

General overview

Corpus	Date	Language
fin_Latn.jsonl.tsv	6/18/2025	Finnish (fi)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
34,815,557	976,395,607	369,343,799 (37.83 %)	22B	154,736,868,874	149.9 GB

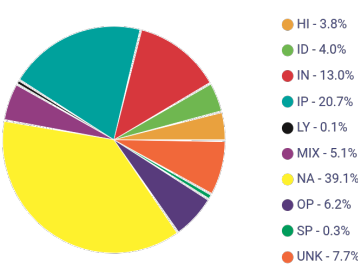
Top 10 domains

Domain	Docs	% of total
blogspot.com	3M	8.60%
blogspot.fi	2.8M	8.09%
wikipedia.org	926K	2.66%
mtv.fi	661K	1.90%
docplayer.fi	531K	1.53%
vuodatus.net	446K	1.28%
suomi24.fi	429K	1.23%
lily.fi	299K	0.86%
wordpress.com	277K	0.80%
yle.fi	238K	0.68%

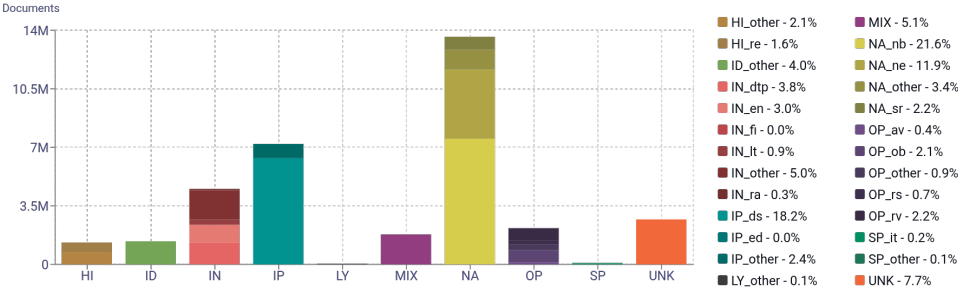
Top 10 TLDs

Domain	Docs	% of total
fi	21M	60.79%
com	9M	25.76%
net	1.6M	4.54%
org	1.3M	3.86%
eu	304K	0.87%
info	232K	0.67%
se	102K	0.29%
de	97K	0.28%
ru	76K	0.22%
no	55K	0.16%

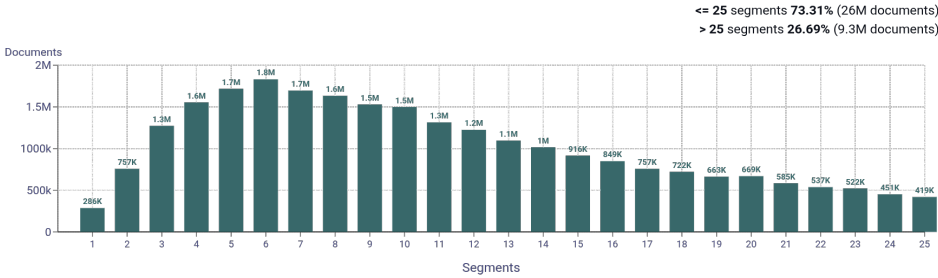
Register labels



MT:4.9% | 1.7M Documents

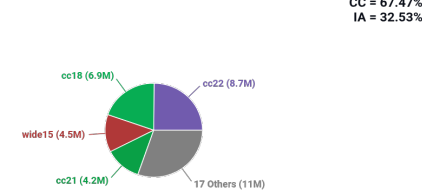


Documents size (in segments)



<= 25 segments 73.31% (26M documents)
> 25 segments 26.69% (9.3M documents)

