

General overview

Corpus	Date	SL	TL
hplt-v2-en-nn.tsv	1/31/2025	English (en)	Norwegian Nynorsk (nn)

Volumes

Segments	SL tokens	SL characters	SL size
563,791	14M	69,346,941	66.38 MB

TL tokens	TL characters	TL size
12M	65,155,921	63.43 MB

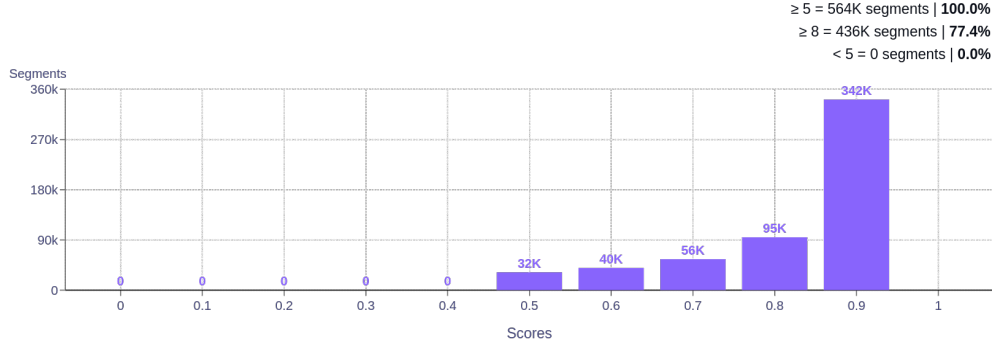
Dataset top 10 domains

SL domain	Segments	TL domain	Segments
wikipedia.org	86.3%	wikipedia.org	72.1%
libreoffice.org	3.4%	uib.no	2.6%
gimp.org	2.8%	bible.com	2.6%
uib.no	2.7%	honsi.org	2.5%
schools-wikipedia.org	2.5%	gimp.org	2.1%
honsi.org	2.4%	ssb.no	2.0%
ssb.no	2.2%	libreoffice.org	1.9%
bible.com	2.2%	visitnorway.no	1.4%
encyclopine.org	1.4%	npd.no	1.2%
visitnorway.com	1.3%	skatteetaten.no	1.1%

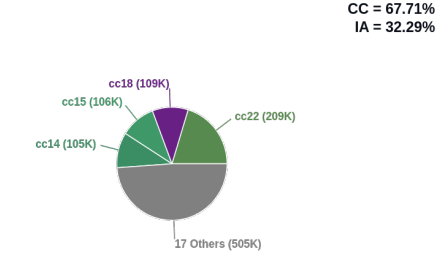
Dataset top 10 TLDs

SL domain	Segments	TL domain	Segments
org	103.5%	org	81.2%
no	32.3%	no	35.4%
com	20.7%	com	12.6%
net	1.8%	net	1.1%
ws	1.2%	ws	0.6%
edu	1.0%	info	0.4%
mobi	0.9%	mobi	0.3%
fm	0.6%	edu	0.2%
info	0.5%	me	0.2%
co.uk	0.2%	museum.no	0.2%

