

General overview

Corpus	Date	Language
afr_Latn.json.tsv	9/6/2024	Afrikaans (af)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
1,457,165	37,737,319	18,802,427 (49.82 %)	1.2B	5,910,906,748	5.56 GB

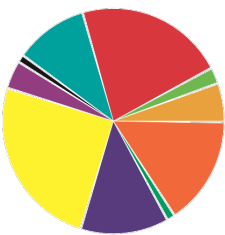
Top 10 domains

Domain	Docs	% of total
wikipedia.org	164K	11.23%
maroelamedia.co.za	68K	4.64%
netwerk24.com	42K	2.87%
landbou.com	35K	2.38%
praag.co.za	33K	2.24%
wordpress.com	32K	2.19%
litnet.co.za	32K	2.17%
sarie.com	28K	1.90%
software.net	20K	1.40%
androware.net	20K	1.39%

Top 10 TLDs

Domain	Docs	% of total
com	529K	36.33%
co.za	451K	30.96%
org	239K	16.38%
net	69K	4.74%
org.za	30K	2.06%
ac.za	28K	1.90%
com.na	19K	1.27%
ca	11K	0.76%
info	9.3K	0.64%
pt	4K	0.28%

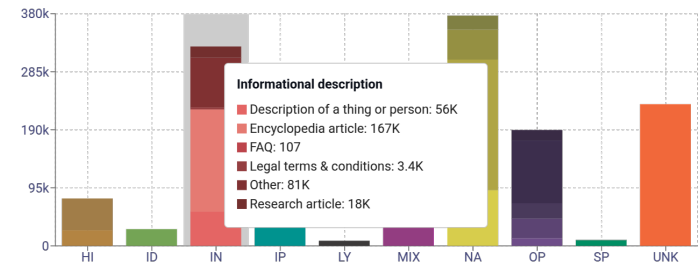
Register labels



- HI - 5.3%
- ID - 1.9%
- IN - 22.4%
- IP - 10.7%
- LY - 0.6%
- MIX - 3.7%
- NA - 25.9%
- OP - 13.0%
- SP - 0.7%
- UNK - 15.9%

MT:12.8% | 186K Documents

Documents

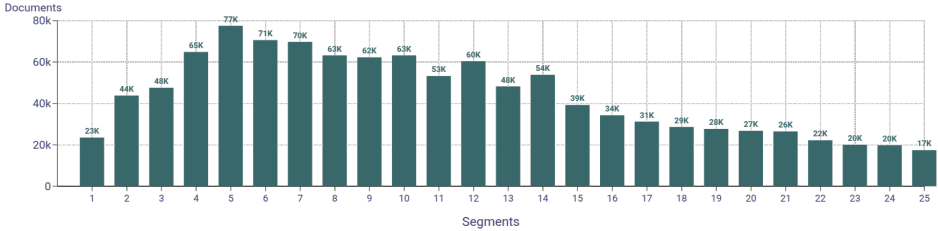


Informational description

- Description of a thing or person: 56K
- Encyclopedia article: 167K
- FAQ: 107
- Legal terms & conditions: 3.4K
- Other: 81K
- Research article: 18K

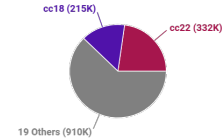
- HI_other - 1.7%
- HI_re - 3.6%
- ID_other - 1.9%
- IN_dtp - 3.8%
- IN_en - 11.5%
- IN_fi - 0.0%
- IN_it - 0.2%
- IN_other - 5.6%
- IN_ra - 1.3%
- IP_ds - 8.6%
- IP_ed - 0.0%
- IP_other - 2.1%
- LY_other - 0.6%
- MIX - 3.7%
- NA_nb - 6.2%
- NA_ne - 14.7%
- NA_other - 3.3%
- NA_sr - 1.6%
- OP_av - 0.9%
- OP_ob - 2.2%
- OP_other - 1.7%
- OP_rs - 7.1%
- OP_rv - 1.2%
- SP_it - 0.6%
- SP_other - 0.1%
- UNK - 15.9%

Documents size (in segments)



<= 25 segments 75.18% (1.1M documents)
> 25 segments 24.82% (362K documents)

