

General overview

Corpus	Date	Language
hplt-v3-ita_Latn	9/18/2025	Italian (it)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
362,985,922	7,535,521,218	4,066,403,281 (53.96 %)	238B	1,294,420,645,275	1.19 TB

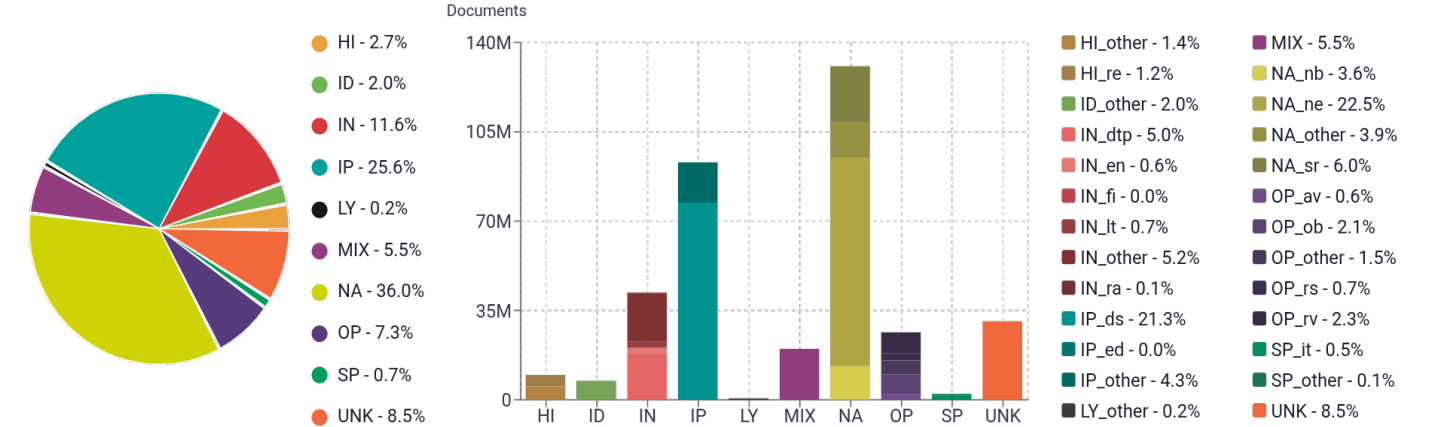
Top 10 domains

Domain	Docs	% of total
blogspot.com	6.8M	1.88%
wordpress.com	3.9M	1.09%
blogspot.it	1.8M	0.51%
kijiji.it	1.5M	0.42%
altervista.org	1.3M	0.36%
docplayer.it	1M	0.28%
repubblica.it	971K	0.27%
corriere.it	969K	0.27%
tripadvisor.it	856K	0.24%
ilsole24ore.com	817K	0.23%

Top 10 TLDs

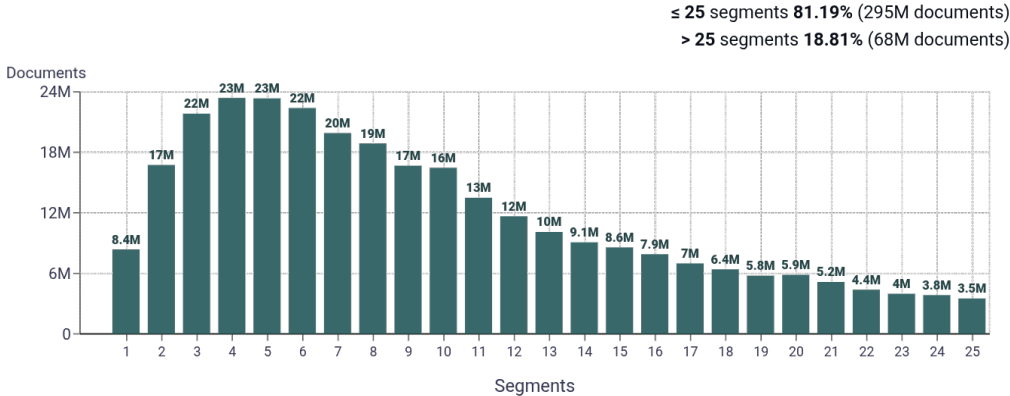
Domain	Docs	% of total
it	207M	57.06%
com	94M	26.01%
org	13M	3.46%
net	12M	3.31%
eu	6.3M	1.74%
info	4.4M	1.21%
ch	3.1M	0.87%
tv	1.7M	0.47%
es	1.2M	0.33%
ru	957K	0.26%

