

General overview

Corpus	Date	SL	TL
hplt-v2-en-az.tsv	1/23/2025	English (en)	Azerbaijani (az)

Volumes

Segments	SL tokens	SL characters	SL size
3,188,231	73M	386,240,961	369.86 MB

TL tokens	TL characters	TL size
62M	387,260,038	426.3 MB

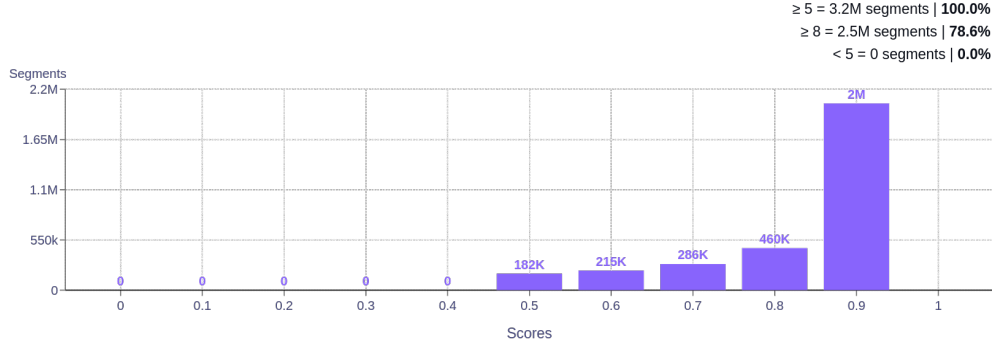
Dataset top 10 domains

SL domain	Segments	TL domain	Segments
wikipedia.org	3.7%	wikipedia.org	3.5%
software.net	3.4%	software.net	2.8%
report.az	2.3%	report.az	2.3%
trend.az	2.1%	trend.az	1.9%
educationbro.com	1.9%	president.az	1.6%
itsmygame.org	1.9%	vsaduidoma.com	1.6%
president.az	1.8%	dualjuridik.org	1.5%
vsaduidoma.com	1.6%	itsmygame.org	1.3%
dualjuridik.org	1.6%	jw.org	1.3%
jw.org	1.3%	aliyev-heritage.org	1.3%

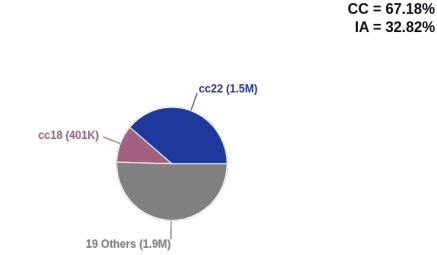
Dataset top 10 TLDs

SL domain	Segments	TL domain	Segments
com	55.7%	com	42.7%
az	22.0%	az	26.6%
org	21.9%	org	19.7%
net	13.4%	net	10.7%
gov.az	4.8%	gov.az	4.9%
co.uk	1.5%	edu.az	1.3%
edu.az	1.3%	info	1.1%
info	1.2%	co.uk	1.1%
eu	1.1%	eu	0.8%
co	0.7%	ru	0.6%

