

General overview

Corpus	Date	Language
hplt-v3-sot_Latn	9/18/2025	Southern Sotho (st)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
152,062	3,618,844	2,830,790 (78.22 %)	128M	601,224,893	576.39 MB

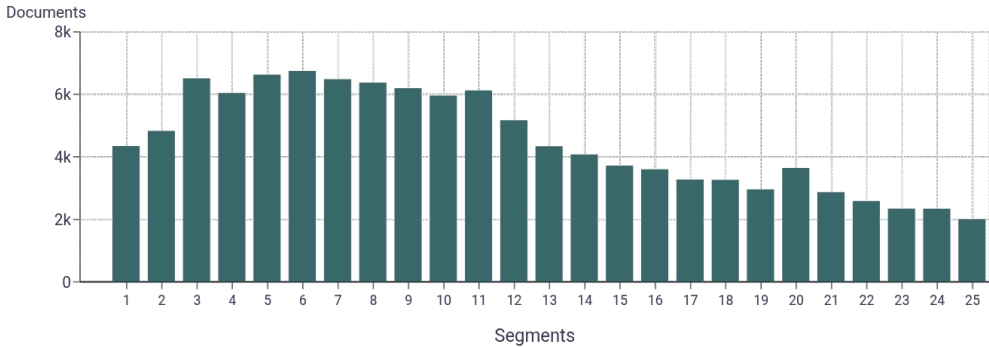
Top 10 domains

Domain	Docs	% of total
eturbonews.com	5.5K	3.62%
martech.zone	5.3K	3.46%
jw.org	3.2K	2.08%
actualidadgadgets.com	2K	1.33%
bitemybun.com	1.9K	1.23%
comme-un-pro.fr	1.4K	0.93%
actualidadviajes.com	1.4K	0.91%
actualidadiphone.com	1.3K	0.84%
desdelinux.net	1.2K	0.78%
hombresconestilo.com	1.2K	0.76%

Top 10 TLDs

Domain	Docs	% of total
com	116K	75.99%
org	7.7K	5.08%
zone	5.3K	3.46%
net	3.8K	2.49%
co.za	2.5K	1.67%
ru	2.3K	1.54%
fr	1.5K	0.96%
news	1.3K	0.84%
info	1.2K	0.81%
co.ls	947	0.62%

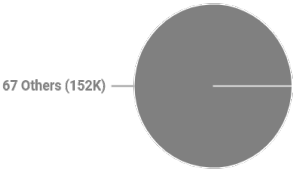
Documents size (in segments) ⓘ



≤ 25 segments **73.94%** (112K documents)
> 25 segments **26.06%** (40K documents)