Translation likelihood

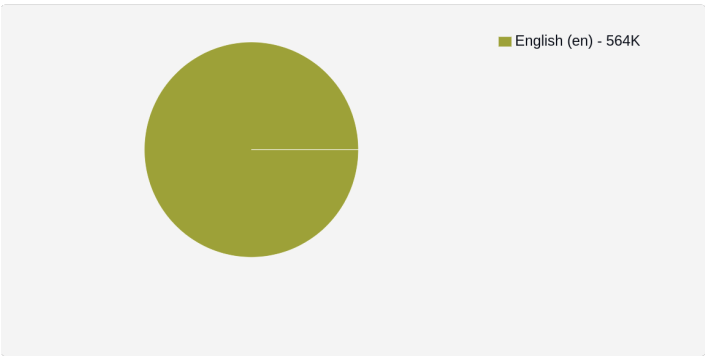


Collections

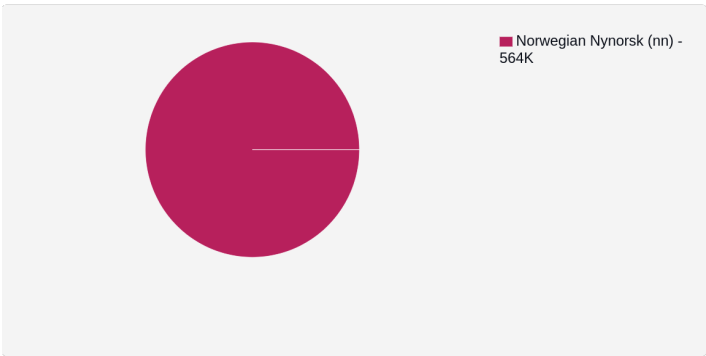


Language Distribution

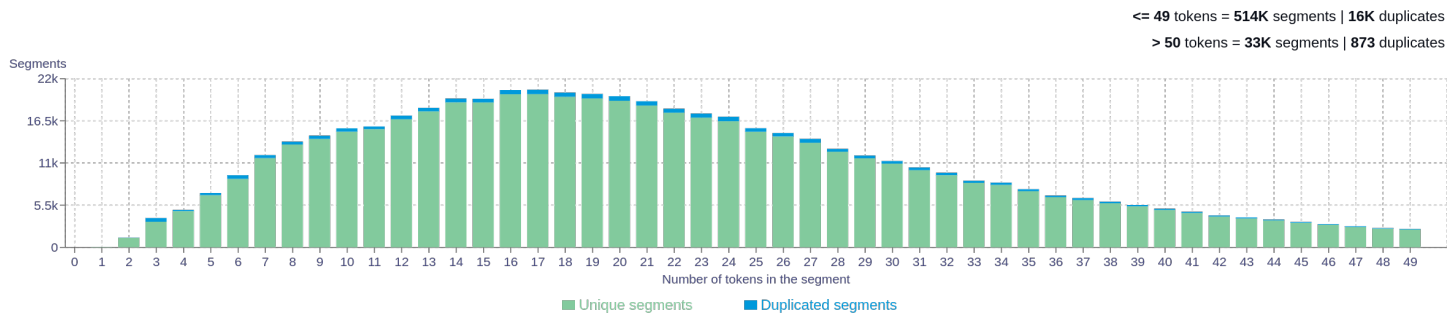
Source



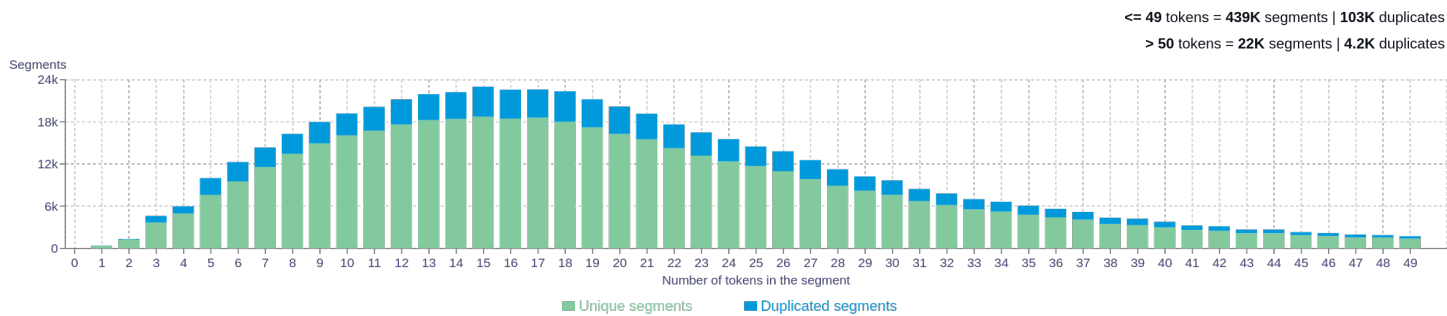
Target



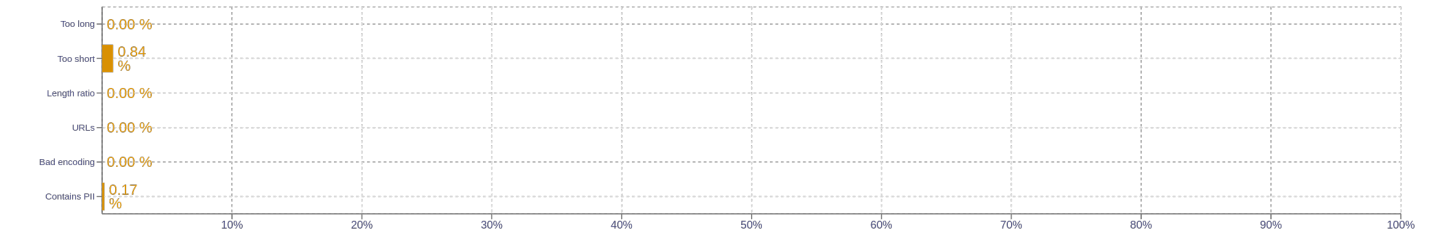
Source segment length distribution by token



Target segment length distribution by token



Segment pair noise distribution



Source n-grams

Size	n-grams
1	also 26911one 25121first 20699new 19634city 16311
2	hotels near 5153united states 3641per cent 2817world war 2381citation needed 2208
3	world war ii 1252rock and roll 803video in good 796new full version 795watch and download 707
4	video in good quality 796hd quality for free 707download in hd quality 707united states geological survey 593incorporates public domain material 564
5	watch and download in hd 707article incorporates public domain material 564material from the united states 546domain material from the united 546united states geological survey document 545

Target n-grams

Size	n-grams
1	endre 22422the 22406under 21757byen 18275ligg 18073
2	endre wikiteksten 9229hotell nær 5354new zealand 1840of the 1650new york 1624
3	wikipedia på engelsk 1246meter over havet 1242video i god 958eurovision song contest 778se og laste 702
4	video i god kvalitet 958laste ned i hd-kvalitet 702geographic names information system 562united states geological survey-artikkelen 557stoff som er offentlig 557
5	stoff som er offentlig eige 557offentleg eige frå united states 557innhald frå geographic names information 557inneheld stoff som er offentlig 557eige frå united states geological 557

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>