

General overview

Corpus	Analytics date	Language
nno_latn.jsonl.tsv	9/21/2024	Norwegian Nynorsk (nn)

Volumes

Docs	Segments	Unique segments	Tokens	Size	Characters
1,423,143	34,603,343	12,530,321 (36.21 %)	983M	5.12 GB	5,371,043,855

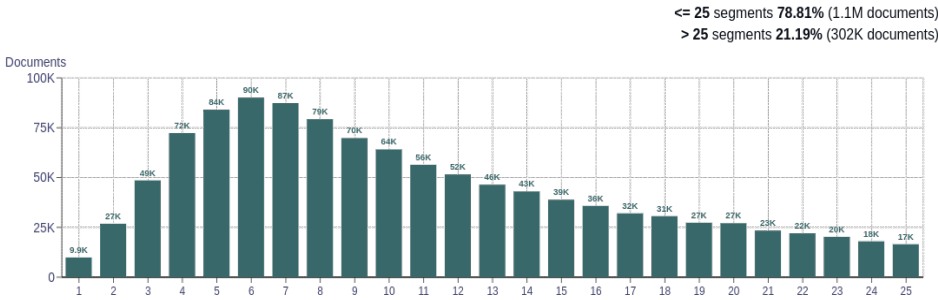
Top 10 domains

Domain	Docs	% of total
wikipedia.org	325K	22.84
blogspot.com	80K	5.60
blogspot.no	49K	3.44
nrk.no	36K	2.52
docplayer.me	31K	2.20
ndla.no	25K	1.76
blogg.no	20K	1.38
wordpress.com	19K	1.36
framtida.no	15K	1.06
allkunne.no	9.7K	0.68

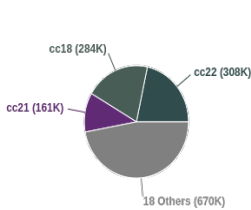
Top 10 TLDs

Domain	Docs	% of total
no	758K	53.29
org	345K	24.28
com	185K	13.03
kommune.no	35K	2.43
me	32K	2.21
net	19K	1.36
info	9.5K	0.67
eu	5K	0.35
vgs.no	4.3K	0.30
dk	1.9K	0.13

Documents size (in segments)

