

General overview

Corpus	Date	Language
lvs_Latn.jsonl.tsv	6/4/2025	Latvian (lvs)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
6,771,633	173,791,676	79,908,770 (45.98 %)	4.2B	25,015,494,704	25.39 GB

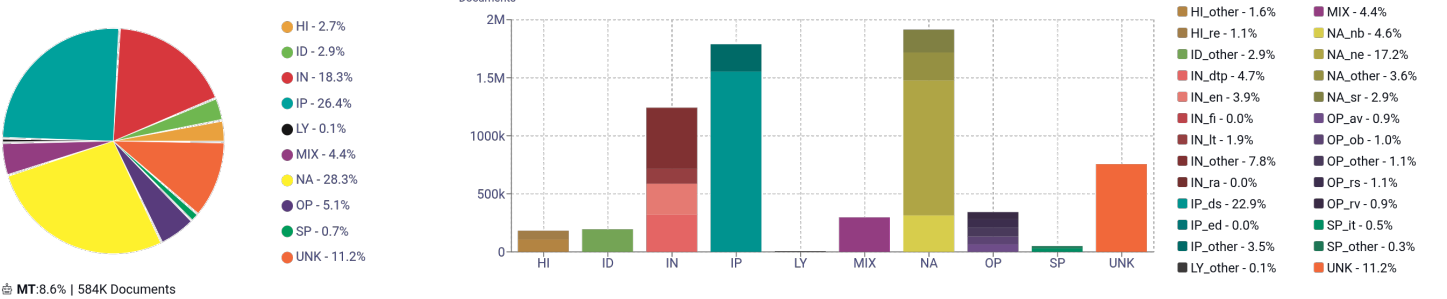
Top 10 domains

Domain	Docs	% of total
wikipedia.org	253K	3.74%
atlants.lv	199K	2.93%
tvnet.lv	122K	1.80%
skaties.lv	93K	1.37%
hotels.com	89K	1.32%
visi.lv	85K	1.26%
delfi.lv	74K	1.09%
ism.lv	72K	1.06%
agoda.com	61K	0.90%
blogspot.com	54K	0.79%

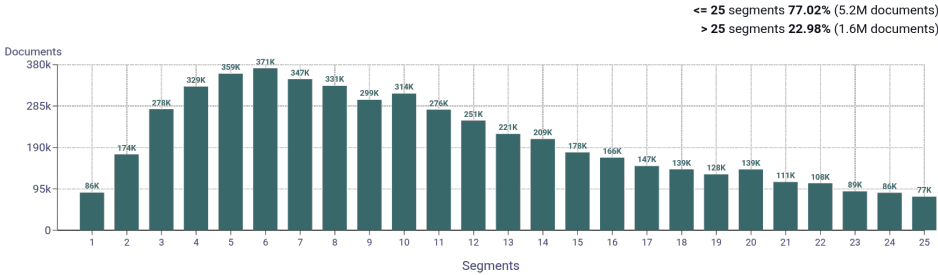
Top 10 TLDs

Domain	Docs	% of total
lv	4.8M	70.33%
com	1.1M	16.32%
org	327K	4.82%
gov.lv	126K	1.86%
eu	122K	1.80%
info	57K	0.84%
net	45K	0.66%
ie	19K	0.27%
edu.lv	17K	0.24%
lt	16K	0.23%

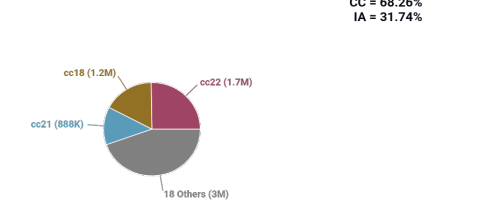
Register labels



Documents size (in segments)

