

General overview

Corpus	Date	Language
hplt-v3-bod_Tibt	9/16/2025	Tibetan (bo)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
27,863	480,612	396,286 (82.45 %)	122M	177,823,758	473.58 MB

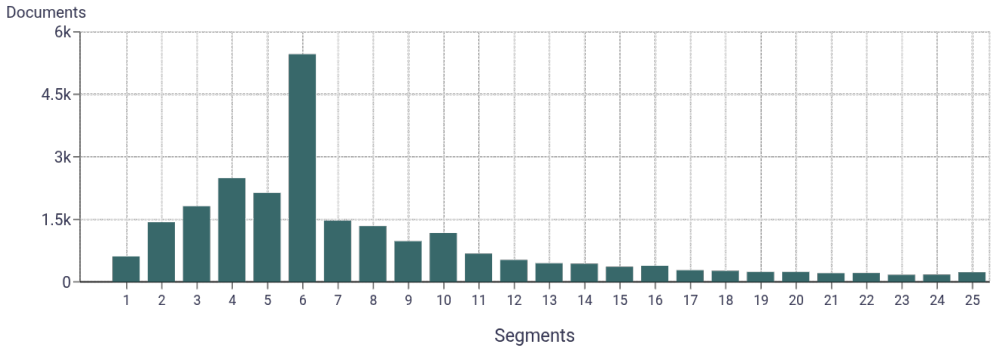
Top 10 domains

Domain	Docs	% of total
84000.co	5K	17.92%
tibettimes.net	2K	7.24%
bod.asia	1.9K	6.97%
tibet3.com	1.7K	6.05%
yongzin.com	1.4K	4.92%
zangdiyg.com	881	3.16%
people.com.cn	873	3.13%
tsadra.org	827	2.97%
tibet.cn	556	2.00%
chithu.org	519	1.86%

Top 10 TLDs

Domain	Docs	% of total
com	9.2K	33.13%
co	5K	17.92%
org	4.8K	17.28%
net	2.8K	9.95%
asia	2K	7.28%
cn	1.8K	6.45%
com.cn	918	3.29%
edu	369	1.32%
gov.cn	252	0.90%
us	74	0.27%

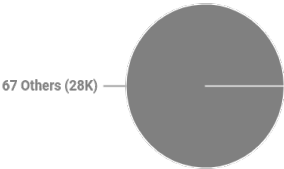
Documents size (in segments) ⓘ



≤ 25 segments **85.3%** (24K documents)
> 25 segments **14.7%** (4.1K documents)