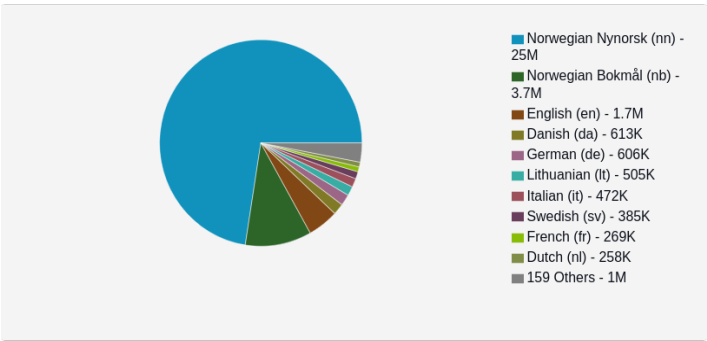


Documents by collection

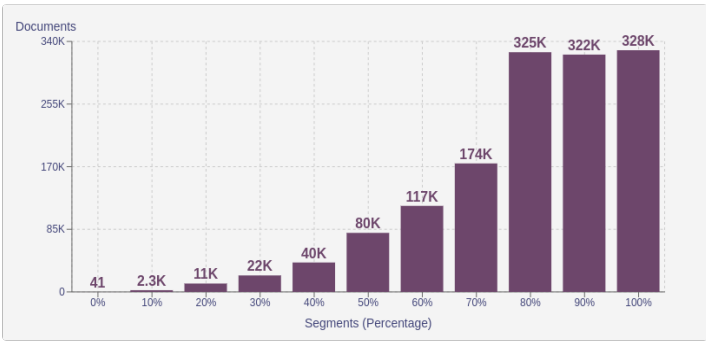


Language Distribution

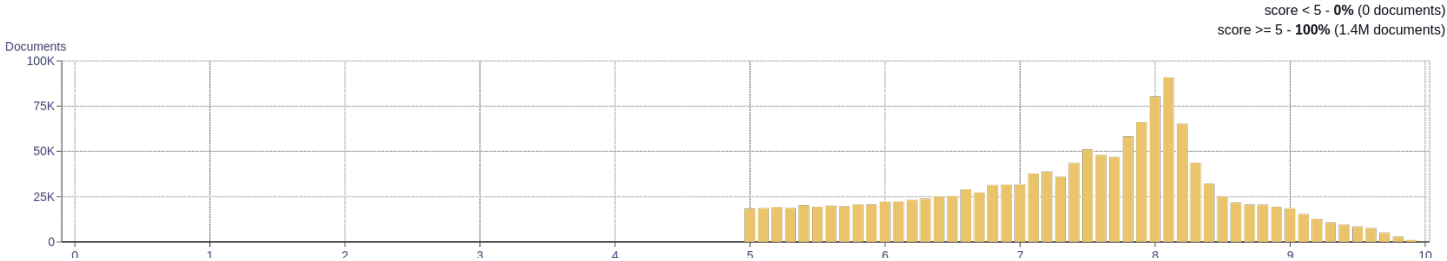
Number of segments



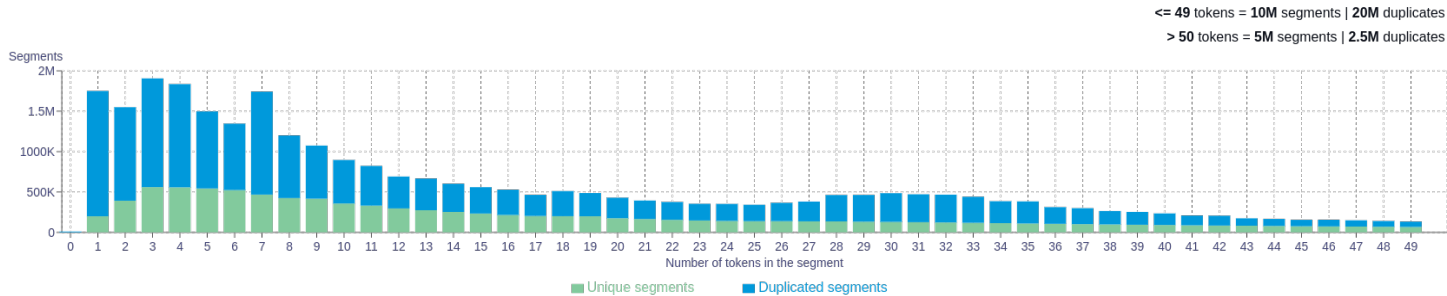
Percentage of segments in Norwegian Nynorsk (nn) inside documents



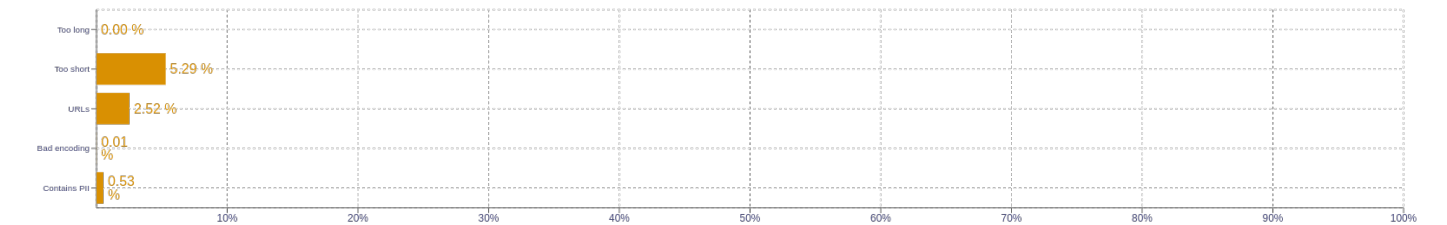
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>kommune   1946572</div> <div>endre   1707517</div> <div>år   1426019</div> <div>må   1266078</div> <div>vert   1193946</div>
2	<div>endre wikiteksten   777797</div> <div>kommune møteprotokoll   223642</div> <div>null null   140176</div> <div>utval møtedato   126664</div> <div>artikkelen bygger   126297</div>
3	<div>sogn og fjordane   262276</div> <div>møre og romsdal   230336</div> <div>null null null   139713</div> <div>wikipedia på engelsk   108166</div> <div>navn funksjon representerer   67805</div>
4	<div>null null null null   139310</div> <div>møre og romsdal fylkeskommune   45201</div> <div>nynorsk/bokmål nynorsk eksamensinformasjon eksamenstid   44853</div> <div>følgjande faste medlemmer møtte   43023</div> <div>medlemmer var til stades   42274</div>
5	<div>null null null null null   138938</div> <div>faste medlemmer var til stades   42266</div> <div>wikipedia på engelsk oppgav desse   32322</div> <div>innkalling til møtet vart gjort   30583</div> <div>møtet vart gjort i samsvar   28997</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>