

General overview

Corpus	Date	Language
hplt-v3-ewe_Latn	9/17/2025	Ewe

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
7,137	218,033	196,843 (90.28 %)	9.5M	39,738,732	41.72 MB

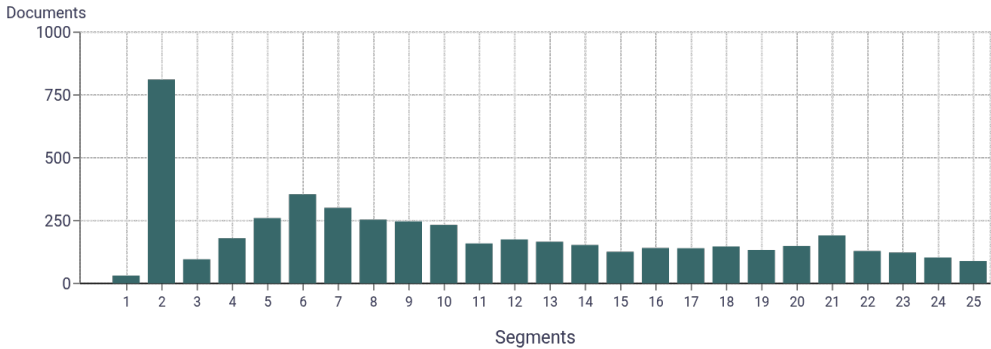
Top 10 domains

Domain	Docs	% of total
jw.org	5K	70.00%
lyfta.app	910	12.75%
eu5.org	227	3.18%
wikipedia.org	125	1.75%
ebible.org	77	1.08%
degbea.org	60	0.84%
zhitov.ru	57	0.80%
ewelyrics.com	56	0.78%
biblearc.com	42	0.59%
bibles.org	35	0.49%

Top 10 TLDs

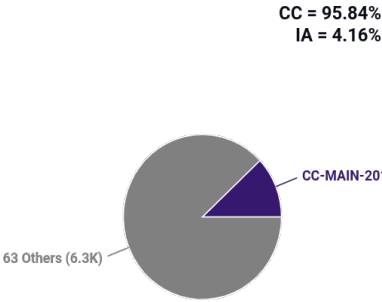
Domain	Docs	% of total
org	5.6K	78.51%
app	910	12.75%
com	400	5.60%
ru	60	0.84%
net	33	0.46%
info	33	0.46%
co.uk	26	0.36%
fr	14	0.20%
is	9	0.13%
news	8	0.11%

Documents size (in segments) ⓘ



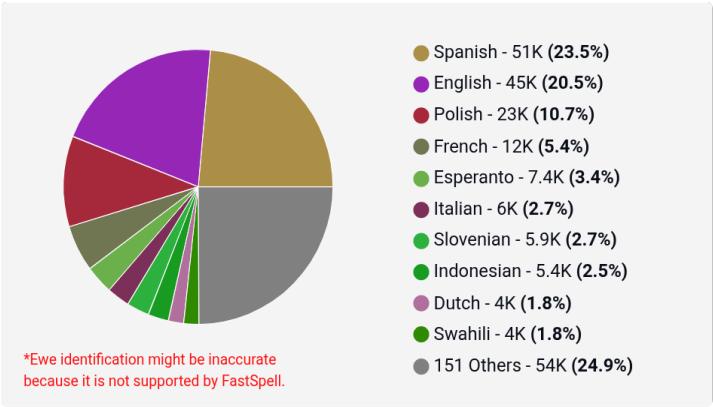
≤ 25 segments **68.56%** (4.9K documents)
> 25 segments **31.44%** (2.2K documents)

