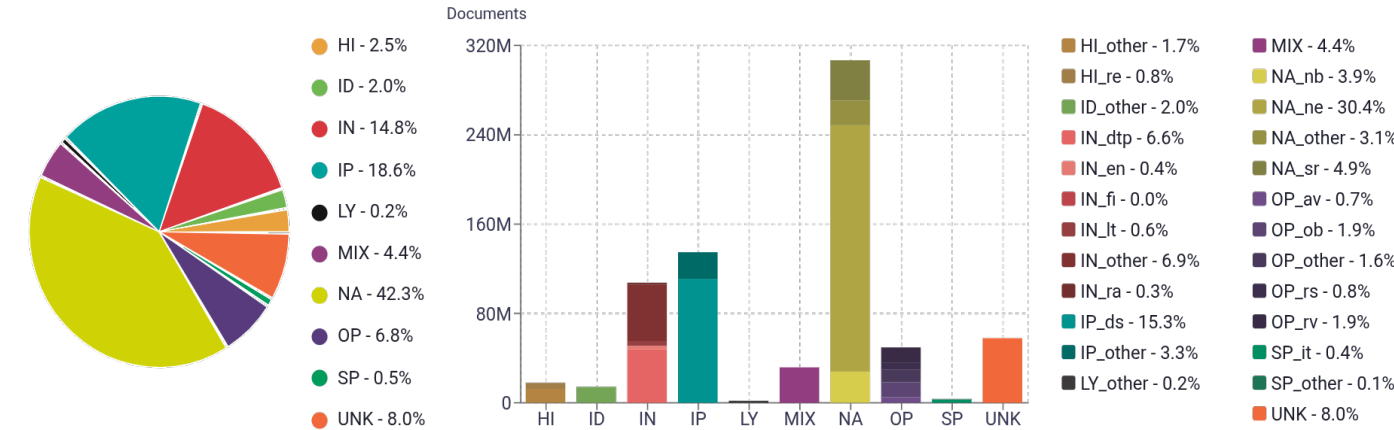
Top 10 domains

Domain	Docs	% of total
blogspot.com	29M	4.05%
wordpress.com	8.5M	1.17%
blogspot.com.es	4.5M	0.62%
buenastareas.com	3.5M	0.48%
blogspot.com.ar	2.4M	0.33%
blogspot.mx	2.2M	0.31%
elpais.com	1.9M	0.26%
web.app	1.4M	0.20%
as.com	1.2M	0.16%
yahoo.com	1.1M	0.16%

Top 10 TLDs

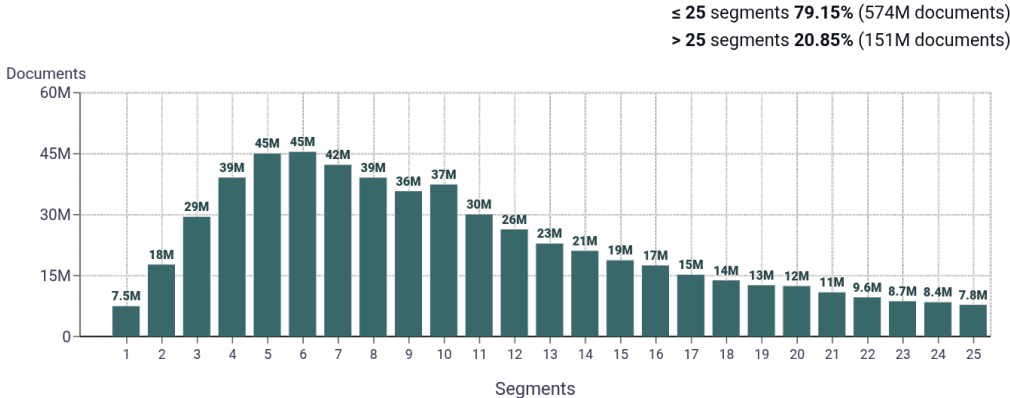
Domain	Docs	% of total
com	371M	51.18%
es	101M	13.90%
com.ar	34M	4.67%
org	27M	3.73%
net	23M	3.18%
com.mx	20M	2.81%
cl	20M	2.70%
mx	16M	2.26%
pe	6.8M	0.93%
info	5.9M	0.82%