Register labels

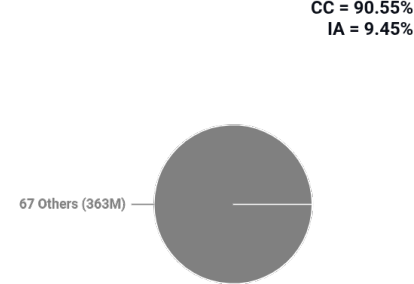


MT:5.3% | 19M Documents

Documents size (in segments) ⓘ

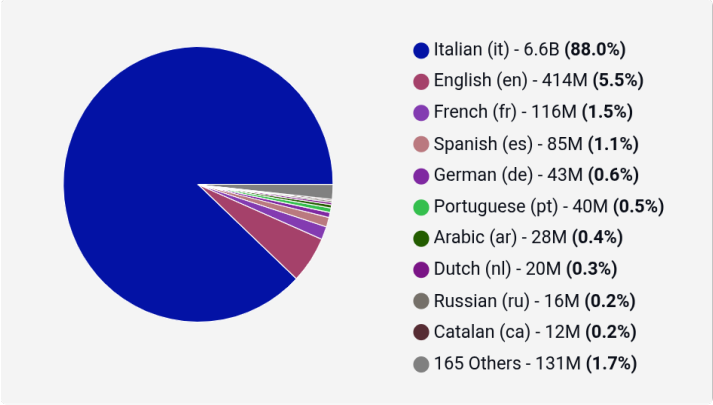


Document collections

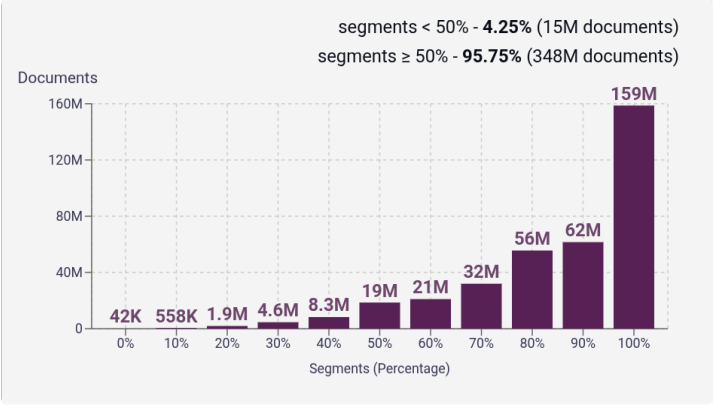


Language Distribution

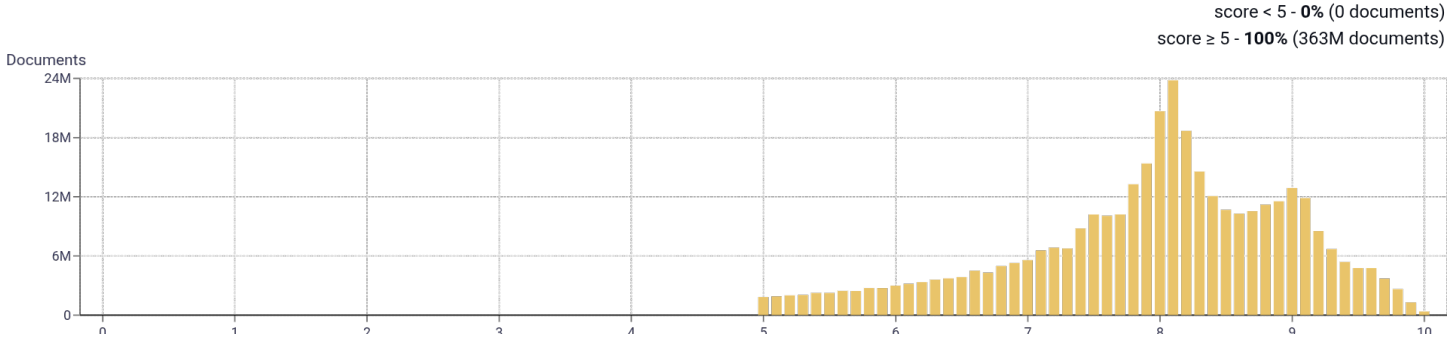
Number of segments in the Italian (it) corpus



