

General overview

Corpus	Date	Language
ckb_Arab.jsonl.tsv	9/20/2024	Central Kurdish (ckb)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
273,745	5,225,884	2,963,957 (56.72 %)	171M	907,857,638	1.55 GB

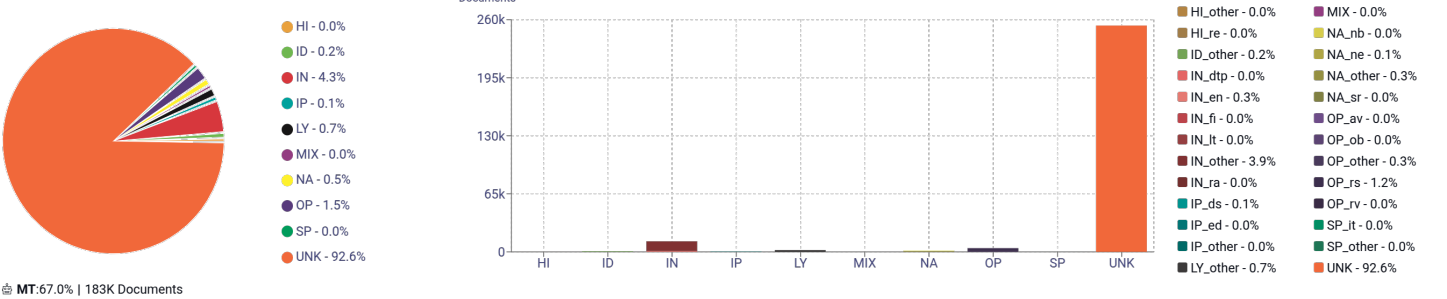
Top 10 domains

Domain	Docs	% of total
denglamerika.com	38K	13.82%
hawlati.co	7.8K	2.85%
blogfa.com	6.2K	2.28%
awene.com	5.1K	1.87%
penusakan.com	4.6K	1.70%
kurdistantv.net	4.4K	1.60%
danglislam.org	4.3K	1.58%
blogspot.com	3.9K	1.42%
payam.tv	3.8K	1.41%
wishe.net	3.7K	1.35%

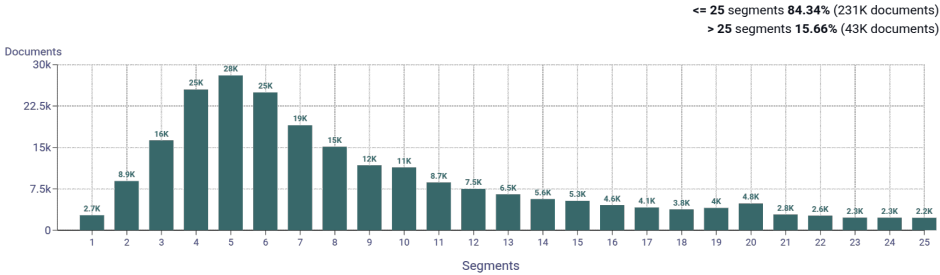
Top 10 TLDs

Domain	Docs	% of total
com	140K	50.97%
net	50K	18.38%
org	33K	12.22%
co	10K	3.80%
info	7.3K	2.68%
krd	6.8K	2.47%
tv	6.4K	2.35%
ir	3.7K	1.33%
ca	3.6K	1.33%
se	2.2K	0.80%

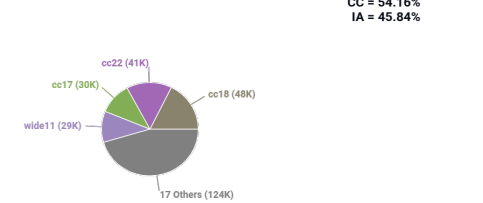
Register labels



Documents size (in segments)