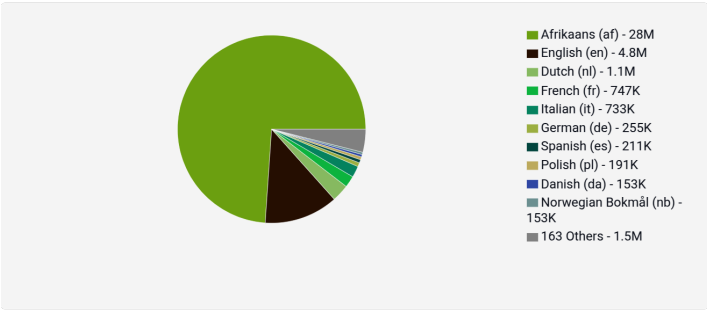
Documents by collection



CC = 59.43%
IA = 40.57%

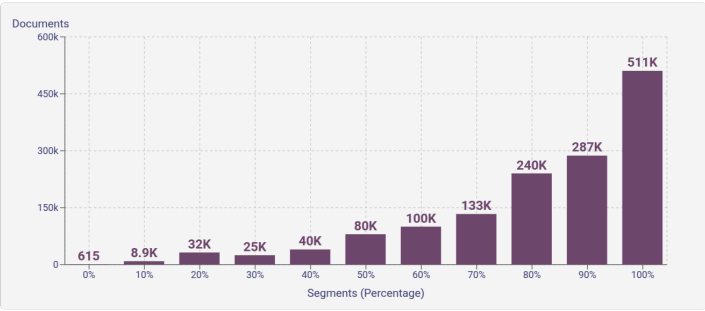
Language Distribution

Number of segments in the Afrikaans (af) corpus

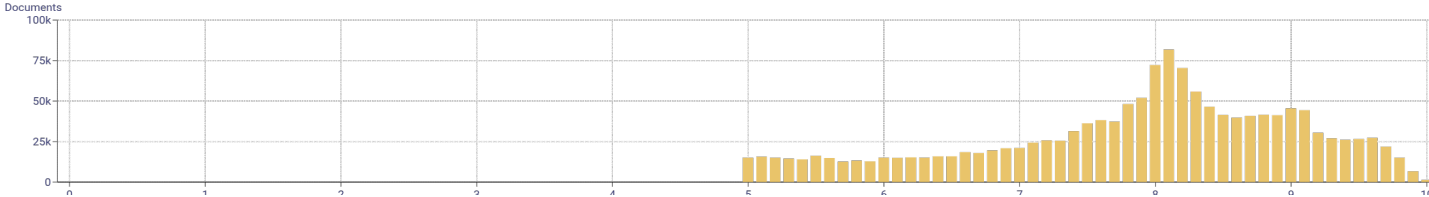


- Afrikaans (af) - 28M
- English (en) - 4.8M
- Dutch (nl) - 1.1M
- French (fr) - 747K
- Italian (it) - 733K
- German (de) - 255K
- Spanish (es) - 211K
- Polish (pl) - 191K
- Danish (da) - 153K
- Norwegian Bokmål (nb) - 153K
- 163 Others - 1.5M

Percentage of segments in Afrikaans (af) inside documents

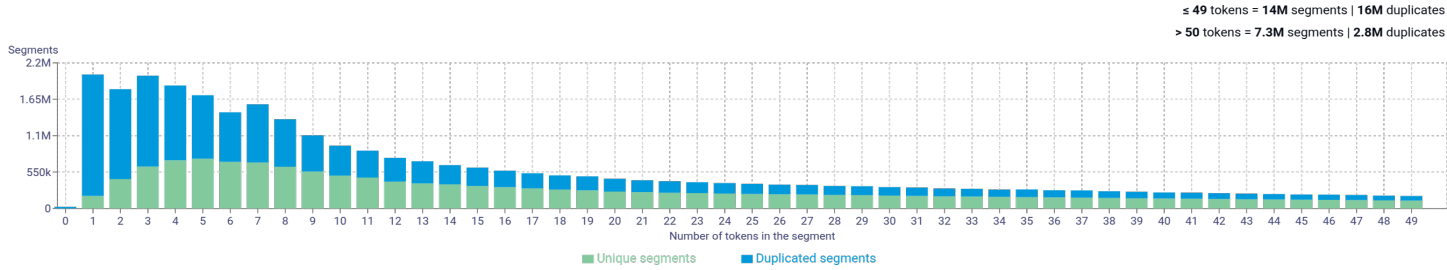


Distribution of documents by document score

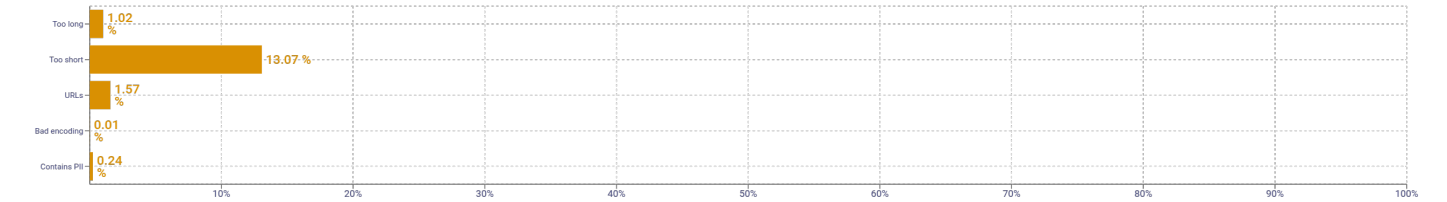


score < 5 - 0% (0 documents)
score >= 5 - 100% (1.5M documents)

Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	word 6406818 of 6184511 deur 3733665 hierdie 3584126 ook 3389903
2	of the 474554 wysig bron 345356 word deur 234204 gebruik word 180407 moet word 173996
3	vanaf die oorspronklike 48240 oor die algemeen 43175 sout en peper 36126 voor te berei 30939 woord van god 30722
4	geargiveer vanaf die oorspronklike 43348 wikimedia commons het meer 24986 commons het meer media 24983 we are searching data 24886 searching data for your 24886
5	wikimedia commons het meer media 24983 we are searching data for 24886 searching data for your request 24886 are searching data for your 24886 speletjie saam met jou beste 24549

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 *number of types (uniques)/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/pablopt6n/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				