

General overview

Corpus	Date	Language
el_Grek.jsonl.tsv	7/14/2025	Greek (el)

Volumes

Docs	Segments	Unique segments	Tokens	Characters	Size
70,328,862	1,847,935,923	608,316,724 (32.92 %)	50B	281,755,294,923	468.4 GB

Top 10 domains

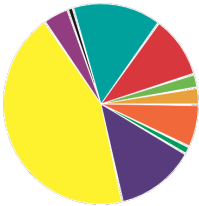
Domain	Docs	% of total
blogspot.gr	8.5M	12.15%
blogspot.com	6.6M	9.32%
wordpress.com	1.5M	2.18%
docplayer.gr	547K	0.78%
liverster.gr	504K	0.72%
inewsgr.com	490K	0.70%
wikipedia.org	446K	0.63%
onsports.gr	314K	0.45%
blogspot.be	291K	0.41%
blogspot.nl	284K	0.40%

Top 10 TLDs

Domain	Docs	% of total
gr	47M	66.34%
com	16M	22.20%
org	1.2M	1.72%
net	828K	1.18%
eu	628K	0.89%
com.cy	619K	0.88%
com.gr	506K	0.72%
be	302K	0.43%
info	254K	0.36%
nl	200K	0.28%

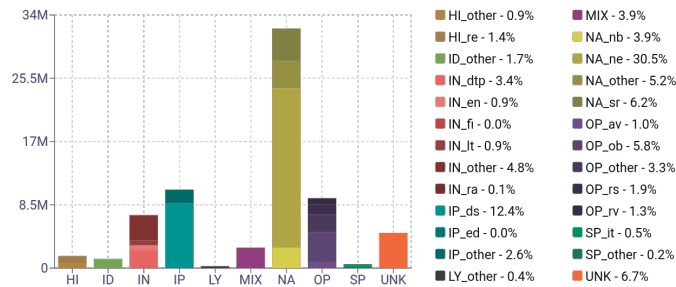
Register labels

ⓘ



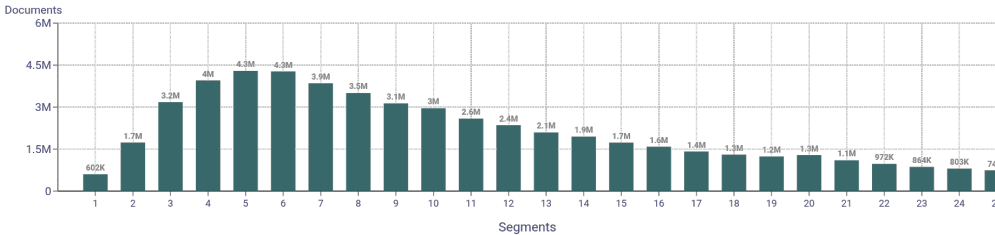
- HI - 2.3%
- ID - 1.7%
- IN - 10.1%
- IP - 15.0%
- LY - 0.4%
- MIX - 3.9%
- NA - 45.7%
- OP - 13.4%
- SP - 0.7%
- UNK - 6.7%

Documents



📊 MT:2.0% | 1.4M Documents

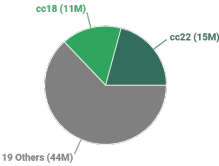
Documents size (in segments) ⓘ



Document collections

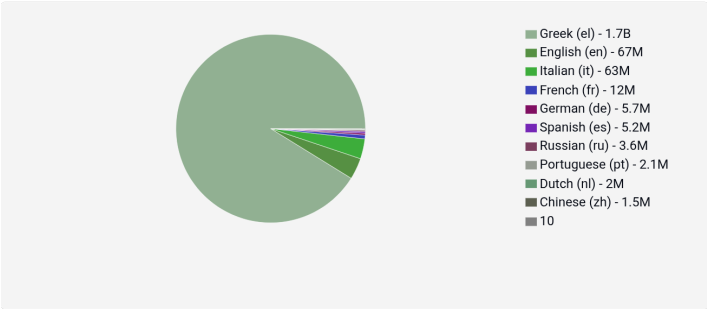
<= 25 segments 100% (53M documents)
> 25 segments 0% (0 documents)

