

General overview

Corpus	Analytics date	Language
lus_Latn.jsonl.tsv	12/3/2024	Mizo (lus)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
160,378	3,433,373	2,006,554 (58.44 %)	146M	624.09 MB	648,740,211

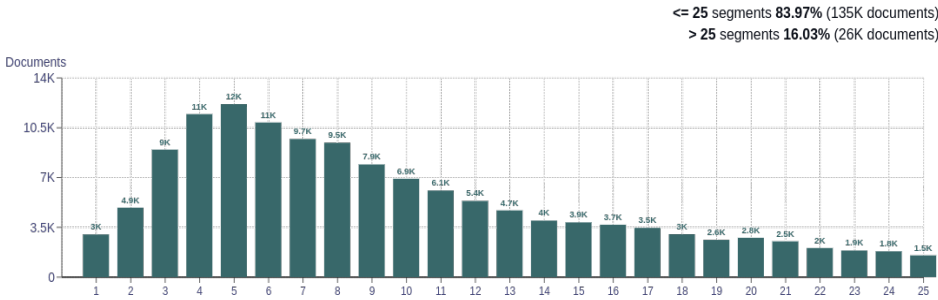
Top 10 domains

Domain	Docs	% of total
<a href="#">misual.com</a>	19K	11.97
<a href="#">khampat.com</a>	7.8K	4.88
<a href="#">zomidaily.org</a>	6.9K	4.30
<a href="#">blogspot.com</a>	6K	3.74
<a href="#">zothlifim.com</a>	3.8K	2.38
<a href="#">virthil.in</a>	3.7K	2.33
<a href="#">zonet.in</a>	3.6K	2.26
<a href="#">mizolyric.com</a>	3.6K	2.26
<a href="#">timesofmizoram.com</a>	3.4K	2.09
<a href="#">exploremizoram.com</a>	3K	1.86

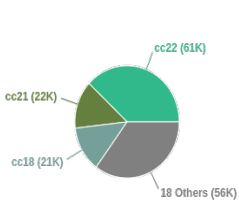
Top 10 TLDs

Domain	Docs	% of total
com	97K	60.35
org	21K	13.19
in	15K	9.15
info	7.5K	4.67
net	5.9K	3.68
gov.in	2.1K	1.32
is	1.9K	1.21
no	1.4K	0.86
nic.in	1.3K	0.83
co.in	1.3K	0.81

Documents size (in segments)

