

General overview

Corpus	Date	Language
hplt-v3-sna_Latn	9/18/2025	Shona

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
183,006	3,779,319	3,019,107 (79.88 %)	93M	624,924,975	597.31 MB

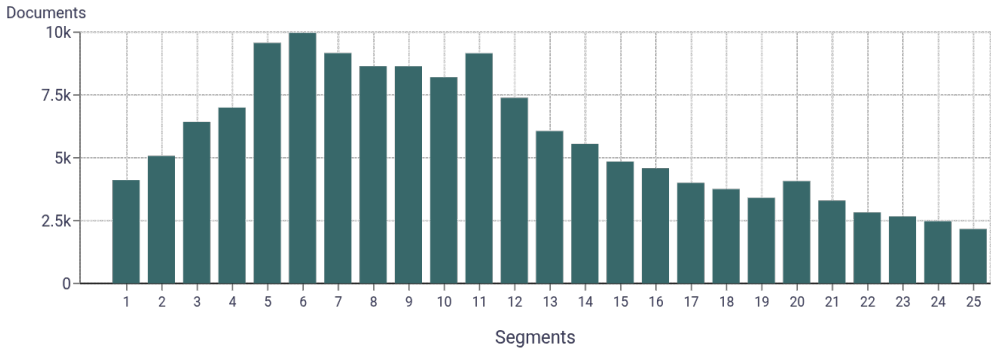
Top 10 domains

Domain	Docs	% of total
voashona.com	17K	9.48%
linuxadictos.com	6.8K	3.74%
jw.org	5.6K	3.03%
martech.zone	5.5K	3.01%
kwayedza.co.zw	5.2K	2.86%
eturbonews.com	4.6K	2.52%
actualidadiphon...	3.9K	2.16%
wikipedia.org	2.5K	1.36%
soydemac.com	2.4K	1.34%
actualidadliter...	1.5K	0.83%

Top 10 TLDs

Domain	Docs	% of total
com	141K	77.16%
org	11K	6.03%
co.zw	6.2K	3.36%
zone	5.5K	3.01%
net	4.4K	2.40%
es	1.3K	0.72%
online	1.3K	0.71%
fr	1.2K	0.67%
ru	1.2K	0.64%
news	862	0.47%

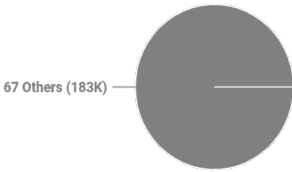
Documents size (in segments) ⓘ



≤ 25 segments 78.21% (143K documents)
> 25 segments 21.79% (40K documents)