CC = 54.93%
IA = 45.07%

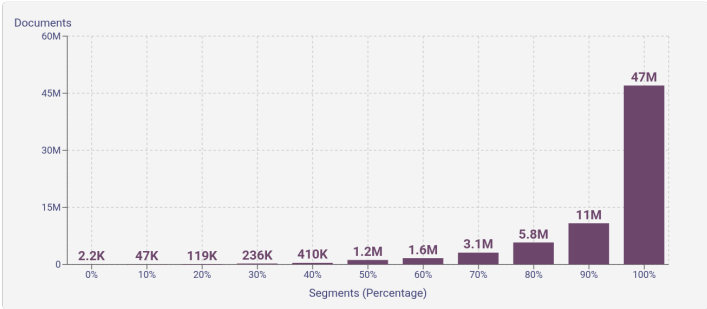


Language Distribution

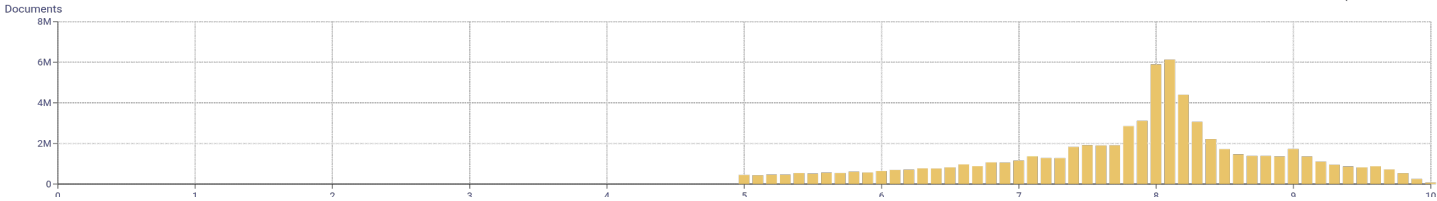
Number of segments in the Greek (el) corpus



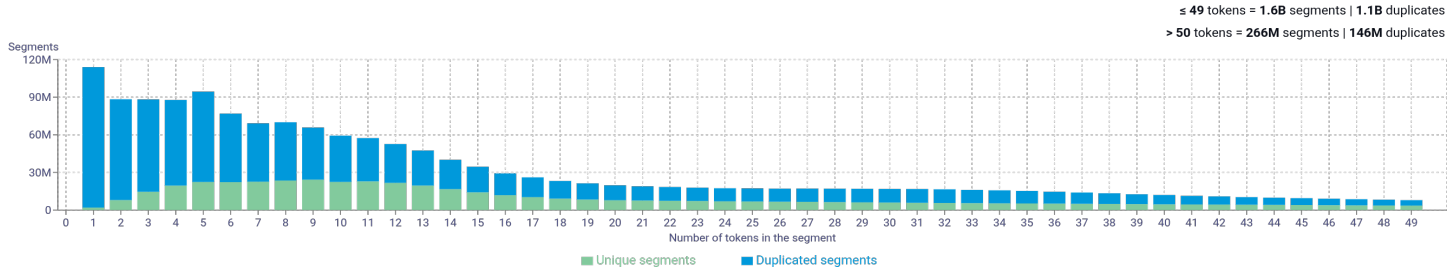
Percentage of segments in Greek (el) inside documents



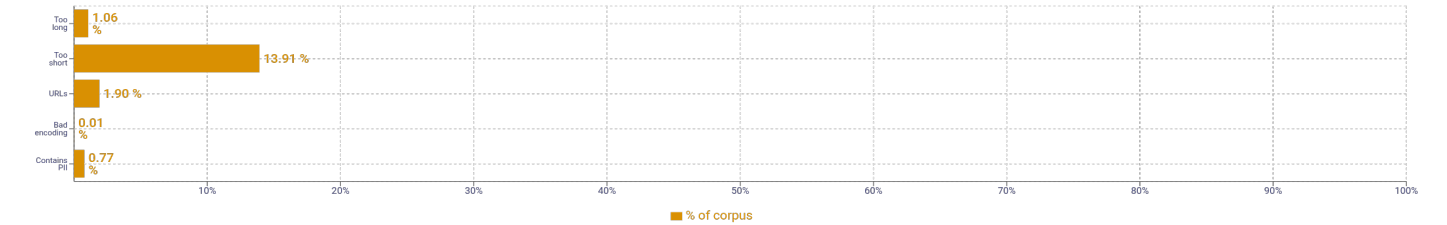
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

Size	n-grams
1	της 709791617 από 500424729 είναι 349978341 τους 269685958 τη 237242716
2	από τη 28406427 από τους 23028260 πριν από 15039016 της χώρας 10230213 διαβάστε περισσότερα 9908523
3	από την άλλη 3109555 ένας από τους 2097234 από την πλευρά 2004434 αυτή τη στιγμή 1772930 από την αρχή 1704974
4	ένα από τα πιο 937480 άδειες χρήσης το παρόν 474848 αξίζει να σημειωθεί ότι 452086 αναρτησα στο διαδικτυο ελληνικη 443930 α π ο ο 422801
5	α π ο ο π 418673 άδειες χρήσης το παρόν εκπαιδευτικό 399666 α π ο φ α 234451 α σ μ α από 204426 google news και μάθετε πρώτοι 185720

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	Enciclopedia article	en	Denominational religious blog or sermon	rs
Sports report	sr	Research article	ra	Advice	av
Narrative blog	nb				