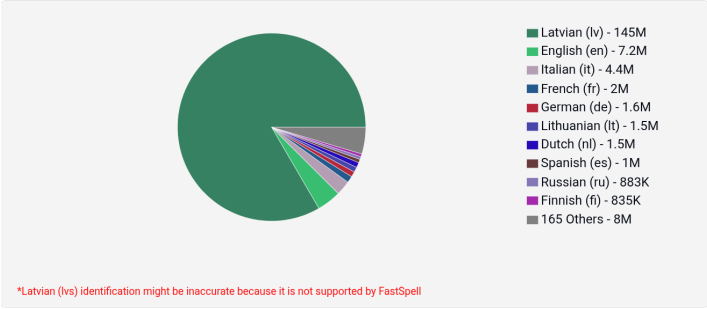


Documents by collection

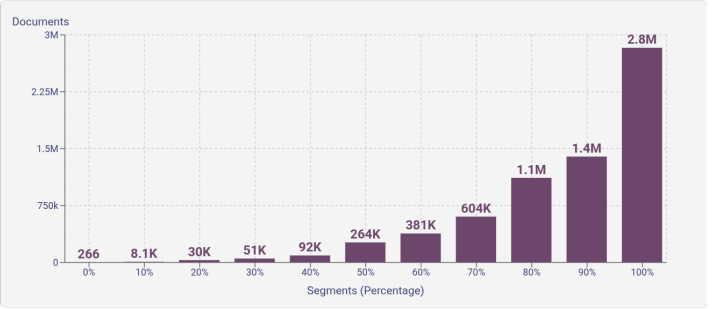


Language Distribution

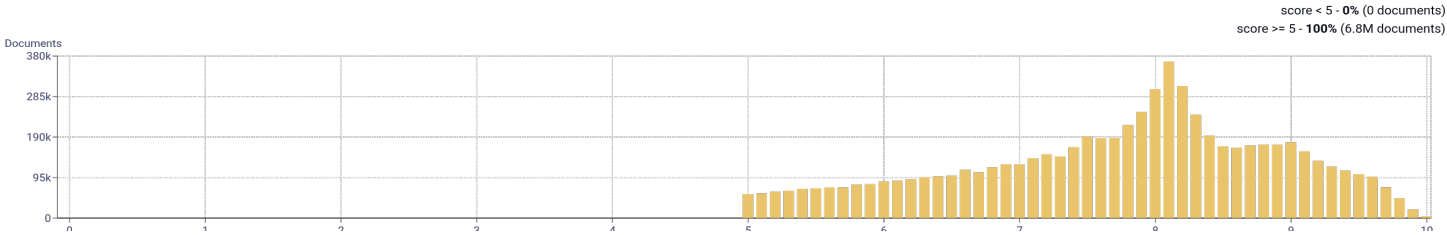
Number of segments in the Latvian (lvs) corpus



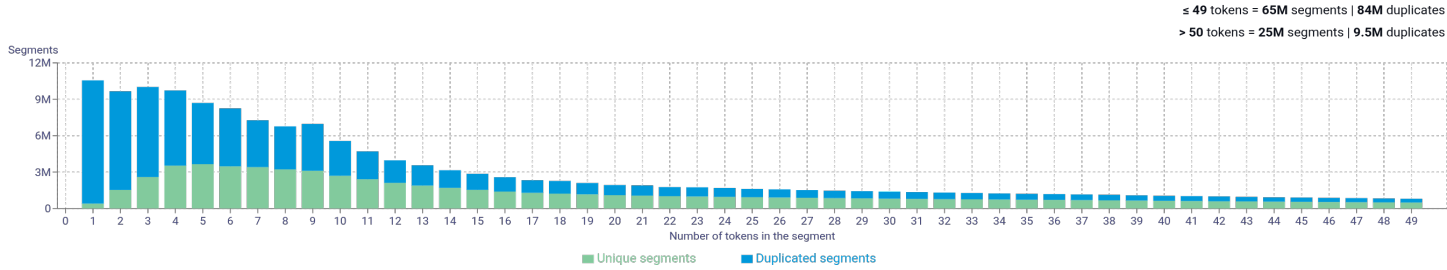
Percentage of segments in Latvian (lvs) inside documents



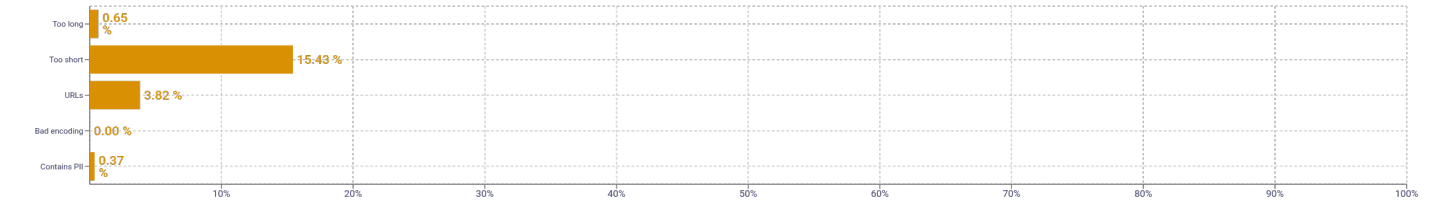
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	<div><div>kas 23005193</div><div>to 15137881</div><div>tas 12521434</div><div>nav 11485252</div><div>latvijas 7400622</div></div>
2	<div><div>šo uzņēmumu 1034889</div><div>šo sadaļu 753720</div><div>labot pirmkodu 681711</div><div>ņemot vērā 595807</div><div>labot šo 585972</div></div>
3	<div><div>labot šo sadaļu 583303</div><div>pēdējās stundas laikā 258318</div><div>stundas laikā šo 258117</div><div>viesnicu ir skatījušas 258109</div><div>laikā šo viesnicu 258109</div></div>
4	<div><div>šo viesnicu ir skatījušas 258109</div><div>stundas laikā šo viesnicu 258109</div><div>pēdējās stundas laikā šo 258109</div><div>kuru nosūtīt darba saiti 184171</div><div>eiropas parlamenta un padomes 70347</div></div>
5	<div><div>pēdējās stundas laikā šo viesnicu 258109</div><div>laikā šo viesnicu ir skatījušas 258109</div><div>vides aizsardzības un reģionālās attīstības 56856</div><div>rakstu pārpublicēšanas noteikumiem lūdzam kontaktēties 43110</div><div>noteikumiem lūdzam kontaktēties ar travelnews.lv 43101</div></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 or 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	Enciclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				