

## General overview

Corpus	Analytics date	Language
ga_1.jsonl.tsv	3/16/2024	Irish (ga)

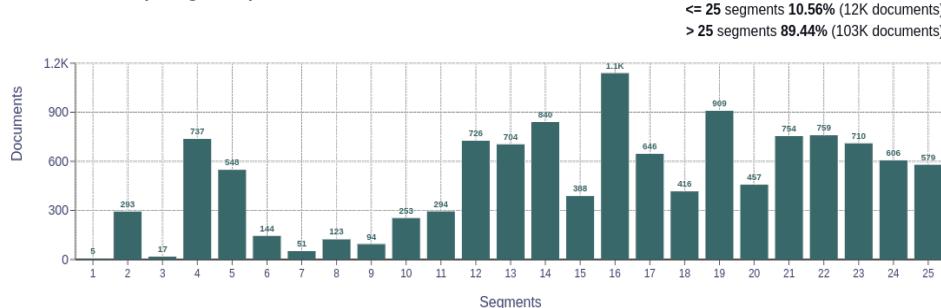
## Volumes

Docs	Segments	Unique segments	Tokens	Size
115,529	13,949,561	16,317 (0.12 %)	152M	810.08 MB

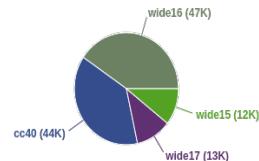
## Type-Token Ratio

Irish (ga)
0.01

## Documents size (in segments)

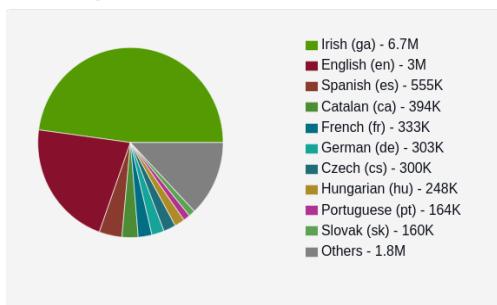


## Documents by collection

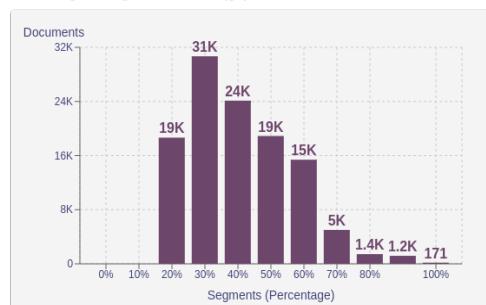


## Language Distribution

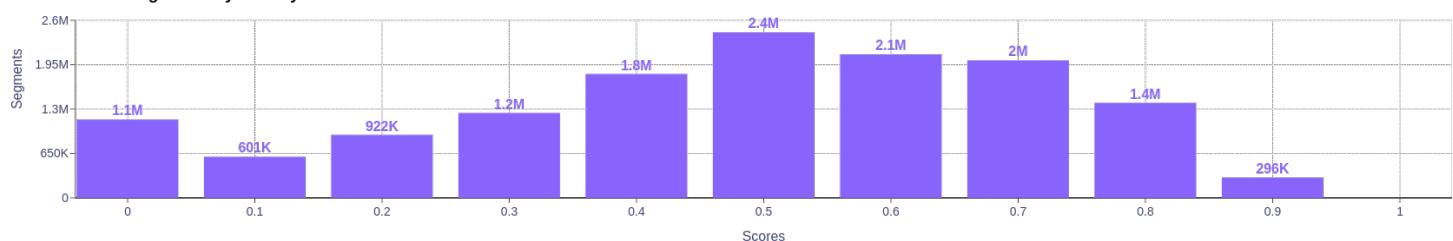
## Number of segments



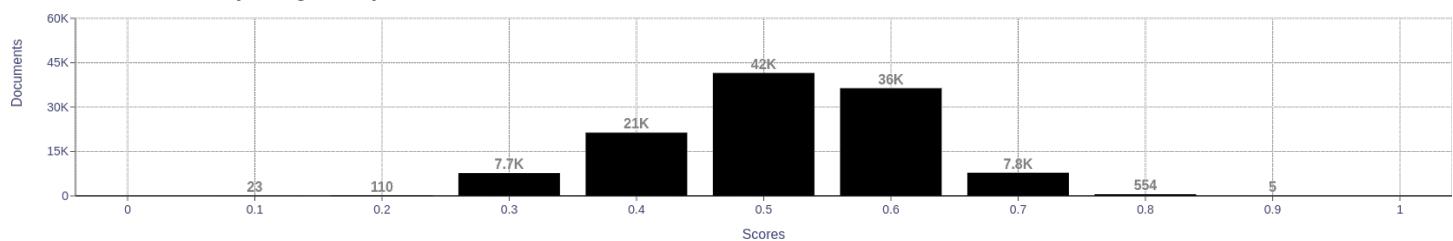
## Percentage of segments in Irish (ga) inside documents