Translation likelihood

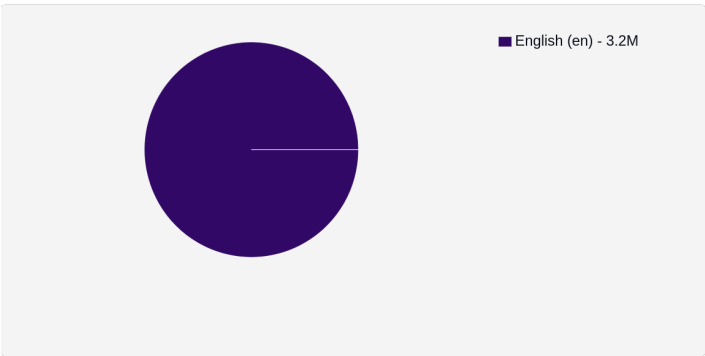


Collections

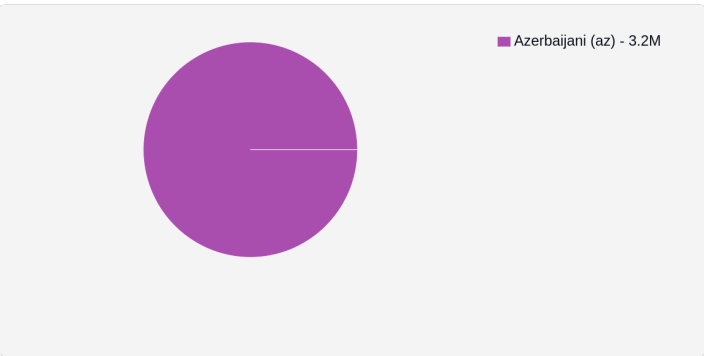


Language Distribution

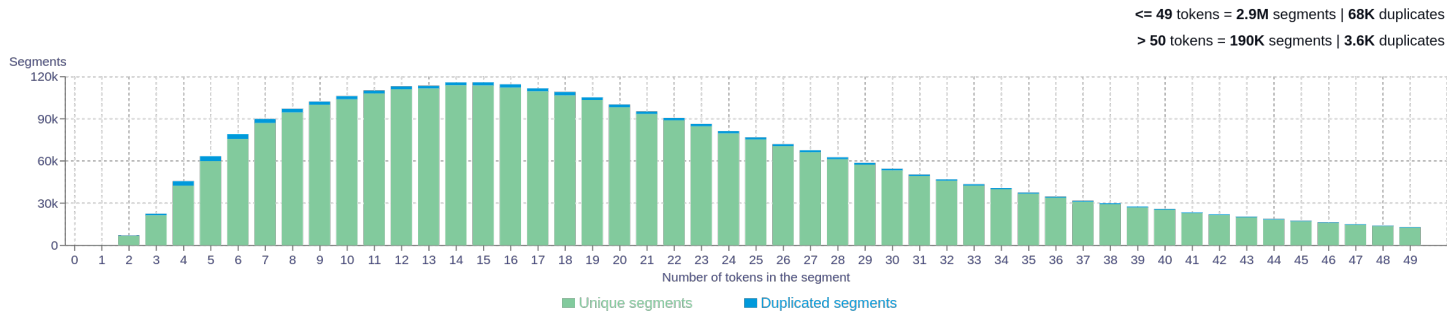
Source



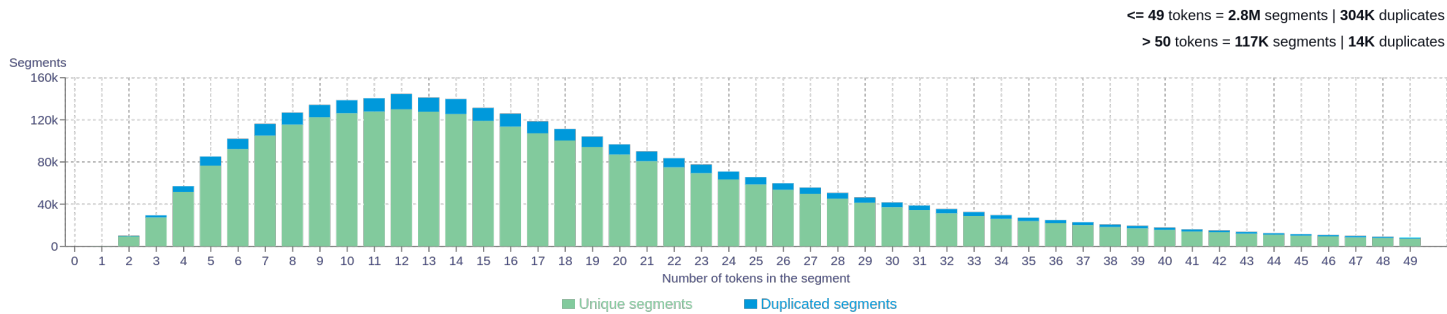
Target



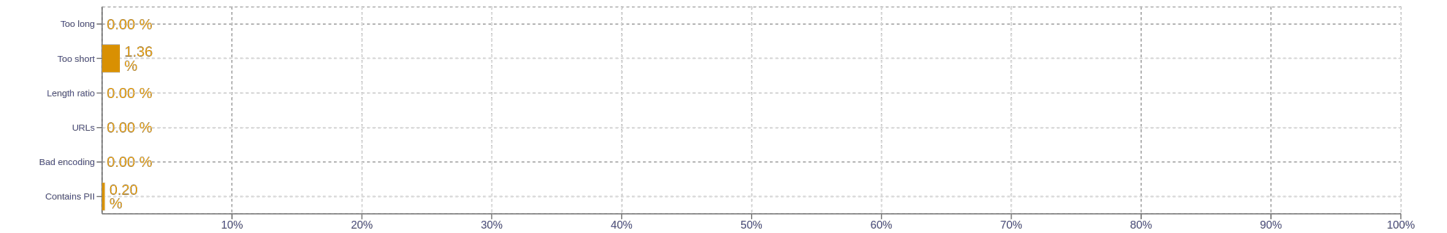
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	azerbaijan   239055also   127582one   119111time   97843new   95567
2	ilham aliyeve   22311heydar aliyeve   15729united states   12811armed forces   11305personal data   11260
3	republic of azerbaijan   42376president ilham aliyeve   10577azerbaijan ilham aliyeve   6968like the game   4653around the world   4521
4	president of the republic   13213republic of azerbaijan ilham   5241games like the game   4418ministry of foreign affairs   3146archived from the original   3006
5	republic of azerbaijan ilham aliyeve   5223technical characteristics of the game   2655general staff of armed forces   2114armed forces of ukraine says   2068foreign affairs of the republic   2048

Target n-grams

Size	n-grams
1	azərbaycan   201024böyük   103937tərəfindən   98791yeni   97624digər   94771
2	azərbaycan respublikasının   33025ədə bilərsiniz   28257azərbaycan respublikası   21615eyni zamanda   19684imkan verir   15814
3	azərbaycan respublikasının prezidenti   9585xarici işlər naziri   6818prezidenti ilham əliyev   6509prezident ilham əliyev   5764respublikasının prezidenti ilham   5509
4	azərbaycan respublikasının prezidenti ilham   5486respublikasının prezidenti ilham əliyev   3926azərbaycan prezidenti ilham əliyev   2393pulsuz mp3 mahnıları dinləyə   2176mp3 mahnıları dinləyə bilərsiniz   2176
5	azərbaycan respublikasının prezidenti ilham əliyev   3910pulsuz mp3 mahnıları dinləyə bilərsiniz   2176musiqi pulsuz mp3 mahnıları dinləyə   2168xarici işlər naziri elmar məmmədیارov   1615ukrayna silahlı qüvvələrinin baş qərargahının   1388

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 or 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>