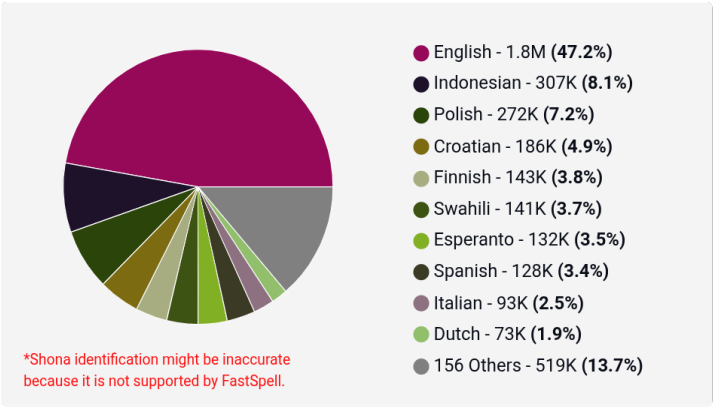
Document collections

CC = 97.66%
IA = 2.34%

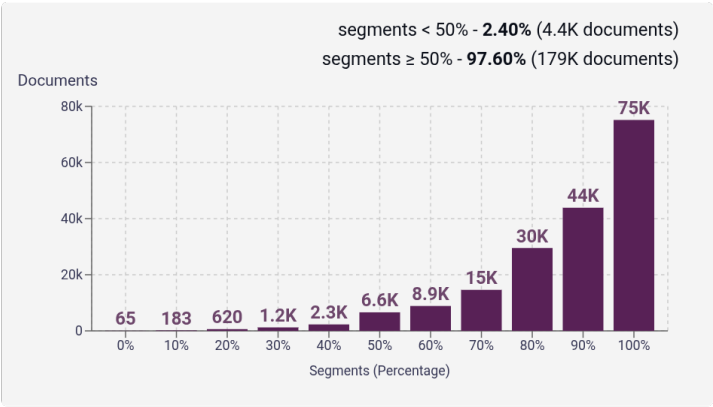


Language Distribution

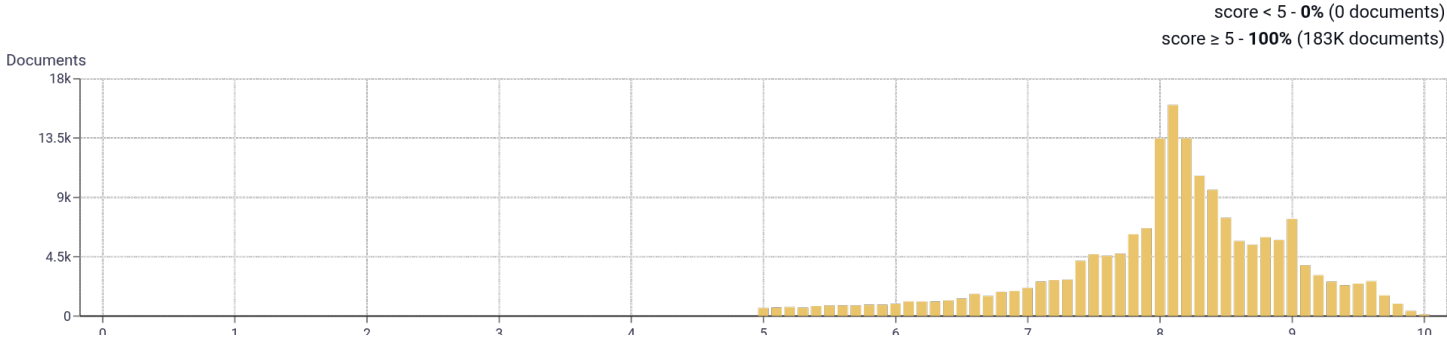
Number of segments in the Shona corpus



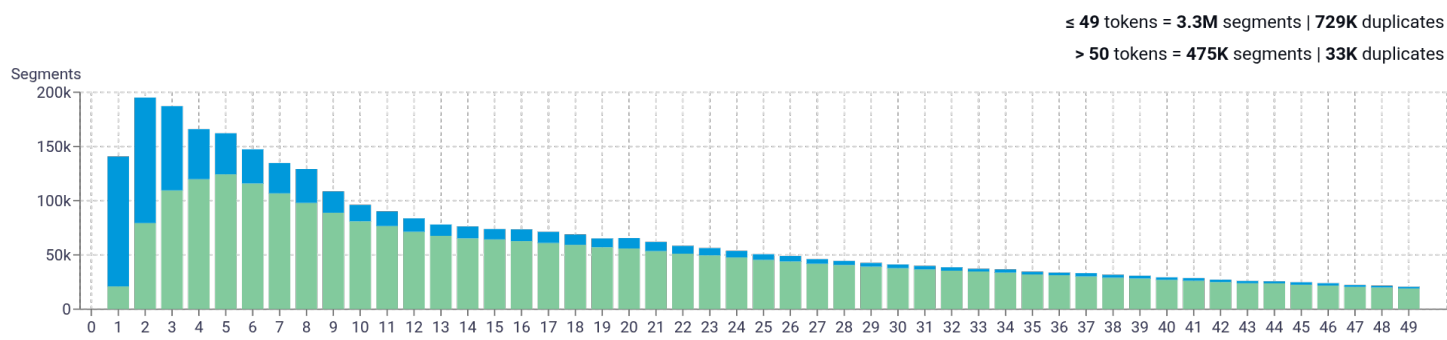
Percentage of segments in Shona inside documents



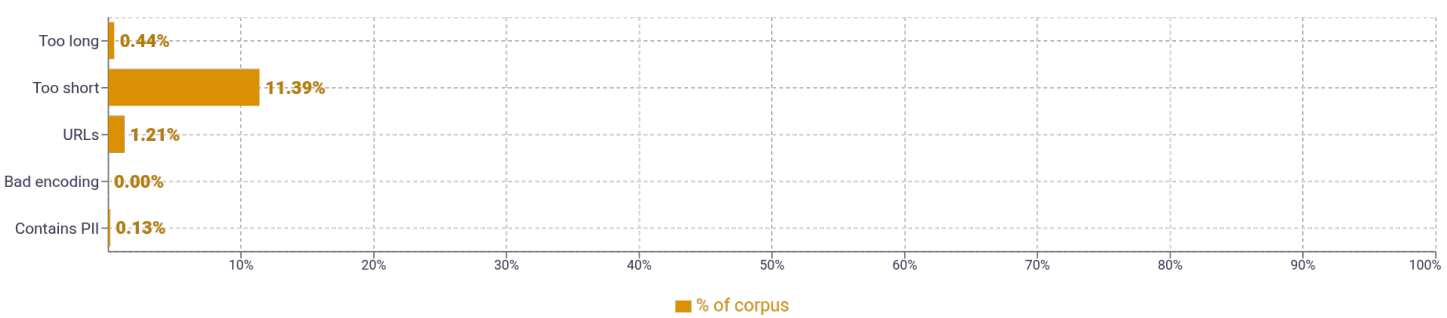
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	kana 923,015iyo 597,966kubva 320,054asi 310,302iri 226,263	📄
2	zviri nyore 28,445makumi maviri 23,767imwe chete 16,409uchinge uchinge 15,799iri nyore 15,561	📄
3	kana iwe uchida 19,504uchinge uchinge uchinge 15,155iva wekutanga kutaura 13,543kana iwe uri 9,554panguva imwe chete 7,624	📄
4	uchinge uchinge uchinge uchinge 14,608plus untold biography facts 2,208munguva pfupi iri kutevera 1,520kana iwe uchida kuziva 1,493kana muchida kupinda muchirongwa 1,368	📄
5	uchinge uchinge uchinge uchinge uchinge 14,108kana muchida kupinda muchirongwa ichi 1,366shanduro yenyika itsva yemagwaro matsvene 1,138munokwanisawo kunzwa studio7 na6 am 779childhood story plus untold biography 726	📄

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				