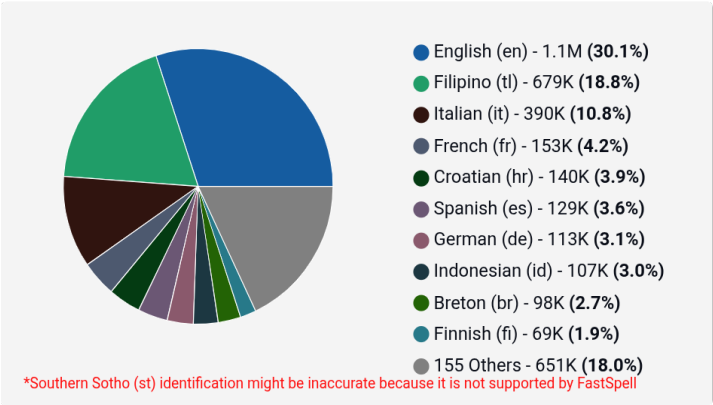
Document collections

CC = **95.63%**
IA = **4.37%**

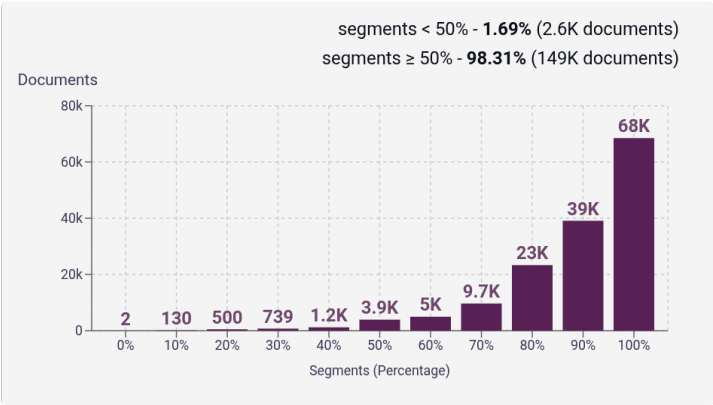


Language Distribution

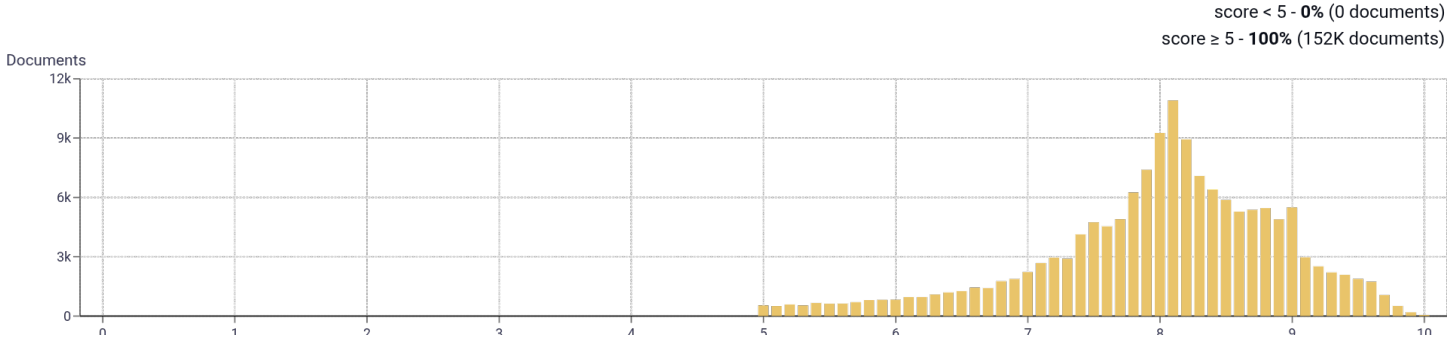
Number of segments in the Southern Sotho (st) corpus



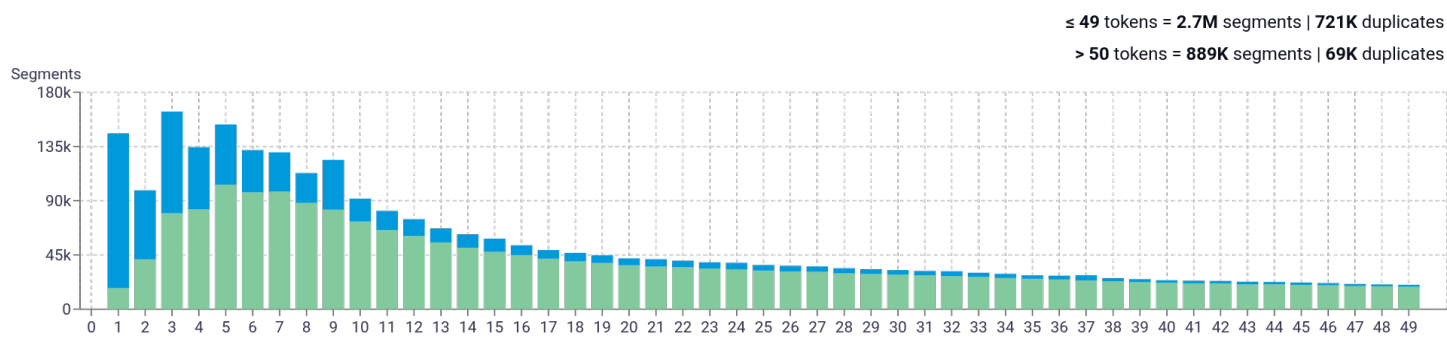
Percentage of segments in Southern Sotho (st) inside documents



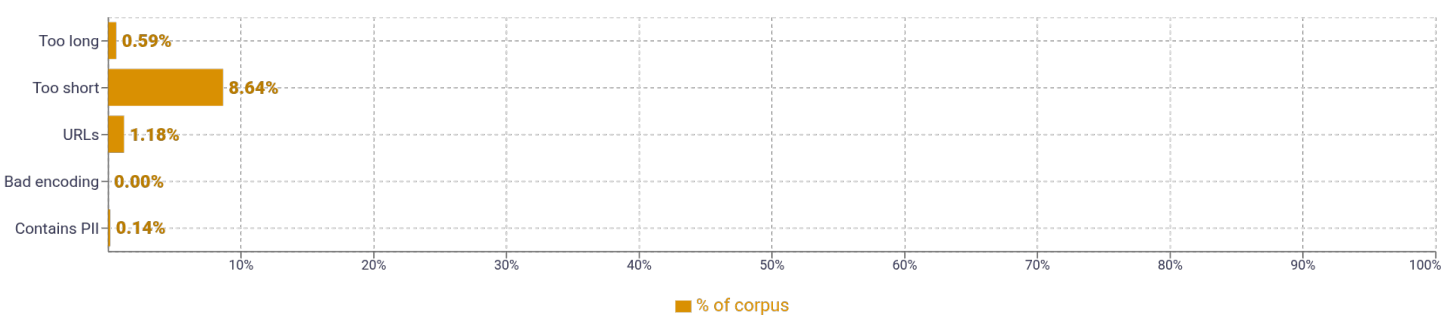
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	u 927,095na 729,609kapa 527,249tla 516,190bakeng 436,607	📄
2	haeba u 87,588boleng bo 52,469u tla 47,777bo botle 42,100hona joale 41,541	📄
3	nako e telele 40,078efe kapa efe 27,448boleng bo holimo 26,495bophelo bo botle 25,470letsatsi le letsatsi 18,805	📄
4	leha ho le joalo 56,230mong le e mong 26,935ntle ka ho fetisisa 19,698u se ke ua 18,470molemo ka ho fetisisa 16,737	📄
5	pele ho fana ka maikutlo 10,018town town town town town 8,714etsa bonnete ba hore u 4,841sebaka sa hau sa marang 4,021na le mefuta e mengata 3,430	📄

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				