

General overview

Corpus	Date	Language
hplt-v3-hye_Armn	9/17/2025	Armenian

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
6,120,714	104,456,778	59,019,021 (56.50 %)	2.8B	16,542,841,938	27.92 GB

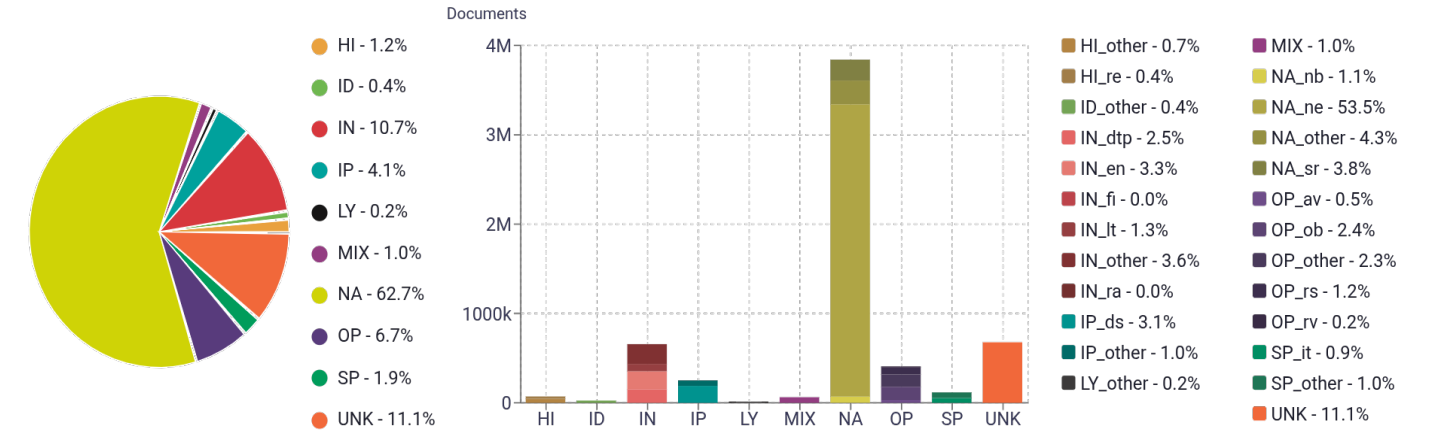
Top 10 domains

Domain	Docs	% of total
armtimes.com	169K	2.75%
wikipedia.org	165K	2.70%
azatutyun.am	145K	2.37%
news.am	130K	2.12%
armeniasputnik.am	103K	1.68%
aravot.am	89K	1.46%
mediamall.am	73K	1.20%
sputniknews.ru	64K	1.05%
7or.am	62K	1.01%
lurer.com	56K	0.91%

Top 10 TLDs

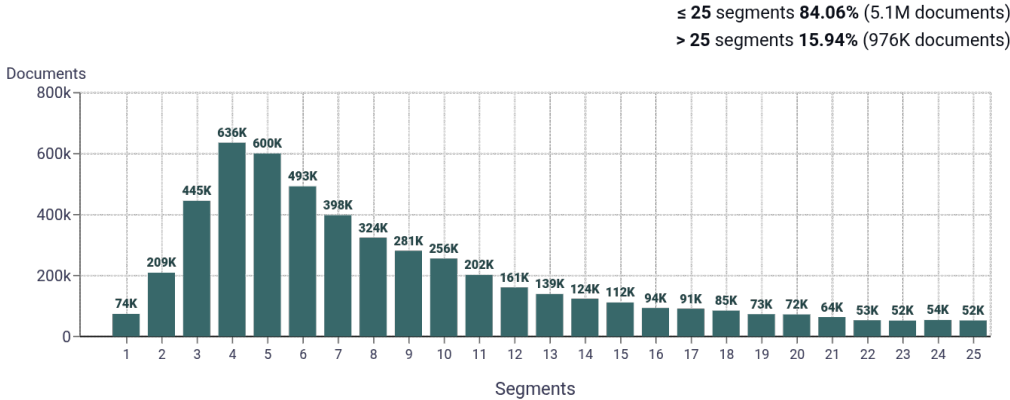
Domain	Docs	% of total
am	3.8M	61.41%
com	1.1M	18.15%
org	409K	6.68%
ru	213K	3.48%
info	105K	1.71%
net	100K	1.63%
news	65K	1.06%
tv	42K	0.69%
today	31K	0.51%
ir	31K	0.50%