Document collections

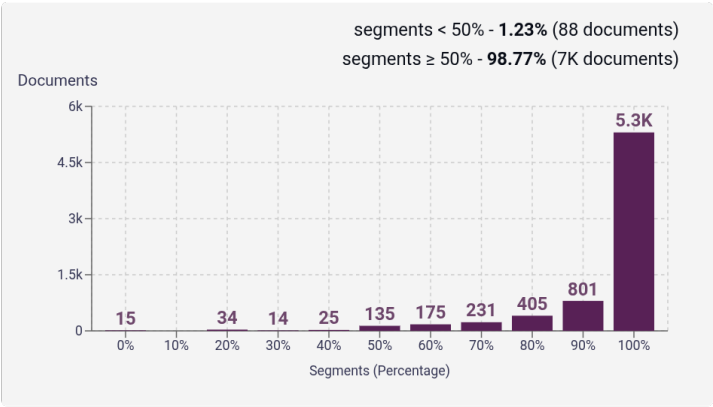


Language Distribution

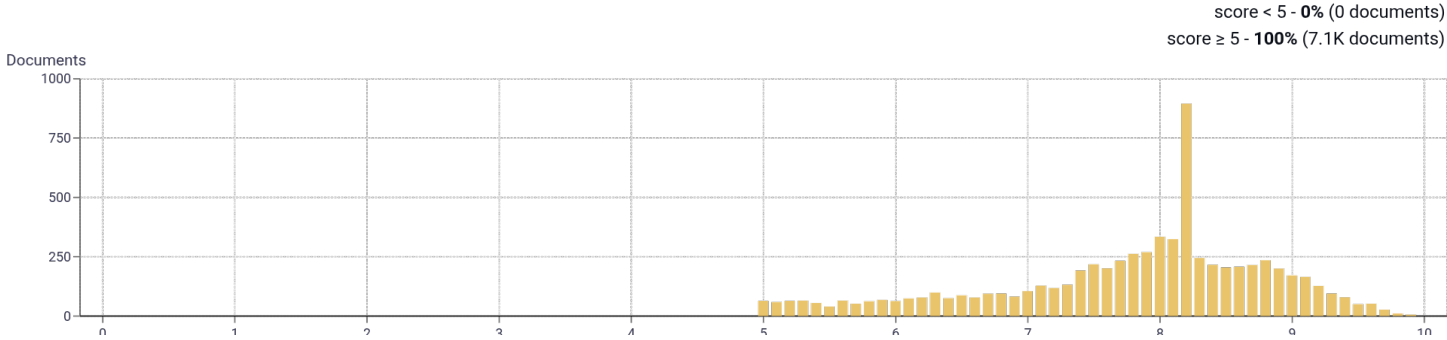
Number of segments in the Ewe corpus



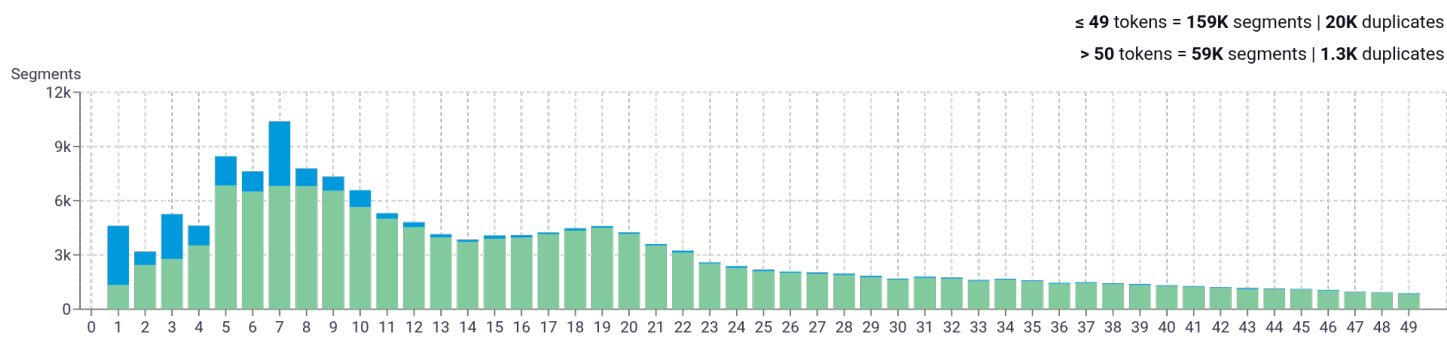
Percentage of segments in Ewe inside documents



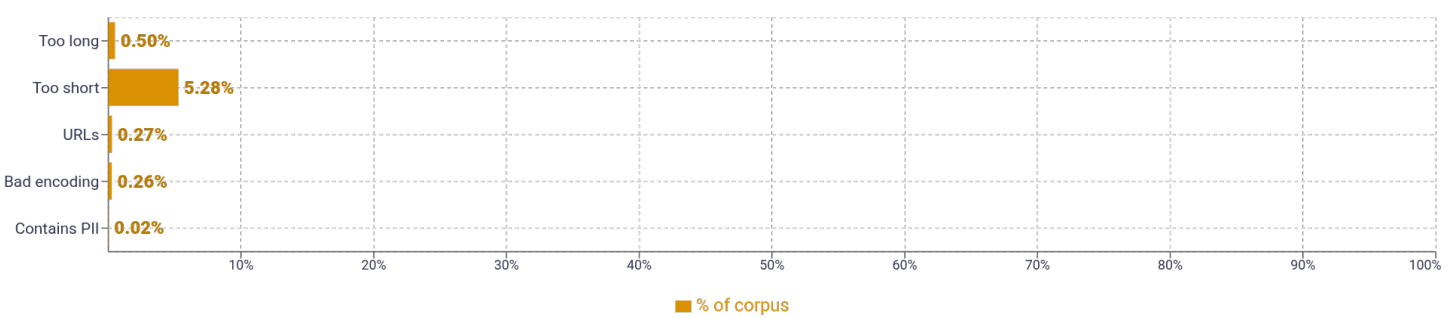
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	yehowa 42,019efe 38,491biblia 20,045wofe 18,500ate 16,477	
2	yehowa ðasefowo 3,297in the 2,511akpa gãṭ 1,583apostolo paulo 1,401nenema kee 1,253	
3	ate ɲu anye 768miate ɲu asɔ 733ate ɲu akpɔ 695xexe yeye gɔmedede 575ate ɲu adzɔ 563	
4	nɔnɔmetata si le axa 2,197nɔnɔmetata siwo le axa 530akpa si le etame 379alesi wòle le kamedede 352fofo si le dzifo 330	
5	ate ɲu akpe de ɲuwò 502alesi wòle le kamedede yeye 183millised on head tàìendavad harjutused 164alesi wòle le kamedede deɕiaɖe 163nɔnɔmetata tsofe si le axa 152	

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				