

General overview

Corpus	Date	Language
hplt-v3-srd_Latn	9/18/2025	Sardinian

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
66,660	791,941	602,167 (76.04 %)	35M	174,508,820	170.06 MB

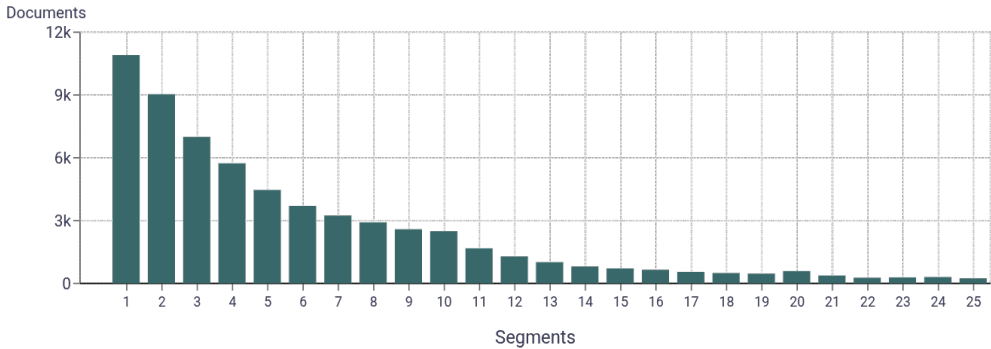
Top 10 domains

Domain	Docs	% of total
sardegna.cultura.it	16K	24.54%
nor-web.eu	5.5K	8.22%
wikipedia.org	4.8K	7.15%
ilminuto.info	4.3K	6.47%
lacanas.it	2.1K	3.13%
salimbassarda.net	1.5K	2.32%
blogspot.com	1.2K	1.81%
wordpress.com	1.1K	1.65%
sagazeta.info	1K	1.54%
academiadesusar...	961	1.44%

Top 10 TLDs

Domain	Docs	% of total
it	33K	49.34%
org	7.2K	10.74%
com	6.6K	9.84%
eu	6.5K	9.79%
info	5.8K	8.72%
net	3.6K	5.36%
co	1.3K	1.93%
or.it	647	0.97%
xyz	277	0.42%
de	108	0.16%

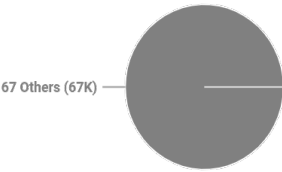
Documents size (in segments) ⓘ



≤ 25 segments 92.71% (62K documents)
> 25 segments 7.29% (4.9K documents)