Register labels

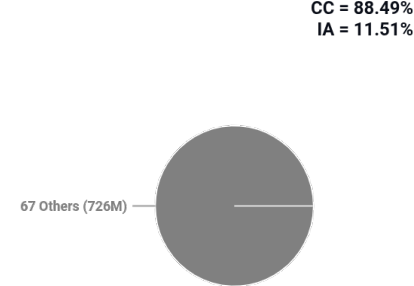


MT:3.8% | 28M Documents

Documents size (in segments) ⓘ

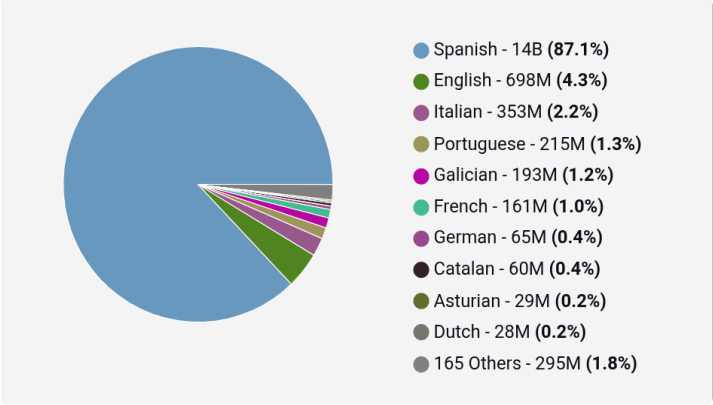


Document collections

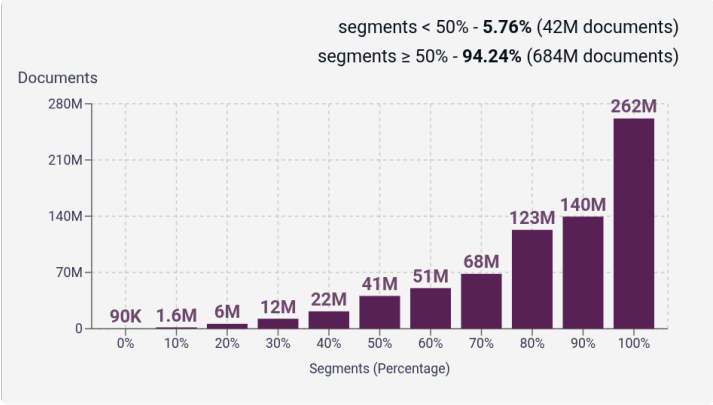


Language Distribution

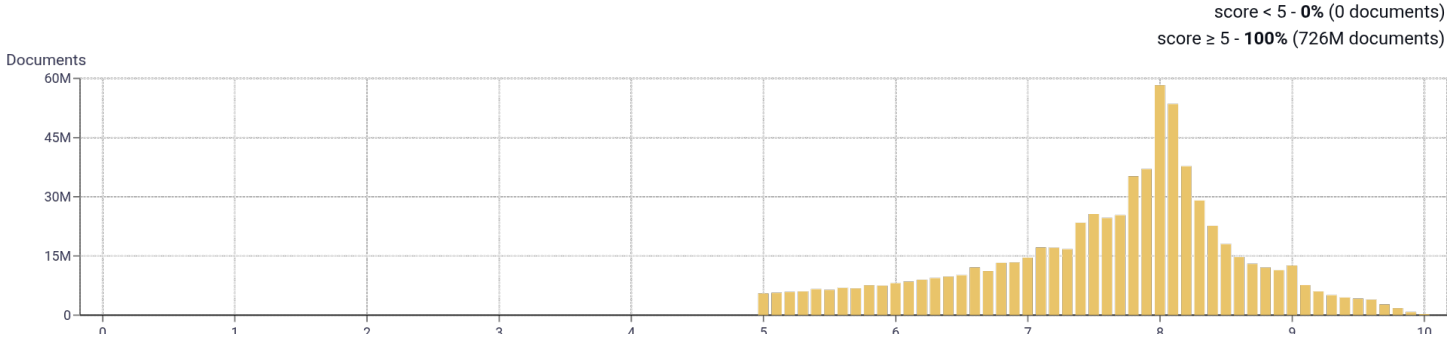
Number of segments in the Spanish corpus



