

General overview

Corpus	Date	Language
hplt-v3-cym_Latn	9/18/2025	Welsh

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
1,081,979	21,098,406	14,720,494 (69.77 %)	651M	3,171,744,689	2.99 GB

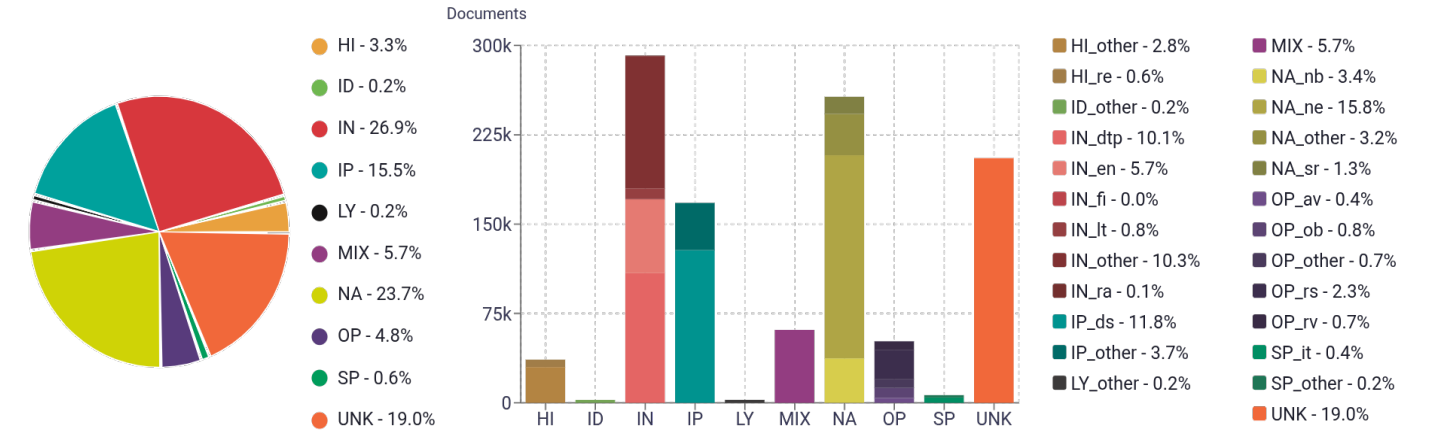
Top 10 domains

Domain	Docs	% of total
wikipedia.org	59K	5.47%
llyw.cymru	27K	2.53%
bbc.co.uk	25K	2.34%
library.wales	20K	1.81%
bbc.com	18K	1.70%
360.cymru	18K	1.68%
testunau.org	15K	1.42%
s4c.cymru	14K	1.30%
cardiff.ac.uk	13K	1.22%
golwg360.cymru	13K	1.19%

Top 10 TLDs

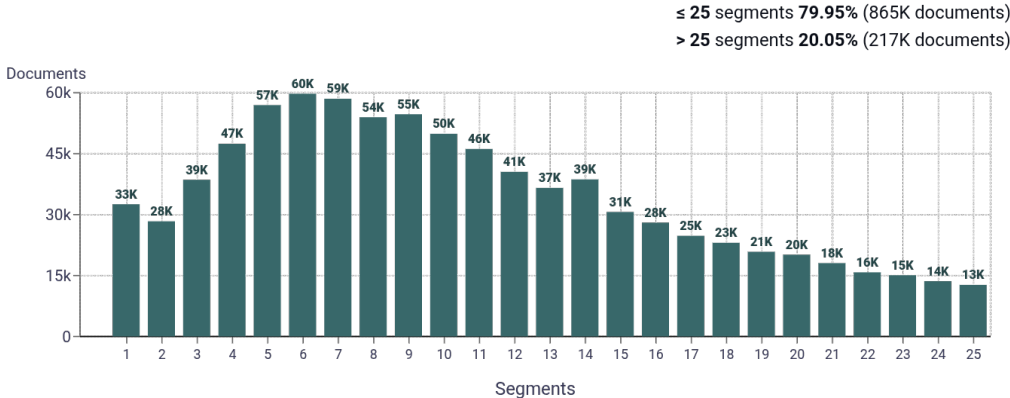
Domain	Docs	% of total
com	316K	29.25%
cymru	185K	17.10%
org	149K	13.80%
co.uk	92K	8.55%
ac.uk	85K	7.87%
gov.uk	65K	6.05%
org.uk	65K	5.99%
wales	46K	4.28%
net	15K	1.40%
zone	6.4K	0.59%