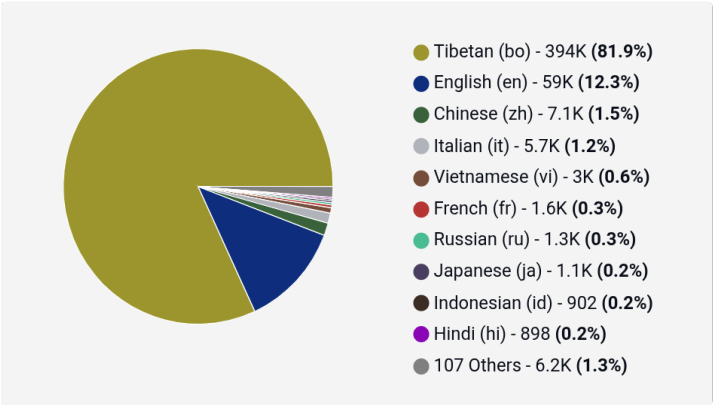
Document collections

CC = **95.21%**
IA = **4.79%**

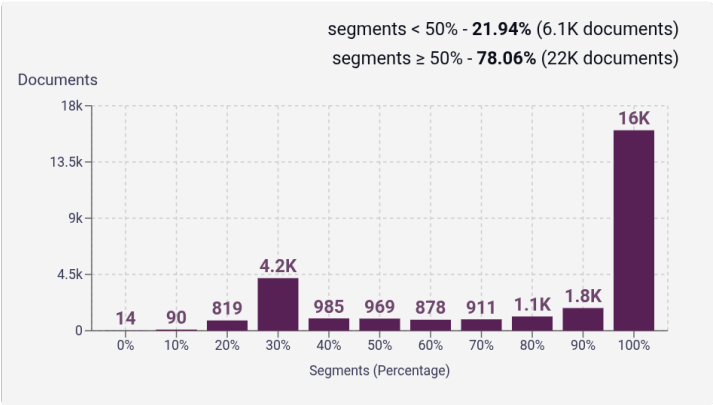


Language Distribution

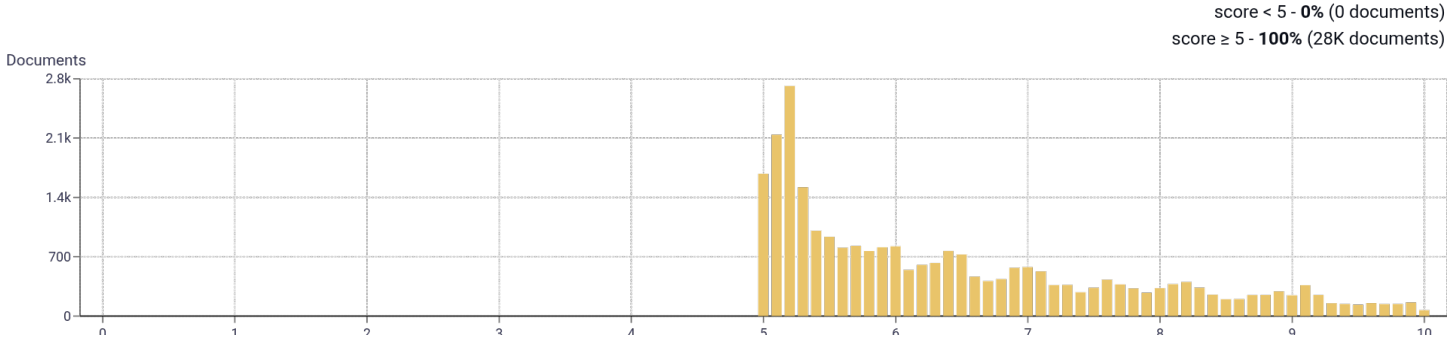
Number of segments in the Tibetan (bo) corpus



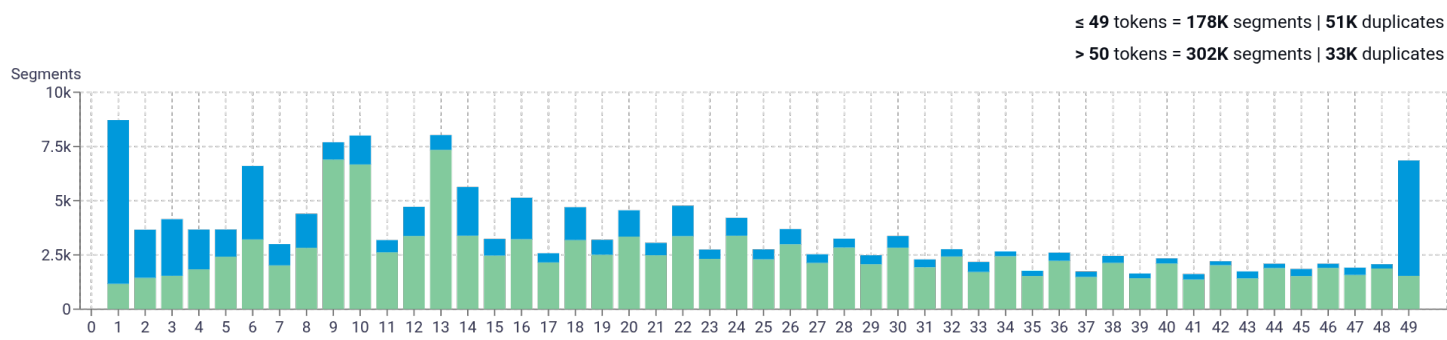
Percentage of segments in Tibetan (bo) inside documents



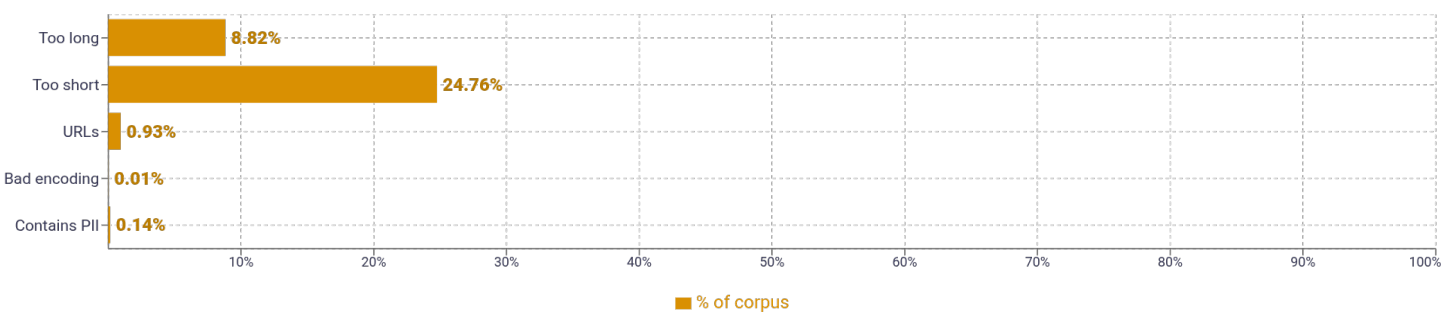
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	<div>གས 1,431,122</div> <div>དང 986,110</div> <div>ལས 665,867</div> <div>ལཱ 611,565</div> <div>ལད 538,381</div>	
2	<div>of the 17,667</div> <div>the tibetan 10,831</div> <div>the u 9,990</div> <div>the ekangyur 9,972</div> <div>in the 9,075</div>	
3	<div>the university of 5,064</div> <div>university of virginia 5,005</div> <div>tibetan and himalayan 5,005</div> <div>the tibetan and 5,001</div> <div>a variety of 4,999</div>	
4	<div>the tibetan and himalayan 4,996</div> <div>the university of virginia 4,992</div> <div>tibetan and himalayan library 4,991</div> <div>used in this translation 4,988</div> <div>translations are made from 4,988</div>	
5	<div>the tibetan and himalayan library 4,990</div> <div>translations are made from a 4,988</div> <div>tibetan page displayed here is 4,988</div> <div>the tibetan page displayed here 4,988</div> <div>the readings used in this 4,988</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				