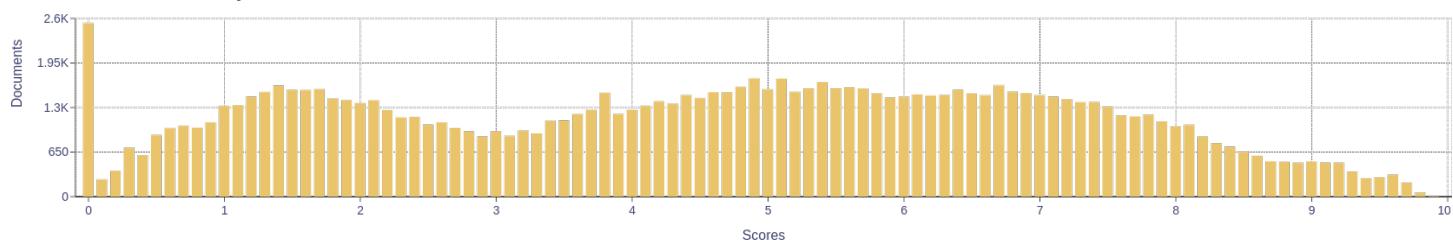
## Distribution of segments by fluency score



## Distribution of documents by average fluency score

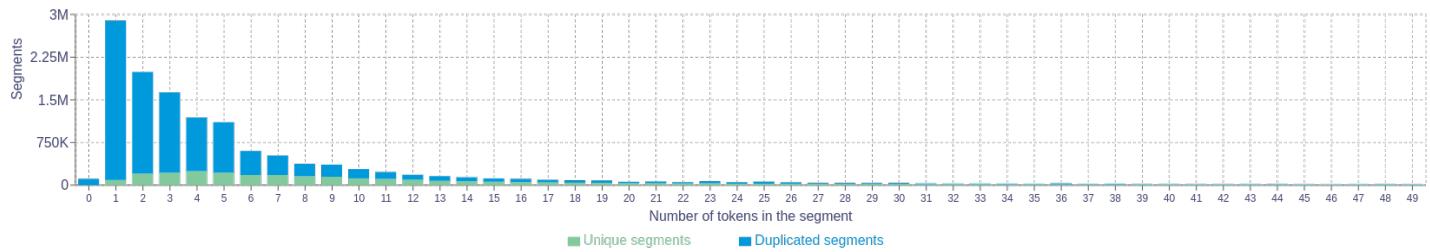


## Distribution of documents by document score

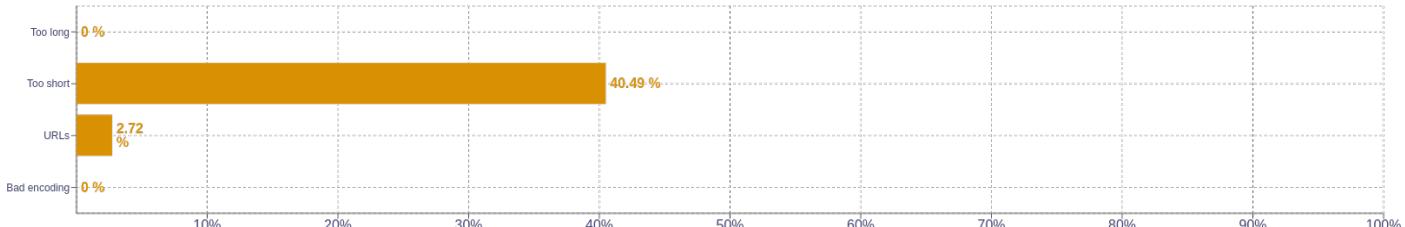


## Segment length distribution by token

<= 49 tokens = 3M segments | 10M duplicates  
 > 50 tokens = 639K segments | 144K duplicates



## Segment noise distribution



## Frequent n-grams

Size	n-grams
1	the   1106634 of   666287 and   525367 to   499552 sin   421573
2	of the   167794 man fmhair   82694 nos m   77263 fdir leat   70355 deireadh fmhair   64487
3	cumann na sagart   268768 lachta an aifrinn   186331 machnamh ar lachta   83362 saor in aisce   65913 baliuchan na scol   38115
4	maidir leis na forbairt   11167 sheoladh lenar liosta riomphoist   11163 gcoimedfai ar an eolas   11161 maith leat go gcoimedfai   11160 thionscadail eile de chuid   11157
5	machnamh ar lachta an aifrinn   81148 macnamh ar lachta an aifrinn   14051 leim go priomhbar an leathanaigh   11406 cuir do sheoladh lenar liosta   11163 mas maith leat go gcoimedfai   11160

## 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>