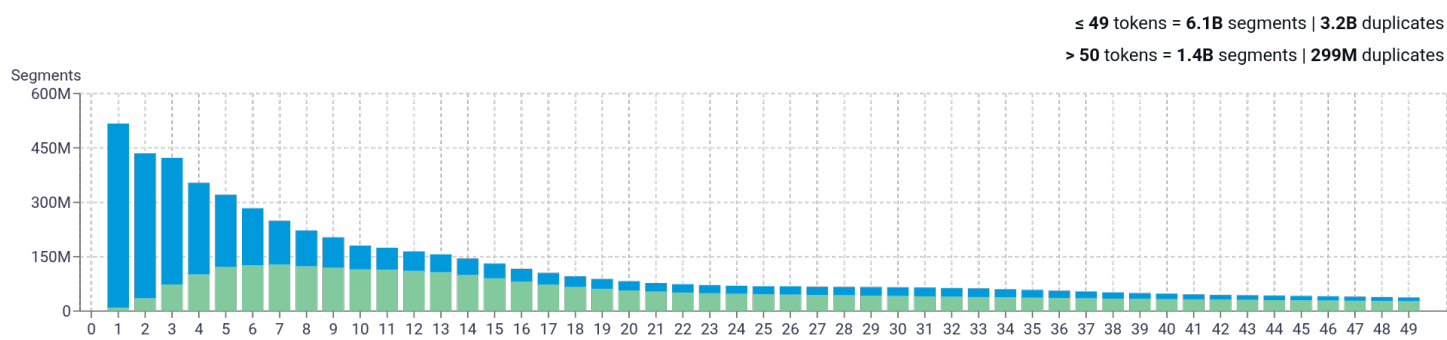
Percentage of segments in Italian (it) inside documents



Distribution of documents by document score

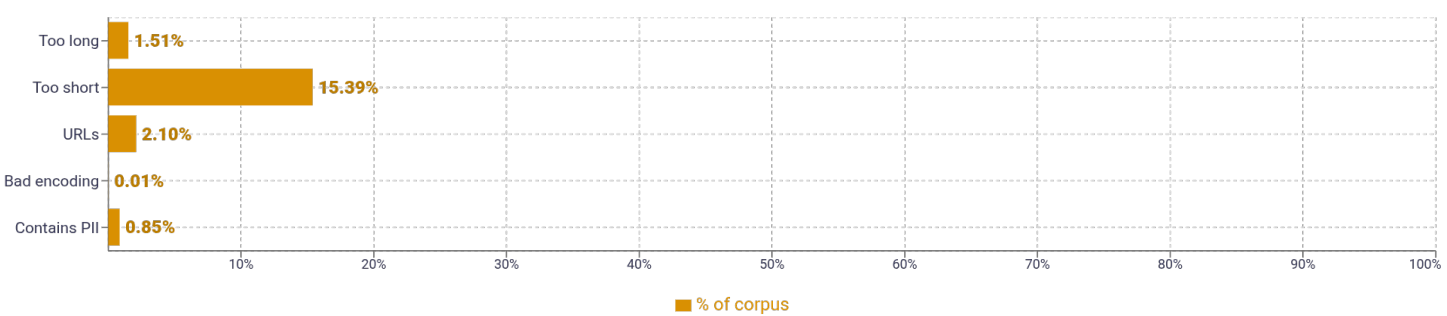


Segment length distribution by token



≤ 49 tokens = 6.1B segments | 3.2B duplicates
> 50 tokens = 1.4B segments | 299M duplicates

Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	<div>essere 376,842,948</div> <div>solo 277,436,536</div> <div>anni 263,607,694</div> <div>due 247,680,609</div> <div>prima 242,934,848</div>	
2	<div>può essere 59,512,639</div> <div>possono essere 31,827,562</div> <div>dopo aver 25,119,554</div> <div>deve essere 24,144,956</div> <div>video porno 22,964,172</div>	
3	<div>continua a leggere 19,516,688</div> <div>punto di vista 16,775,493</div> <div>donna cerca uomo 12,954,258</div> <div>milioni di euro 12,190,188</div> <div>perdita di peso 11,118,312</div>	
4	<div>maggior parte dei casi 1,861,710</div> <div>soggiorno con angolo cottura 1,711,764</div> <div>donne che cercano uomini 1,588,728</div> <div>due camere da letto 1,537,450</div> <div>trattamento dei dati personali 1,330,810</div>	
5	<div>ancora non ci sono recensioni 1,682,214</div> <div>donne in cerca di uomini 1,610,348</div> <div>slots free spins no deposit 1,296,405</div> <div>vita di tutti i giorni 1,113,824</div> <div>stimato in base alle ricerche 1,061,137</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 $\frac{\text{number of types (uniques)}}{\text{number of tokens}}$, 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	dtp
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				