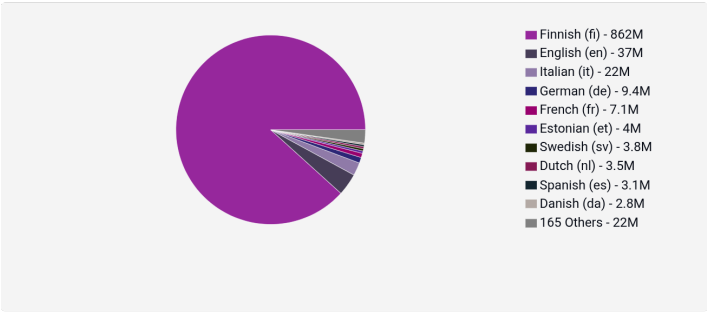
Documents by collection



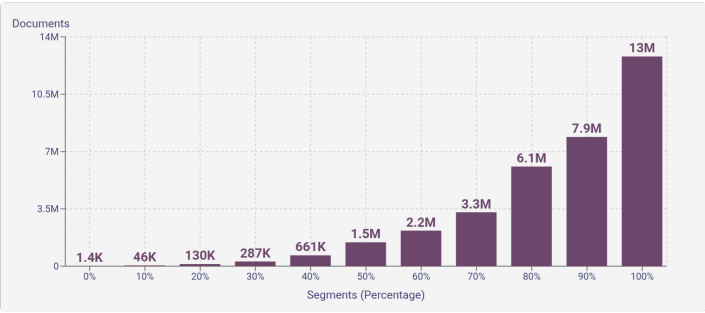
CC = 67.47%
IA = 32.53%

Language Distribution

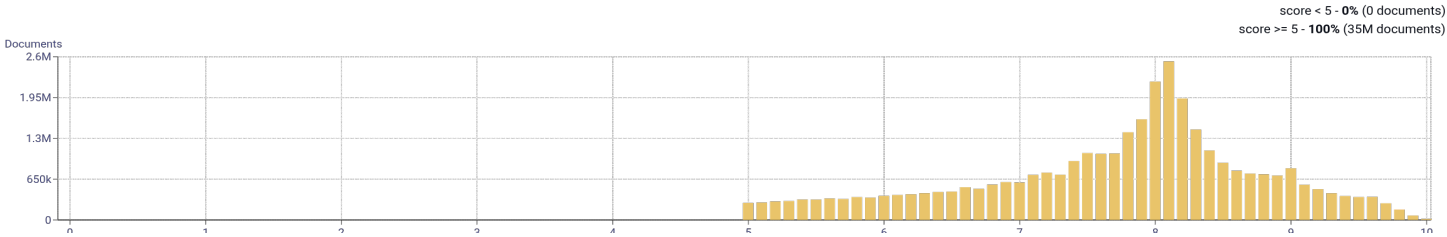
Number of segments in the Finnish (fi) corpus



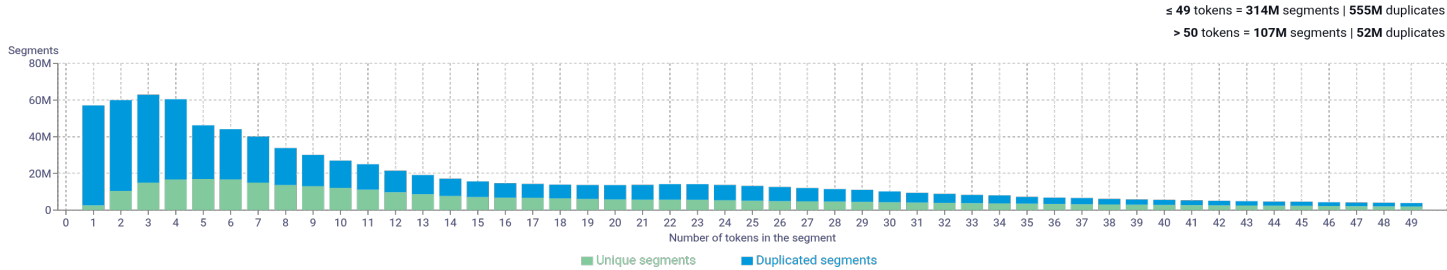
Percentage of segments in Finnish (fi) inside documents



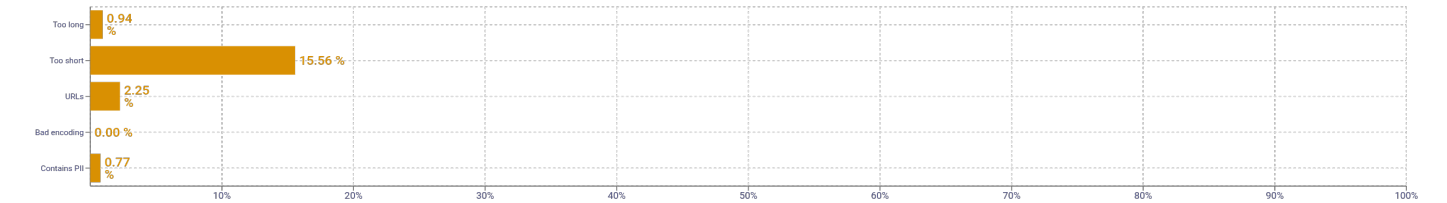
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div>voi 45786156</div> <div>jo 32904042</div> <div>vain 30085692</div> <div>sitten 29553005</div> <div>vielä 26025452</div>
2	<div>muun muassa 4164197</div> <div>tällä hetkellä 3484605</div> <div>muokkaa wikitekstiä 3162248</div> <div>lue lisää 2938806</div> <div>tällä kertaa 2168233</div>
3	<div>lasten ja nuorten 763765</div> <div>hotellia viimeisen tunnin 436579</div> <div>ketjusta on poistettu 396132</div> <div>a b c 388571</div> <div>hallituksen esitys eduskunnalle 374676</div>
4	<div>hotellia viimeisen tunnin sisällä 436579</div> <div>tarkasteli tätä hotellia viimeisen 355308</div> <div>henkilöä tarkasteli tätä hotellia 355308</div> <div>hallituksen esitys eduskunnalle laiksi 239730</div> <div>a b c d 231039</div>
5	<div>henkilöä tarkasteli tätä hotellia viimeisen 355308</div> <div>aineen tai seoksen ja yhtiön 162106</div> <div>kokouksen laillisuuden ja päätösvaltaisuuden toteaminen 149027</div> <div>a b c d e 144817</div> <div>esityksen pääasiallinen sisältö esityksessä ehdotetaan 142340</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	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	Encyclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				