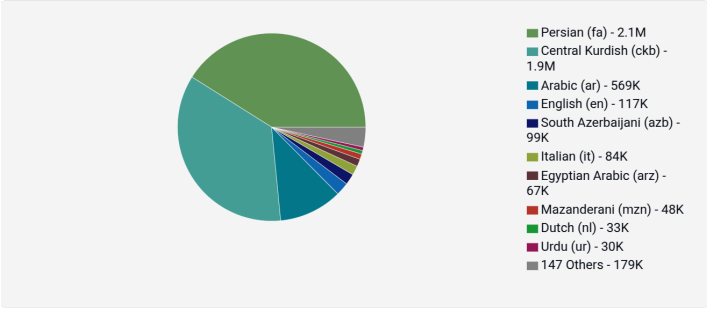


Documents by collection

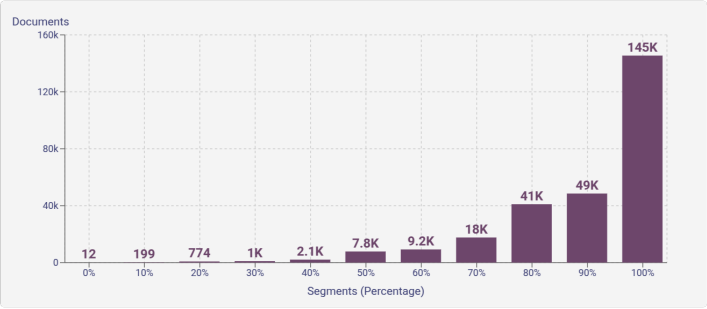


Language Distribution

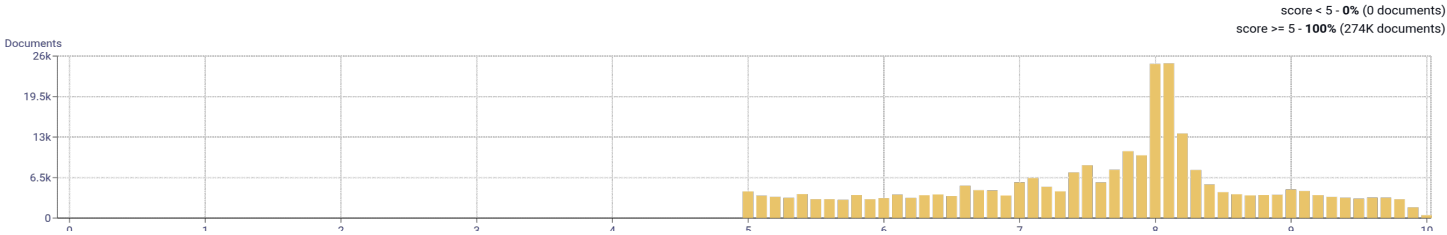
Number of segments in the Central Kurdish (ckb) corpus



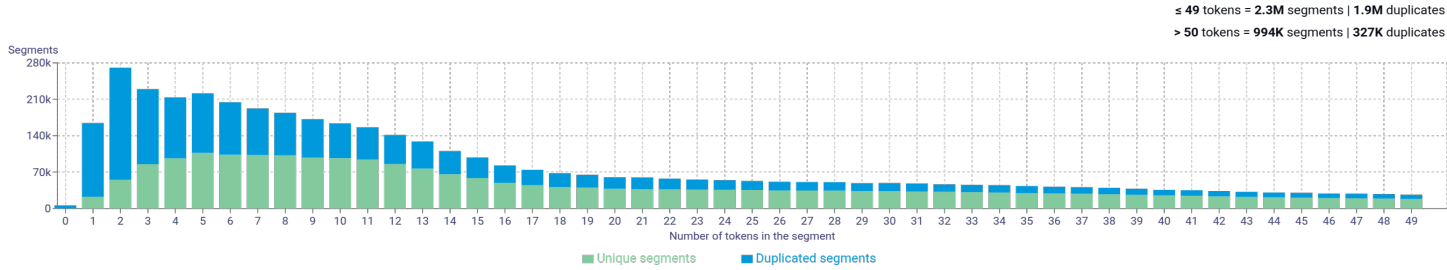
Percentage of segments in Central Kurdish (ckb) inside documents



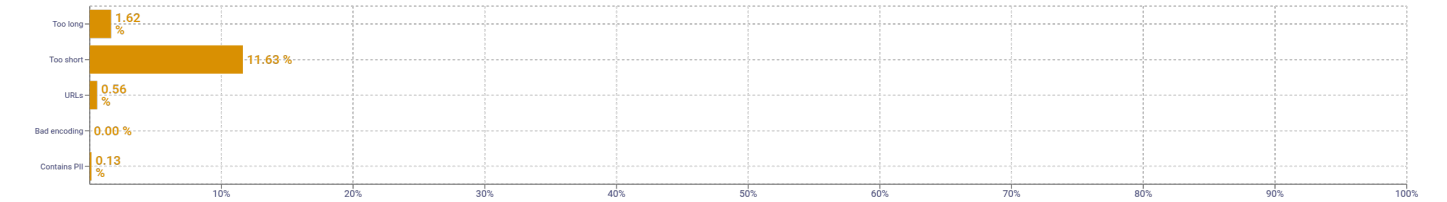
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	6089957 ی 4498158 له 1204616 که 1168762 نهو 933437 که
2	181781 ان ی 169315 یک ی 128083 ی ی 105117 ی نه 89001 که له
3	43848 له سال ی 41560 ولات ان ی 29221 خوا ی گهوره 26341 الله عليه وسلم 24269 صلى الله عليه
4	22939 صلى الله عليه وسلم 18964 صلى الله عليه وسلم 7223 له ولات ان ی 5153 شا بان ی باسه 5148 له رنگهی کلیککردنی نهو
5	5113 3358 له رنگهی کلیککردنی نهو فایله 3358 رنگهی کلیککردنی نهو فایله دهنگیبا نهو 3346 خوا ی لن بن ت 2494 خوا ی لهسهر بن ت 2072 گفتوگوکه بن له رنگهی کلیککردنی

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 *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>, , , 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				