

General overview

Corpus	Date	Language
hplt-v3-zul_Latn	9/18/2025	Zulu

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
336,440	8,015,560	5,321,756 (66.39 %)	168M	1,113,609,767	1.04 GB

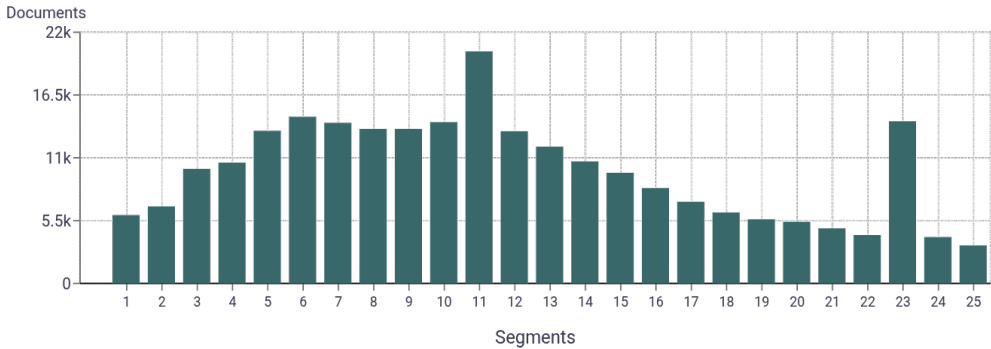
Top 10 domains

Domain	Docs	% of total
airbnb.com	55K	16.29%
voandebele.com	18K	5.25%
news24.com	14K	4.21%
scrolla.africa	10K	3.04%
unanseas.com	9K	2.68%
bayedenews.com	7.3K	2.18%
jw.org	6.4K	1.91%
isolezwe.co.za	6.2K	1.85%
martech.zone	5K	1.50%
impempe.com	4.9K	1.46%

Top 10 TLDs

Domain	Docs	% of total
com	235K	69.73%
co.za	34K	9.98%
org	21K	6.12%
africa	10K	3.05%
net	5.9K	1.74%
zone	5K	1.50%
co.zw	3.1K	0.91%
info	2.2K	0.66%
xyz	1.9K	0.56%
ru	1.3K	0.39%

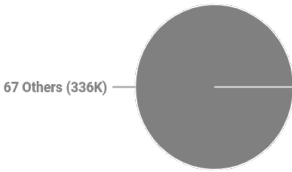
Documents size (in segments) ⓘ



≤ 25 segments 73.1% (246K documents)
> 25 segments 26.9% (91K documents)

