

General overview

Corpus	Date	Language
Ind_Latn.jsonl.tsv	7/26/2025	Indonesian (id)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
98,141,797	2,387,131,737	1,067,290,415 (44.71 %)	64B	381,930,808,361	357.31 GB

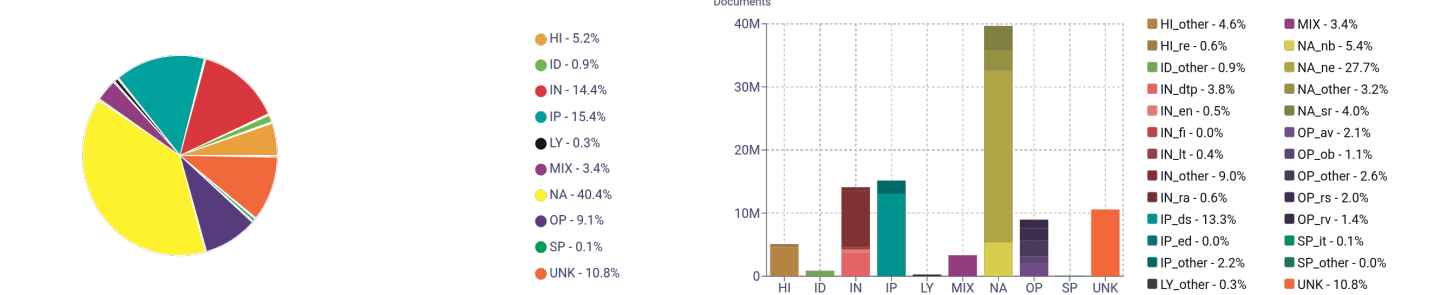
Top 10 domains

Domain	Docs	% of total
blogspot.com	8M	8.17%
wordpress.com	3.7M	3.77%
tribunnews.com	1.6M	1.68%
blogspot.co.id	1.3M	1.31%
kompas.com	783K	0.80%
blogspot.sg	470K	0.48%
us.com	406K	0.41%
detik.com	404K	0.41%
wikipedia.org	359K	0.37%
okezone.com	332K	0.34%

Top 10 TLDs

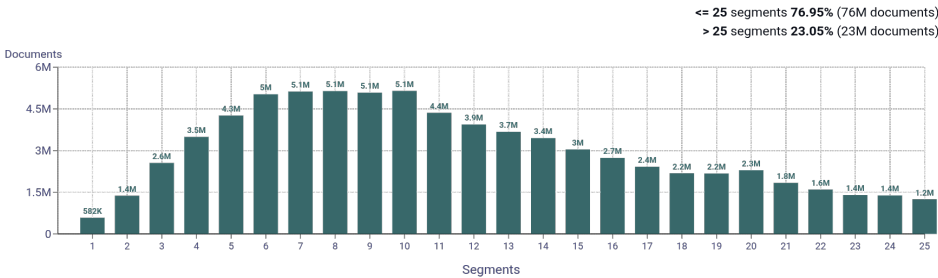
Domain	Docs	% of total
com	66M	66.77%
co.id	6M	6.07%
net	4M	4.09%
id	3.9M	4.02%
org	3.2M	3.24%
info	1.8M	1.85%
co	1.7M	1.72%
ac.id	1.6M	1.68%
go.id	1.5M	1.51%
web.id	1.1M	1.07%

Register labels



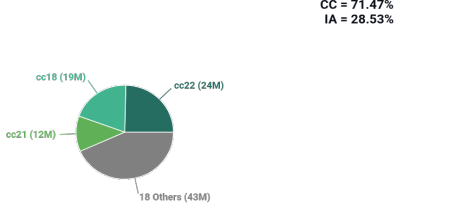
MT:4.4% | 4.4M Documents

Documents size (in segments)



<= 25 segments 76.95% (76M documents)
> 25 segments 23.05% (23M documents)