Register labels

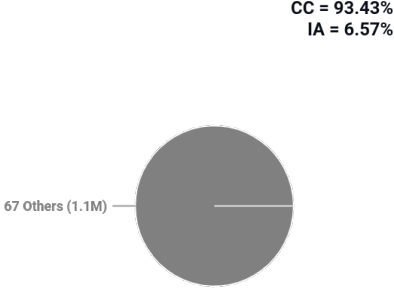


MT:15.4% | 167K Documents

Documents size (in segments)

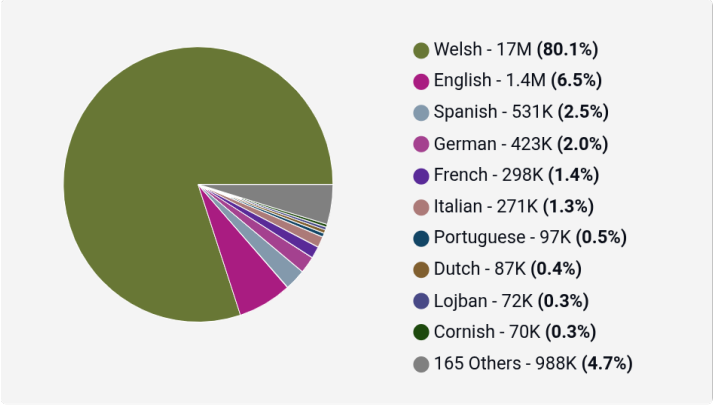


Document collections

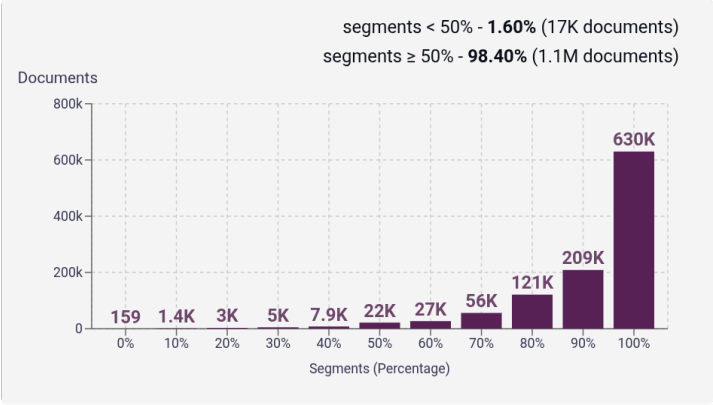


Language Distribution

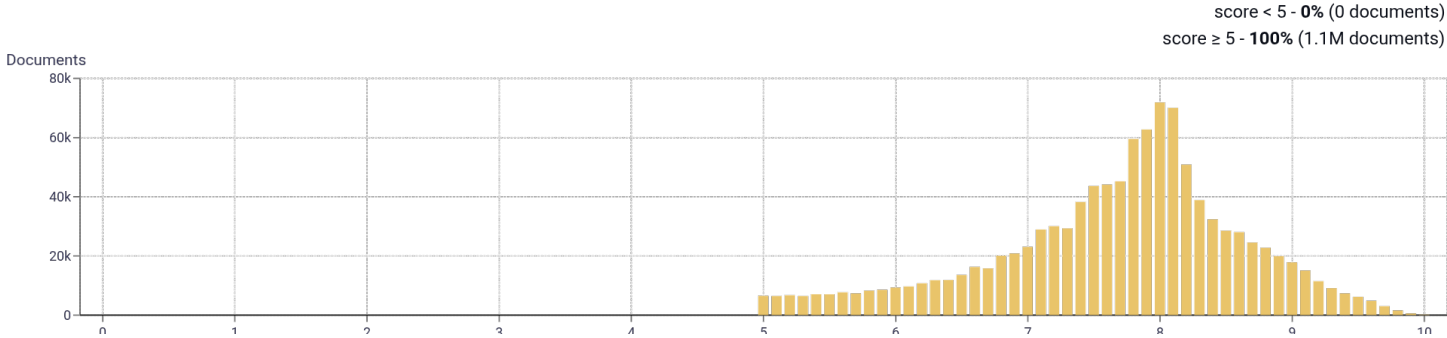
Number of segments in the Welsh corpus



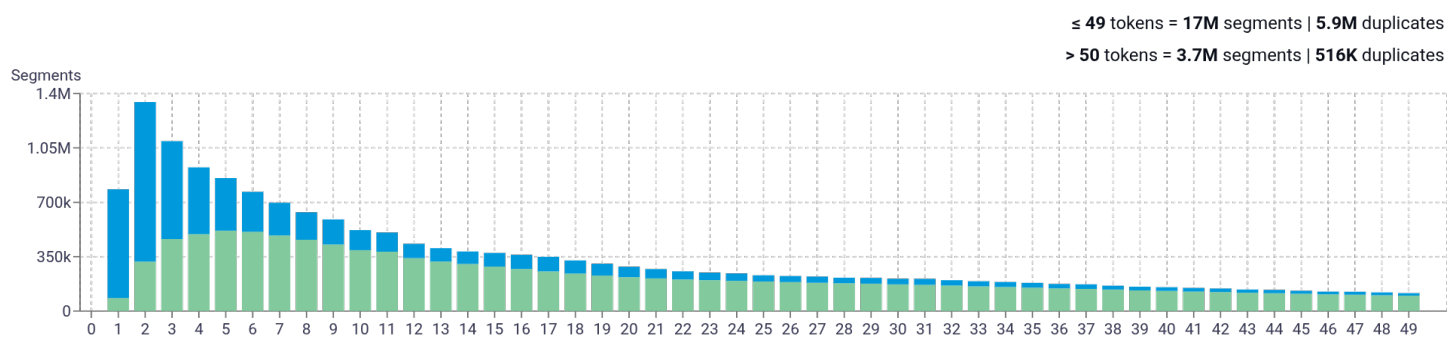
Percentage of segments in Welsh inside documents



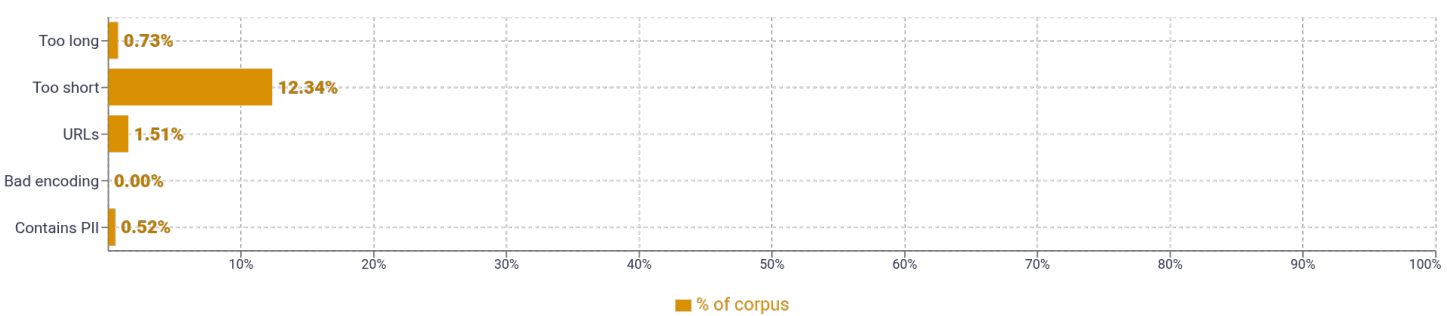
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS					
1	cymru 1,259,774	newydd 938,436	gwaith 677,413	amser 656,153	cynnwys 617,162	
2	llywodraeth cymru 124,047		golygu cod 89,410		awyr agored 58,998	
3	cod y dudalen 48,176		rhagor o wybodaeth 41,413		cynnal a chadw 25,907	
4	golygu cod y dudalen 48,172		poblogaidd y flwyddyn honno 24,282		ganwyd y cyfarwyddwr ffilm 21,854	
	newid yn yr hinsawdd 16,561				actorion yn y ffilm 16,970	
5	ffilmiau gan gynnwys y canlynol 24,303		ffilm fwyaf poblogaidd y flwyddyn 24,289		sgwennwyd y sgript yn wreiddiol 21,565	
	ellir dweud fod y gwaith 8,453		gwaith wedi cyrraedd enwogrwydd byd-eang 8,451			

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				