

General overview

Corpus	Analytics date	Language
kaz_Cyrl.jsonl.tsv	9/21/2024	Kazakh (kk)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
2,637,363	81,006,479	42,567,471 (52.55 %)	1.8B	18.78 GB	11,053,418,553

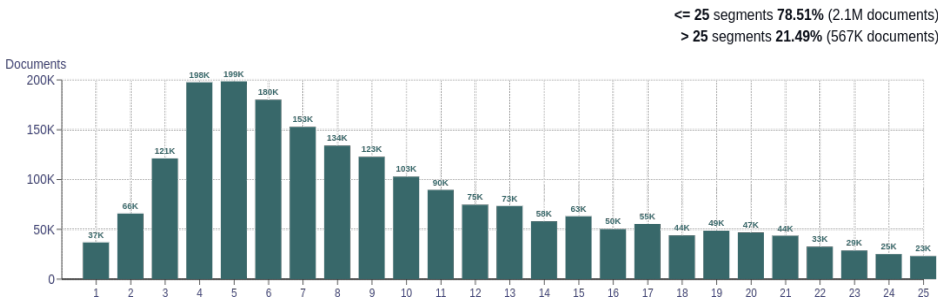
Top 10 domains

Domain	Docs	% of total
wikipedia.org	154K	5.84
azattyq.org	143K	5.43
strategy2050.kz	116K	4.41
kodeksy-kz.com	62K	2.37
tengrinews.kz	58K	2.22
nur.kz	57K	2.17
inform.kz	50K	1.91
stud.kz	47K	1.77
massaget.kz	37K	1.40
baq.kz	35K	1.32

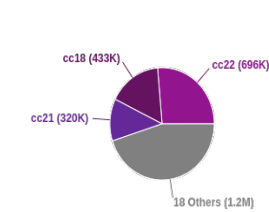
Top 10 TLDs

Domain	Docs	% of total
kz	1.8M	67.51
org	350K	13.25
com	216K	8.21
ru	77K	2.93
gov.kz	64K	2.41
net	29K	1.10
info	29K	1.08
edu.kz	23K	0.86
mobi	6.4K	0.24
uz	6K	0.23

Documents size (in segments)