Register labels



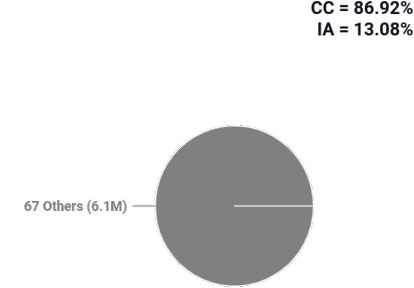
MT:4.7% | 285K Documents

Documents size (in segments) ⓘ



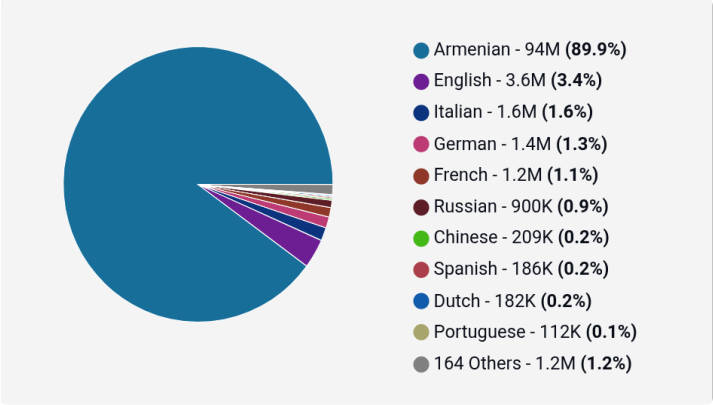
≤ 25 segments 84.06% (5.1M documents)  
> 25 segments 15.94% (976K documents)

Document collections

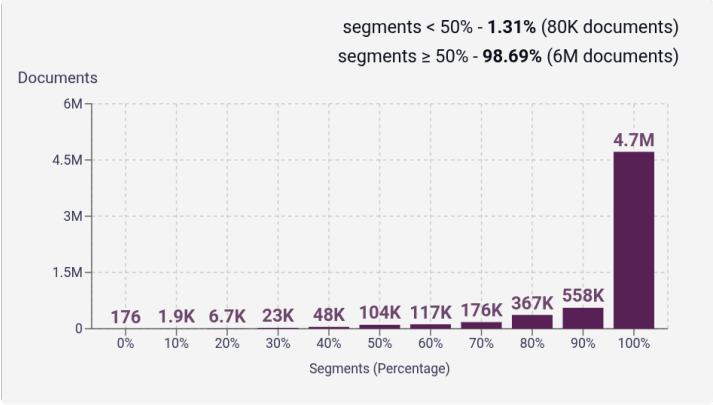


Language Distribution

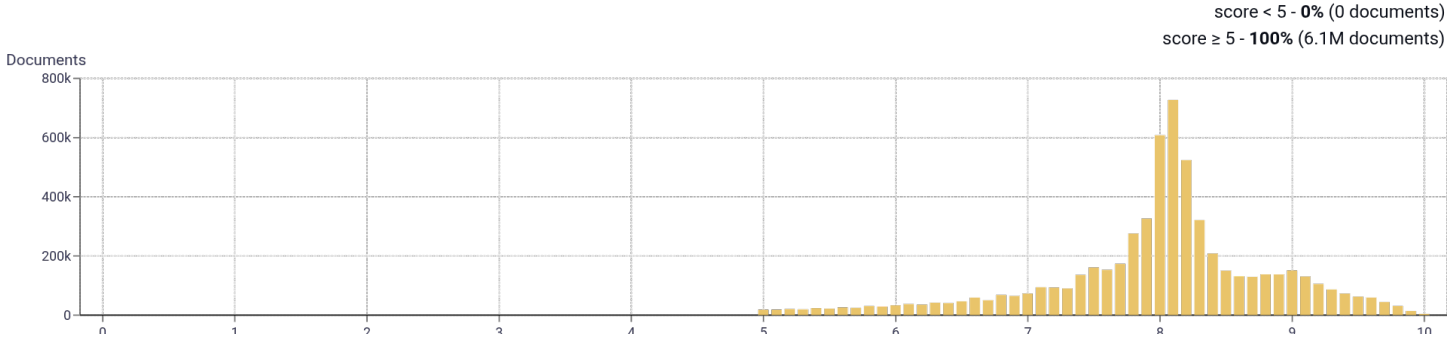
Number of segments in the Armenian corpus



