

General overview

Corpus	Analytics date	Language
tso_Latn.jsonl.tsv	9/19/2024	Tsonga (ts)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
11,008	221,245	136,723 (61.80 %)	10M	47.58 MB	49,075,139

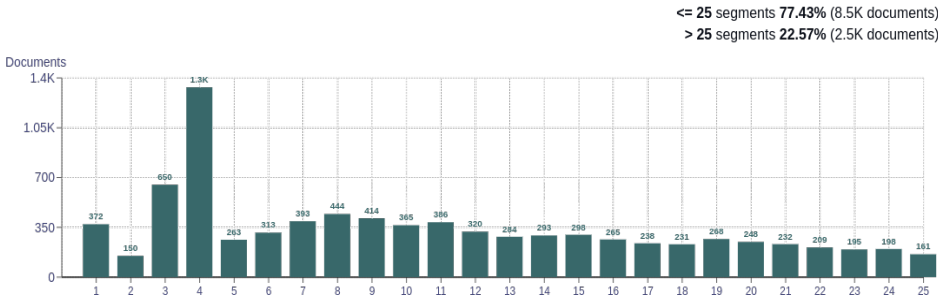
Top 10 domains

Domain	Docs	% of total
jw.org	6.5K	58.92
biblesa.co.za	1K	9.49
wikipedia.org	999	9.08
bible.is	653	5.93
southafrica.co.za	427	3.88
vivmag.co.za	243	2.21
rivoni.org	54	0.49
munghanalonenefm.co.za	45	0.41
myconstitution.co.za	40	0.36
matimunews.co.za	38	0.35

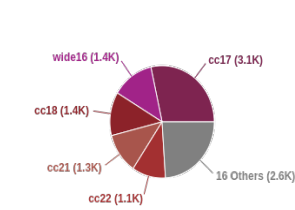
Top 10 TLDs

Domain	Docs	% of total
org	7.8K	70.60
co.za	2K	18.10
is	653	5.93
com	252	2.29
net	99	0.90
gov.za	82	0.74
org.za	44	0.40
ac.za	33	0.30
africa	13	0.12
ch	11	0.10

Documents size (in segments)