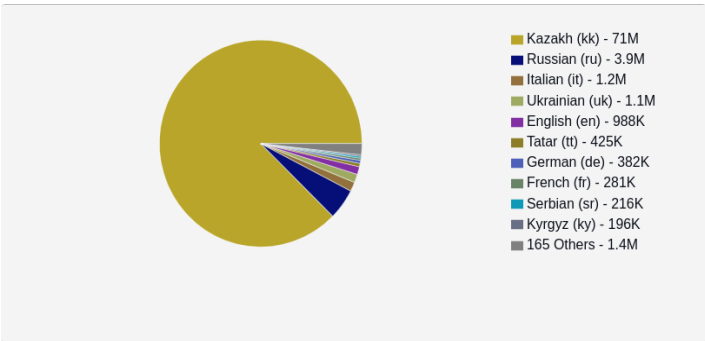


Documents by collection

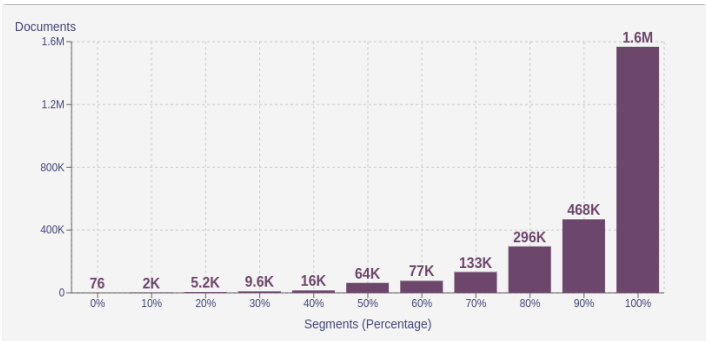


Language Distribution

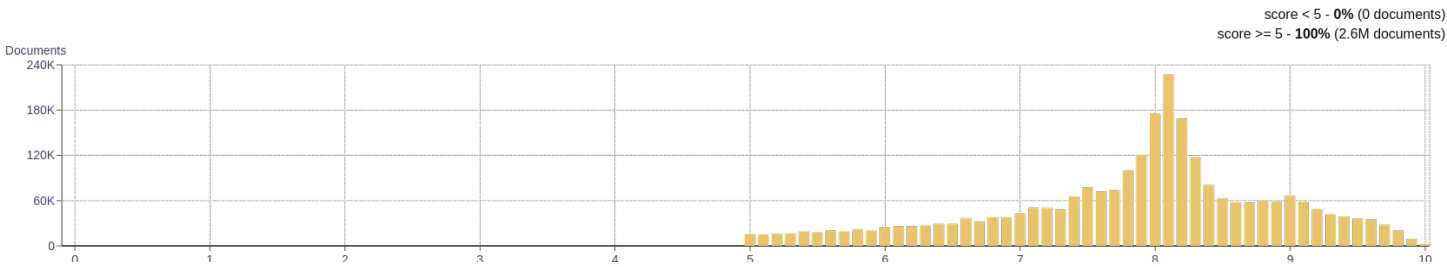
Number of segments



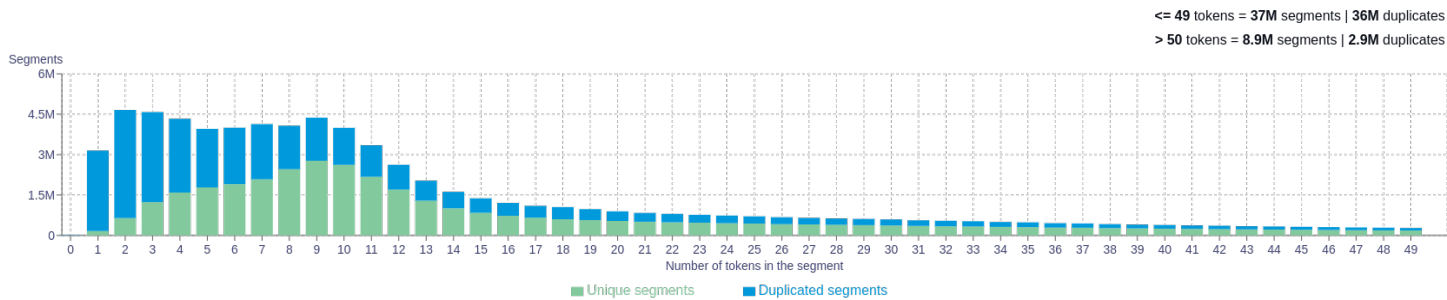
Percentage of segments in Kazakh (kk) inside documents



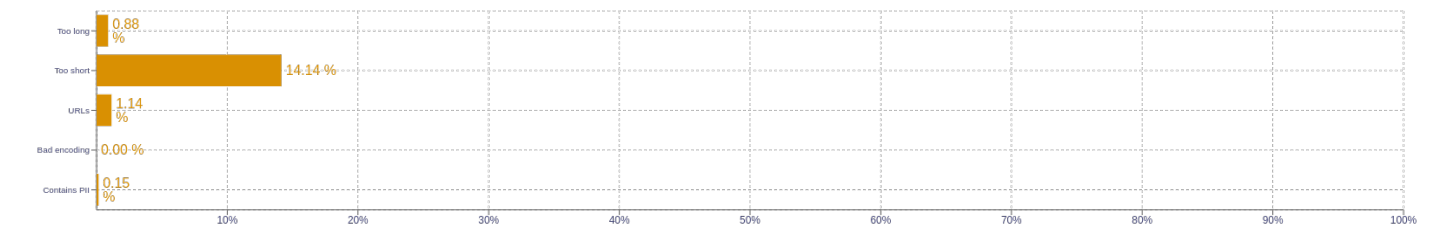
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>және   17900830<div>да   5290698</div>бір   5108395<div>бойынша   4820816</div>қазақстан   4770376</div>
2	<div>қазақстан республикасының   1407502<div>болып табылады   1009956</div>қазақстан республикасы   946957<div>білім беру   703405</div>басқа да   560917</div>
3	<div>өткен соң қолданысқа   193307<div>және басқа да   191300</div>он күн өткен   166733<div>күнтізбелік он күн   146253</div>ресми жарияланған күнінен   144430</div>
4	<div>күн өткен соң қолданысқа   102442<div>өткен соң қолданысқа енгізіледі   158912</div>күнтізбелік он күн өткен   137602<div>алғашқы ресми жарияланған күнінен   134652</div>жарияланған күнінен кейін күнтізбелік   93685</div>
5	<div>он күн өткен соң қолданысқа   166096<div>күн өткен соң қолданысқа енгізіледі   150197</div>ресми жарияланған күнінен кейін күнтізбелік   93542<div>күнінен кейін күнтізбелік он күн   85030</div>жарияланған күнінен кейін күнтізбелік он   84740</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>