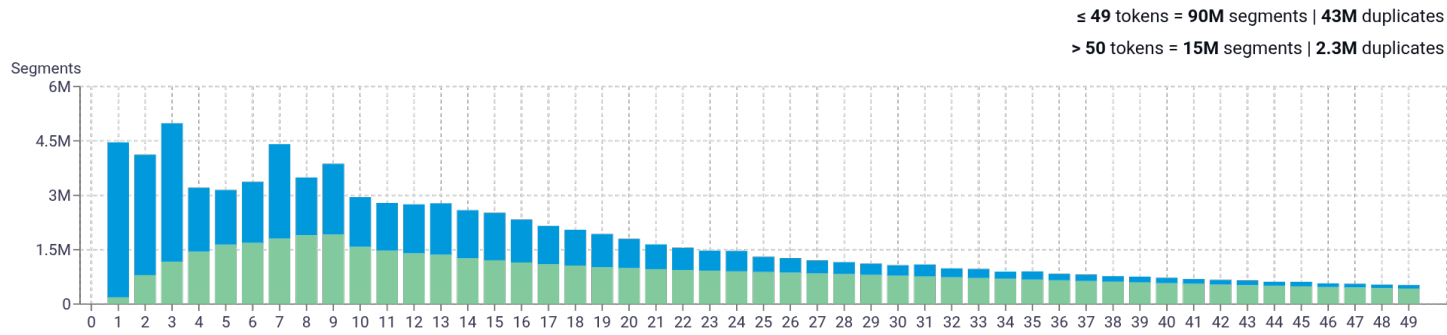
Percentage of segments in Armenian inside documents



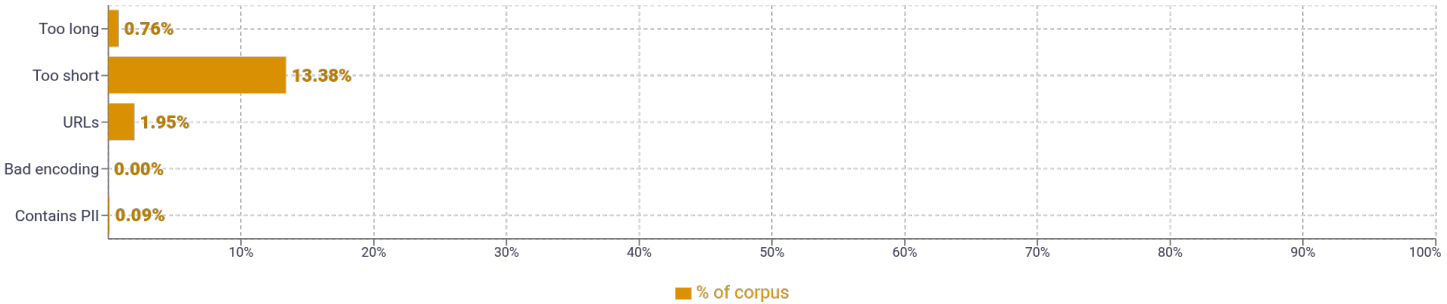
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

SIZE	N-GRAMS	
1	ել   11,858,661մի   8,104,134կարող   6,705,086չի   5,810,350ել   5,625,432	
2	մի քանի   1,416,176հայաստանի հանրապետության   1,022,280դիտել ավելին   790,714ոչ միայն   759,969մի շարք   740,663	
3	տեղի է ունեցել   308,984թույլ է տալիս   188,763պետք է լինի   181,587կարող է լինել   172,456ցույց է տալիս   152,083	
4	հարուցվել է քրեական գործ   55,826հի վարչապետ Նիկոլ փաշինյանը   54,138կայքում արտահայտված կարծիքները կարող   50,585ֆեյսբուքյան իր էջում գրել   48,789կարող են չհամընկնել խմբագրության   48,398	
5	արտահայտված կարծիքները կարող են չհամընկնել   48,066կարծիքները կարող են չհամընկնել խմբագրության   47,998կարող են չհամընկնել խմբագրության տեսակետի   47,979մասնակի կամ ամբողջական հեռուստաընթերցումներ կատարելիս   47,959Նկատմամբ վստահության եւ հարգանքի ձեւավորմանը   44,586	

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>

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				