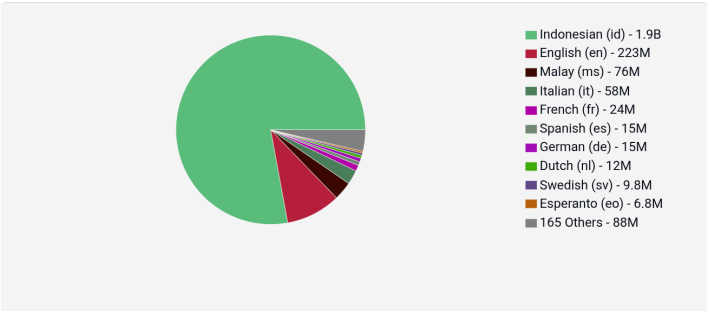
Documents by collection



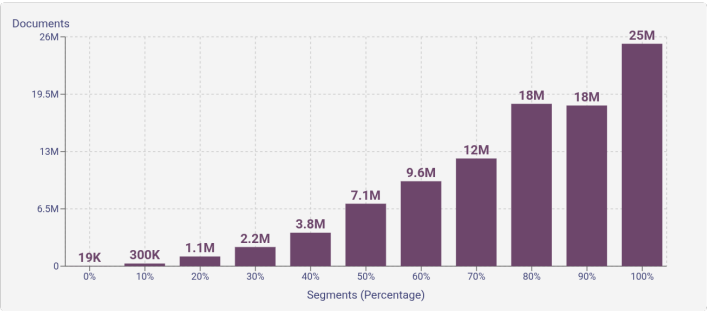
CC = 71.47%
IA = 28.53%

Language Distribution

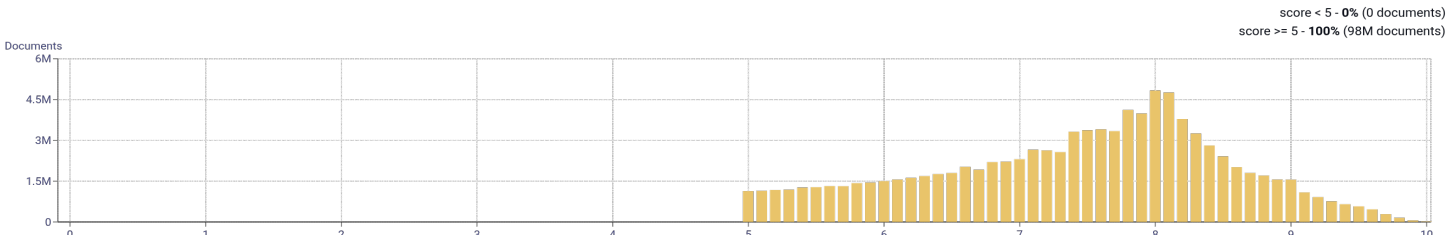
Number of segments in the Indonesian (id) corpus



Percentage of segments in Indonesian (id) inside documents

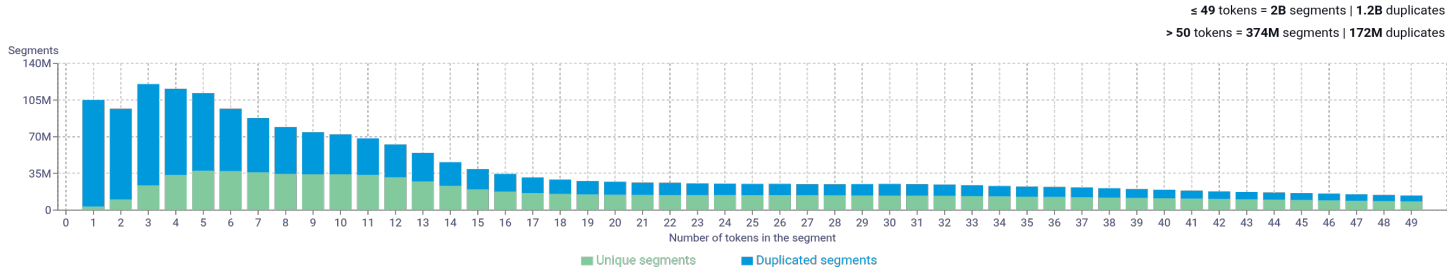


Distribution of documents by document score

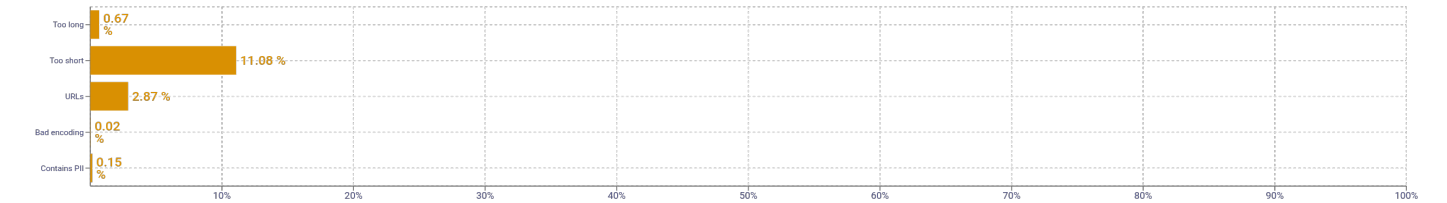


score < 5 - 0% (0 documents)
score >= 5 - 100% (98M documents)

Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>nya 919650811</div> <div>orang 134380517</div> <div>Indonesia 117256773</div> <div>memiliki 100406959</div> <div>salah 78546807</div>
2	<div>arti nya 11130962</div> <div>orang tua 10517856</div> <div>slot online 8471962</div> <div>nama nya 8036915</div> <div>judi online 7947441</div>
3	<div>salah satu nya 7814827</div> <div>bab i pendahuluan 5658857</div> <div>judi slot online 2631678</div> <div>tuhan yang maha 2386592</div> <div>bolak balik sifon 2335581</div>
4	<div>jilbab bolak balik sifon 2335581</div> <div>tuhan yang maha esa 2122539</div> <div>rahmat tuhan yang maha 1689813</div> <div>bolak balik sifon polos 1556883</div> <div>bab i pendahuluan a 1484150</div>
5	<div>rahmat tuhan yang maha esa 1636772</div> <div>jilbab bolak balik sifon polos 1556883</div> <div>jilbab bolak balik pricilla warna 811984</div> <div>jilbab bolak balik sifon motif 778696</div> <div>elSa anna dan olaf ready 731809</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 *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	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				