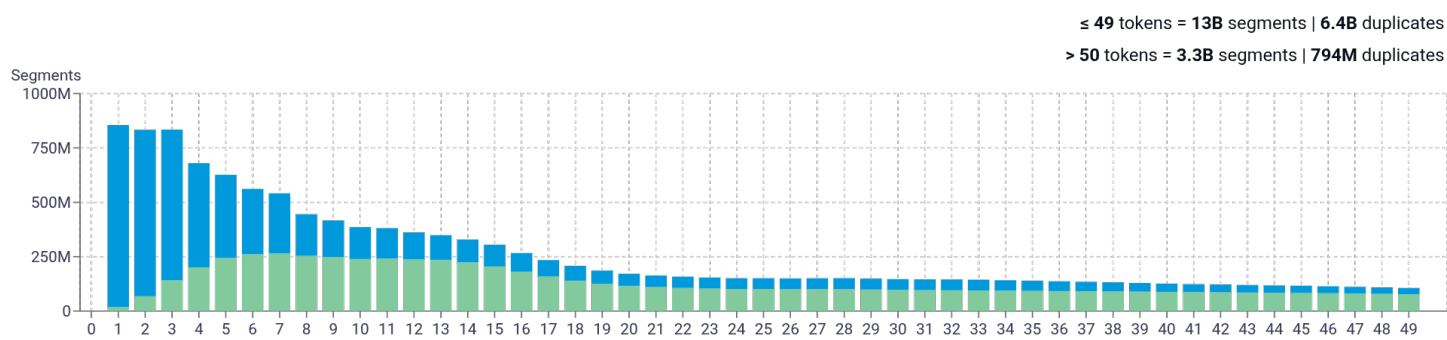
Percentage of segments in Spanish inside documents



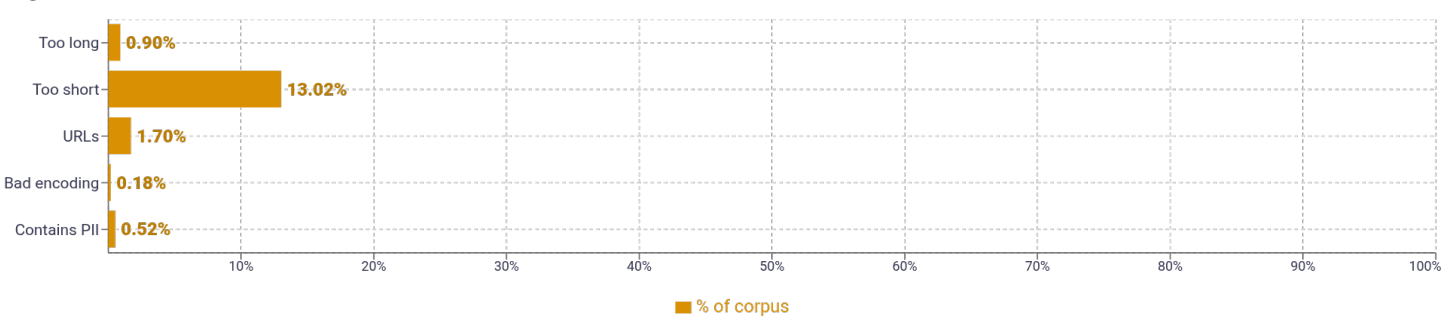
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	<div>ser 643,143,061</div> <div>prostitutas 548,934,410</div> <div>años 501,140,618</div> <div>dos 455,327,367</div> <div>así 428,749,066</div>	
2	<div>redes sociales 51,660,916</div> <div>alta calidad 45,661,142</div> <div>sitio web 43,361,566</div> <div>prostitutas prostitutas 42,781,955</div> <div>primera vez 38,493,484</div>	
3	<div>molino de bolas 32,120,939</div> <div>tener en cuenta 28,348,056</div> <div>fin de semana 27,532,359</div> <div>trituradora de piedra 27,041,319</div> <div>publicar un comentario 24,649,442</div>	
4	<div>opinión de un comprador 10,623,690</div> <div>comentario en la entrada 5,530,382</div> <div>importante tener en cuenta 5,093,737</div> <div>presidente de la república 4,844,115</div> <div>provincia de buenos aires 4,752,698</div>	
5	<div>organización mundial de la salud 3,046,120</div> <div>mejor forma de comprar online 2,470,833</div> <div>servicio de atención al cliente 2,461,309</div> <div>américa latina y el caribe 2,358,188</div> <div>mejorar la calidad de vida 2,131,368</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>

Register labels

Name	Abbr.	Name	Abbr.	Name	Abbr.
Machine-translated	MT	How-to or instructions	HI	Description of a thing or person	ntp
Lyrical	LY	Recipe	re	FAQ	fi
Spoken	SP	Informational persuasion	IP	Legal terms & conditions	lt
Interview	it	Description with intent to sell	ds	Opinion	OP
Interactive discussion	ID	News & opinion blog or editorial	ed	Review	rv
Narrative	NA	Informational description	IN	Opinion blog	ob
News report	ne	Encyclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				