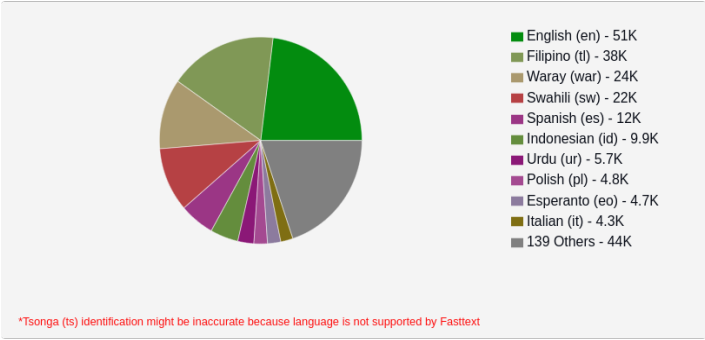


Documents by collection

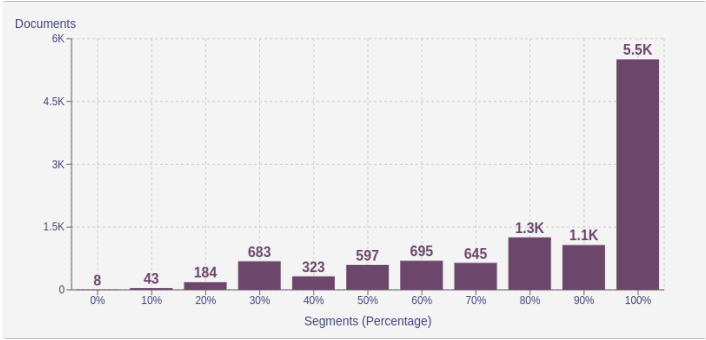


Language Distribution

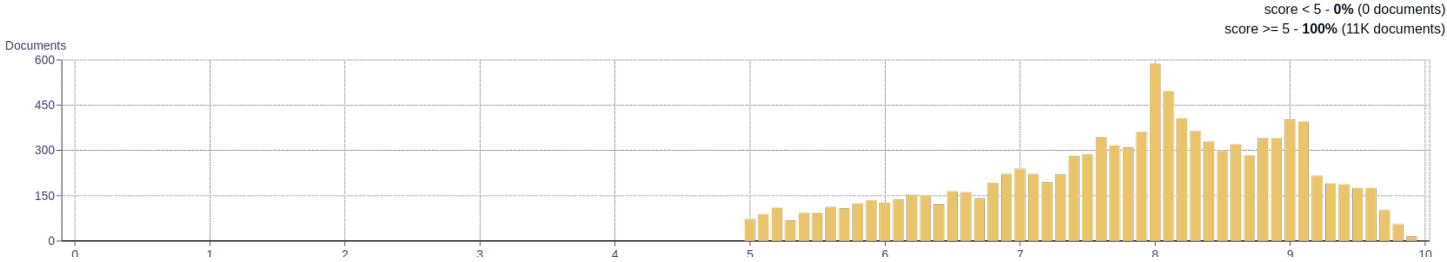
Number of segments



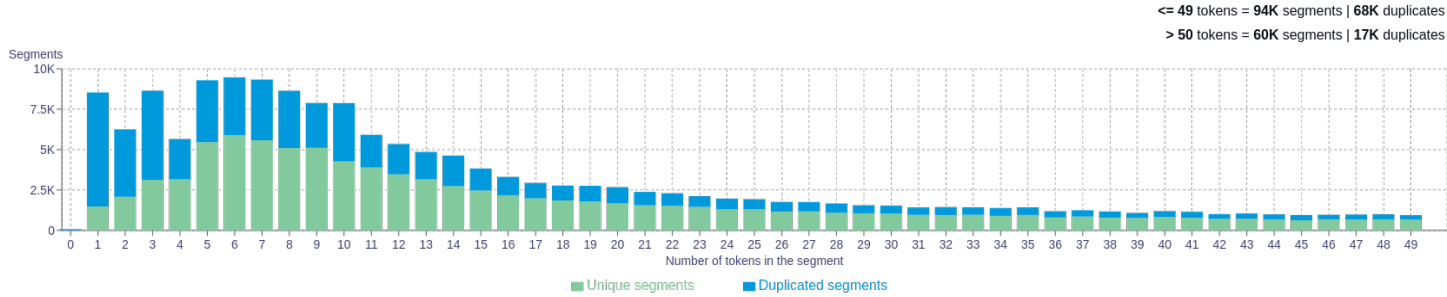
Percentage of segments in Tsonga (ts) inside documents



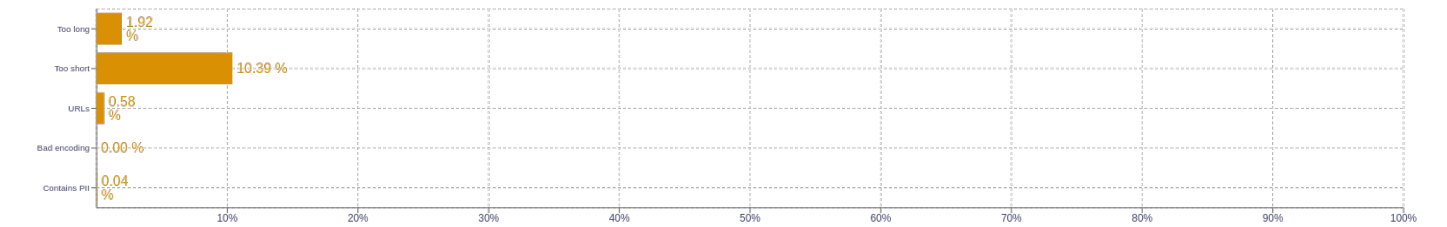
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>ni   165530</div> <div>vha   63263</div> <div>ha   59393</div> <div>wana   52896</div> <div>n   46262</div>
2	<div>ndlela leyi   10450</div> <div>ha yona   8602</div> <div>ha yini   6061</div> <div>wana ni   5889</div> <div>vanhu lava   5636</div>
3	<div>timbhoni ta yehovha   3870</div> <div>wana ni un   2301</div> <div>vanhu vo tala   2125</div> <div>bya misava leyintshwa   1677</div> <div>vuhundzuluxeri bya misava   1671</div>
4	<div>vuhundzuluxeri bya misava leyintshwa   1664</div> <div>bya misava leyintshwa bya   1294</div> <div>misava leyintshwa bya matsalwa   1271</div> <div>leyintshwa bya matsalwa yo   1250</div> <div>bya matsalwa yo kwetsima   1250</div>
5	<div>vuhundzuluxeri bya misava leyintshwa bya   1294</div> <div>bya misava leyintshwa bya matsalwa   1271</div> <div>misava leyintshwa bya matsalwa yo   1250</div> <div>leyintshwa bya matsalwa yo kwetsima   1250</div> <div>tsakela ku hlaya xihloko lexi   1188</div>

About HPLT Analytics

Volumes - Segments

Segments correspond to paragraph and list boundaries as defined by HTML elements (<p>, <ul>, <ol>, 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>, <ul>, <ol>, 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>