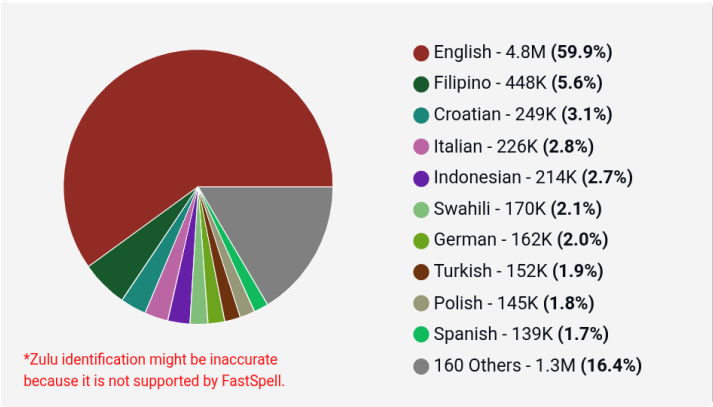
Document collections

CC = 94.85%
IA = 5.15%

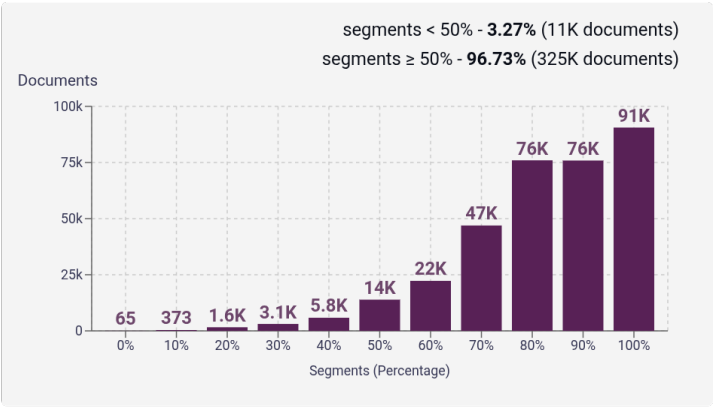


Language Distribution

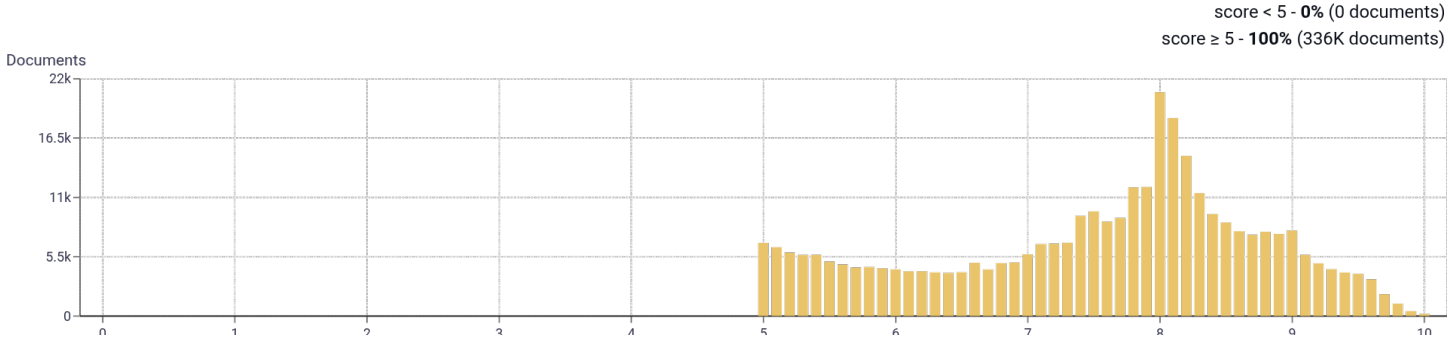
Number of segments in the Zulu corpus



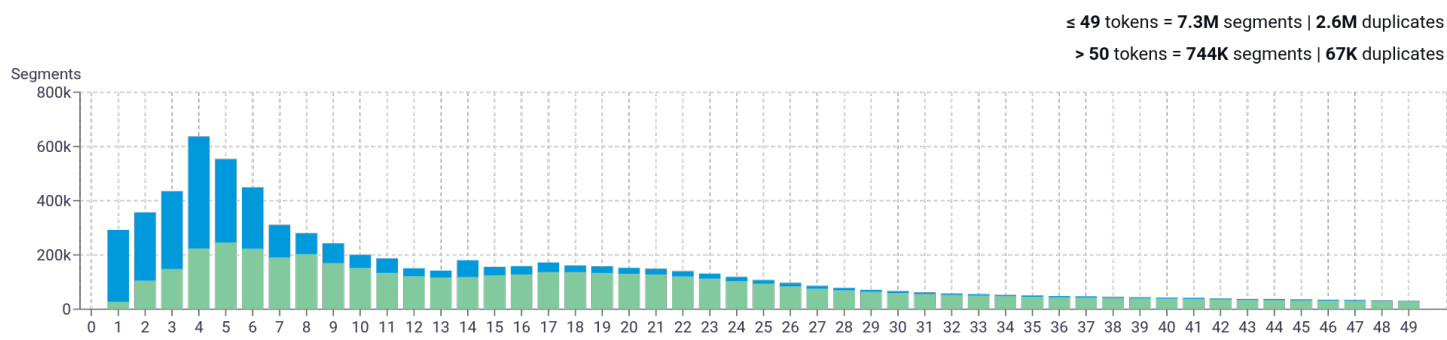
Percentage of segments in Zulu inside documents



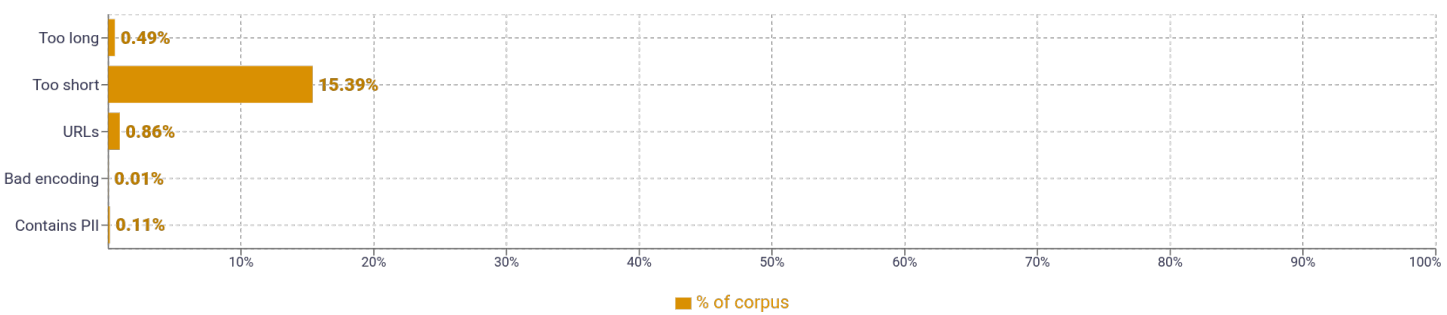
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	<div>i 1,659,953</div> <div>e 1,044,287</div> <div>okungu 523,781</div> <div>ukuphawula 521,149</div> <div>isilinganiso 515,187</div>	
2	<div>ukuphawula okungu 506,543</div> <div>isilinganiso esingu 491,609</div> <div>ikhaya e 193,809</div> <div>ifulethi e 100,895</div> <div>iyunithi yokuqasha 83,414</div>	
3	<div>iyunithi yokuqasha e 81,908</div> <div>izindawo eziqashwayo zeholide 54,403</div> <div>thola futhi ubhukhe 37,049</div> <div>zeholide ezifanela sonke 37,023</div> <div>ubungako bendawo obukwenele 37,023</div>	
4	<div>zeholide ezifanela sonke isitayela 37,023</div> <div>thola ubungako bendawo obukwenele 37,023</div> <div>izindawo eziqashwayo zeholide ezifanela 37,023</div> <div>eziqashwayo zeholide ezifanela sonke 37,023</div> <div>thola futhi ubhukhe izindawo 26,650</div>	
5	<div>izindawo eziqashwayo zeholide ezifanela sonke 37,023</div> <div>eziqashwayo zeholide ezifanela sonke isitayela 37,023</div> <div>thola futhi ubhukhe izindawo zokuhlala 17,769</div> <div>ubhukhe izindawo zokuhlala eziyingqayivele ku 17,397</div> <div>zindawo zokuhlala zinconywa kakhulu ngokwendawo 17,394</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	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				