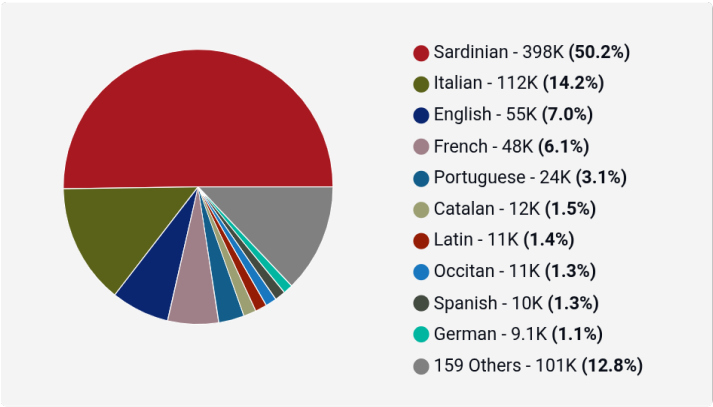
Document collections

CC = 96.33%
IA = 3.67%

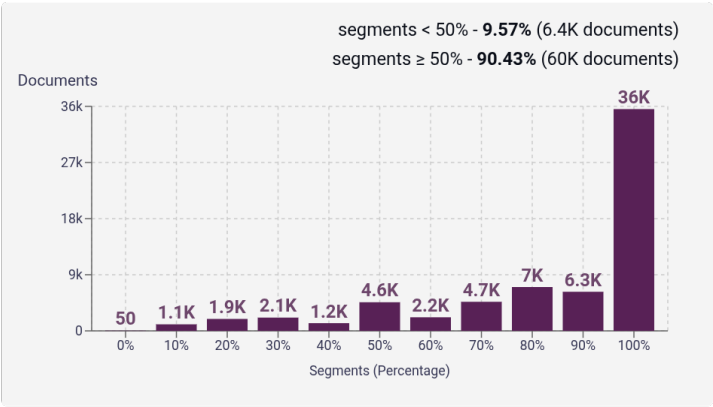


Language Distribution

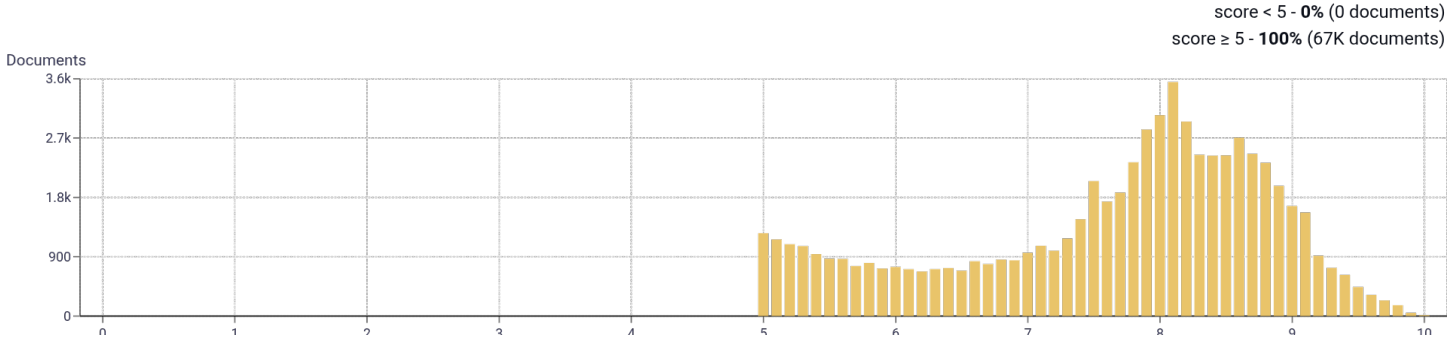
Number of segments in the Sardinian corpus



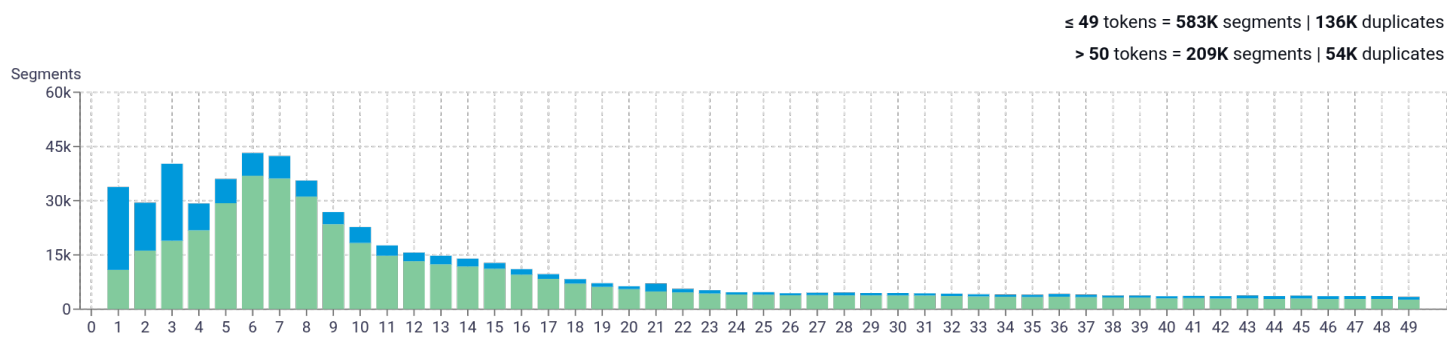
Percentage of segments in Sardinian inside documents



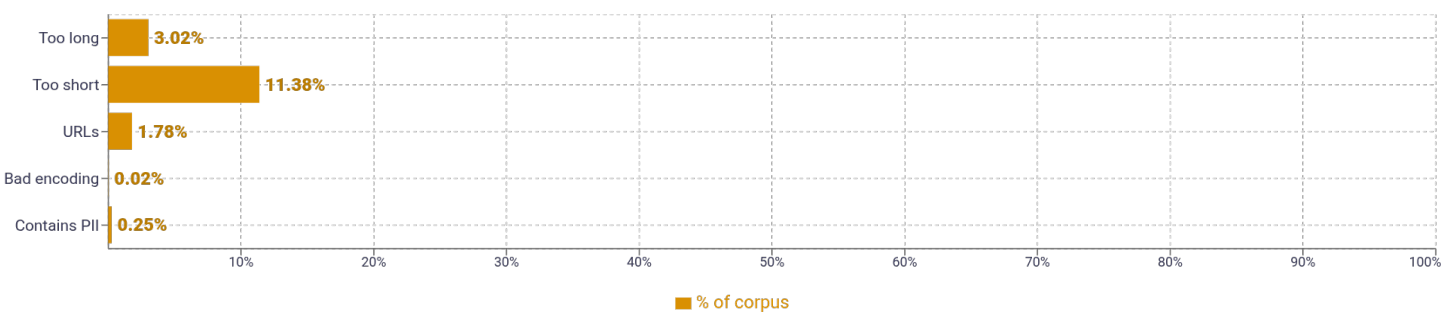
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	sa 666,682s 485,915est 284,938si 232,719sos 197,588	
2	traduzione frantzesu 41,342cun sa 21,194ètimu srd 18,232ingresu to 16,839sa limba 16,645	
3	sinònimos e contràrios 40,029synonyms e antonyms 11,934sinonimi e contrari 11,482sa limba sarda 5,954còdighe de origine 5,732	
4	maneras de nàrrere csns 4,916realizzata col contributo della 2,777contributo della regione sardegna 2,777col contributo della regione 2,777attività realizzata col contributo 2,776	
5	modifica su còdighe de origine 5,675realizzata col contributo della regione 2,777col contributo della regione sardegna 2,777attività realizzata col contributo della 2,776acadèmia de su sardu onlus 1,243	

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 $\frac{\text{number of types (uniques)}}{\text{number of tokens}}$, 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	Enciclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				