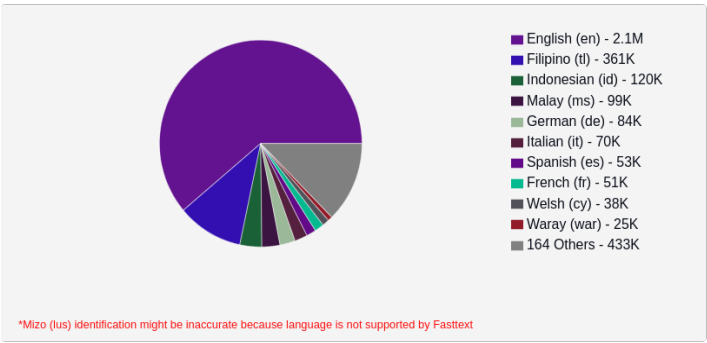


Documents by collection

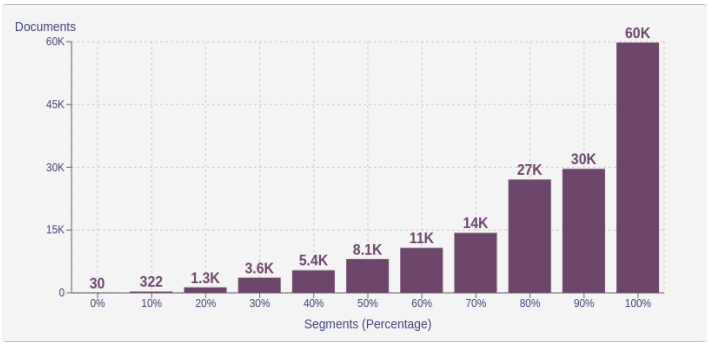


Language Distribution

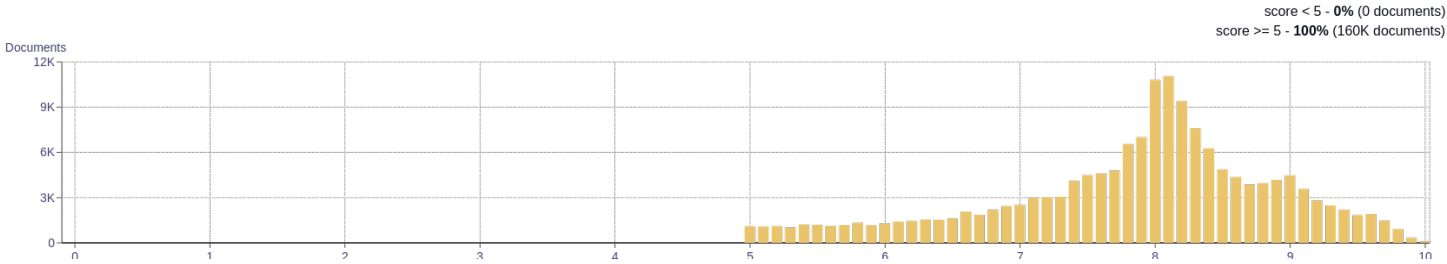
Number of segments



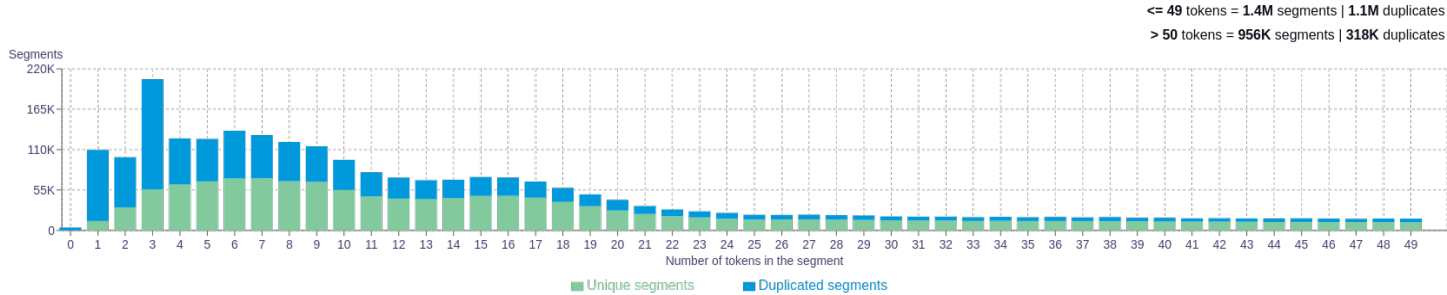
Percentage of segments in Mizo (lus) inside documents



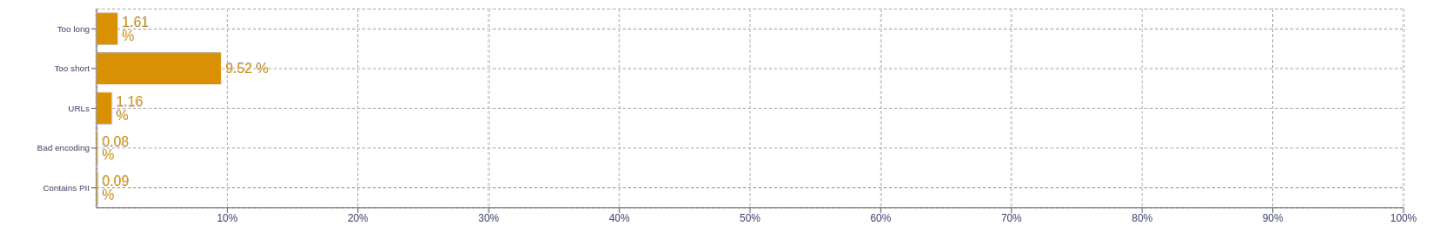
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>hi   2135292</div> <div>lo   1359152</div> <div>leh   1319316</div> <div>chu   1040545</div> <div>i   914983</div>
2	<div>uh hi   102534</div> <div>ding hi   71853</div> <div>em em   52047</div> <div>te hi   50118</div> <div>te chu   45190</div>
3	<div>thu a sawi   31666</div> <div>tiah a chim   12026</div> <div>chi hrang hrang   10653</div> <div>tiah a ti   10297</div> <div>ding uh hi   9887</div>
4	<div>positive contact atanga hri   7227</div> <div>thu a sawi bawk   6934</div> <div>contact atanga hri kai   6664</div> <div>hri kaina chhui ngai   6119</div> <div>pakhat nih a chim   4716</div>
5	<div>positive contact atanga hri kai   6661</div> <div>bang hang hiam cih leh   3068</div> <div>paragraph paragraph paragraph paragraph paragraph   2189</div> <div>chu lawmawm a tih thu   2113</div> <div>leh hri kaina chhui ngai   1424</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>