

General overview

| Corpus             | Analytics date | Language            |
|--------------------|----------------|---------------------|
| hne_Deva.jsonl.tsv | 11/27/2024     | Chhattisgarhi (hne) |

Volumes

| Docs  | Segments | Unique segments     | Tokens | Size    | Characters |
|-------|----------|---------------------|--------|---------|------------|
| 2,806 | 54,999   | 40,771<br>(74.13 %) | 2.5M   | 24.6 MB | 10,541,012 |

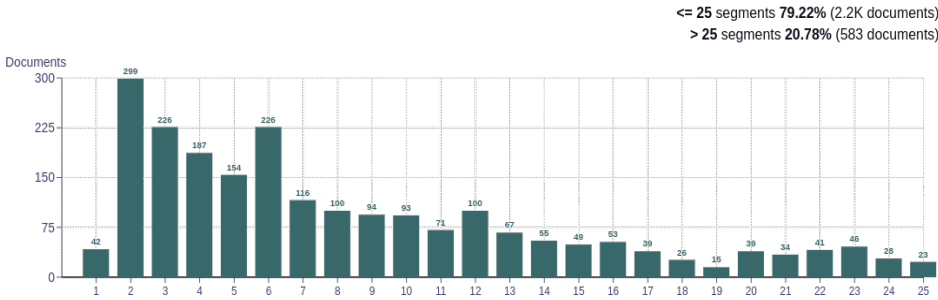
Top 10 domains

| Domain                               | Docs | % of total |
|--------------------------------------|------|------------|
| <a href="#">jayjohar.com</a>         | 560  | 19.96      |
| <a href="#">gurturgoth.com</a>       | 462  | 16.46      |
| <a href="#">blogspot.com</a>         | 313  | 11.15      |
| <a href="#">hi-takedrivers.com</a>   | 179  | 6.38       |
| <a href="#">blogspot.in</a>          | 120  | 4.28       |
| <a href="#">ruralindiaonline.org</a> | 117  | 4.17       |
| <a href="#">news18.com</a>           | 100  | 3.56       |
| <a href="#">vithivichar.com</a>      | 90   | 3.21       |
| <a href="#">biblica.com</a>          | 82   | 2.92       |
| <a href="#">silverweed.club</a>      | 67   | 2.39       |

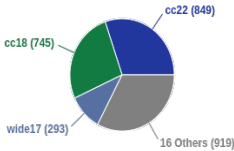
Top 10 TLDs

| Domain | Docs | % of total |
|--------|------|------------|
| com    | 2.1K | 75.37      |
| in     | 185  | 6.59       |
| org    | 169  | 6.02       |
| club   | 159  | 5.67       |
| co.in  | 64   | 2.28       |
| online | 55   | 1.96       |
| ae     | 17   | 0.61       |
| xyz    | 15   | 0.53       |
| net    | 12   | 0.43       |
| jp     | 2    | 0.07       |

Documents size (in segments)

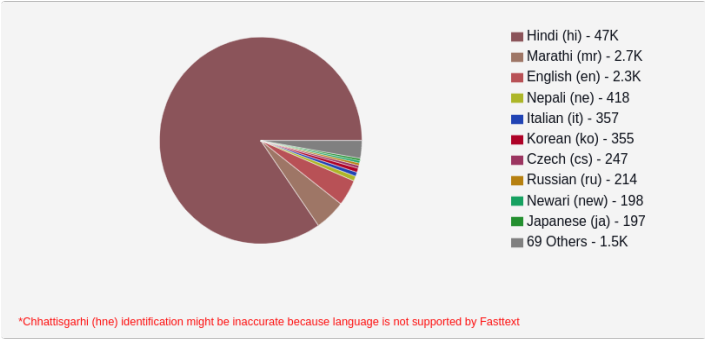


Documents by collection

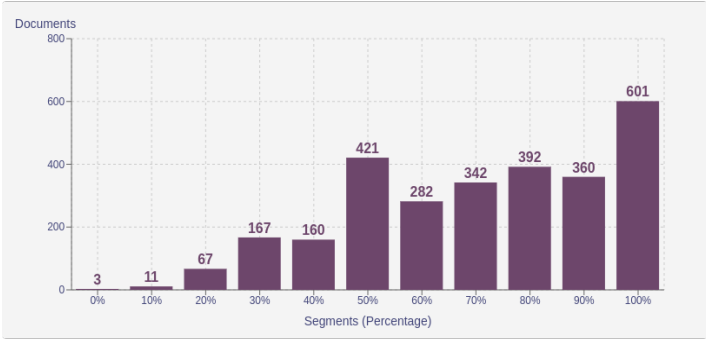


Language Distribution

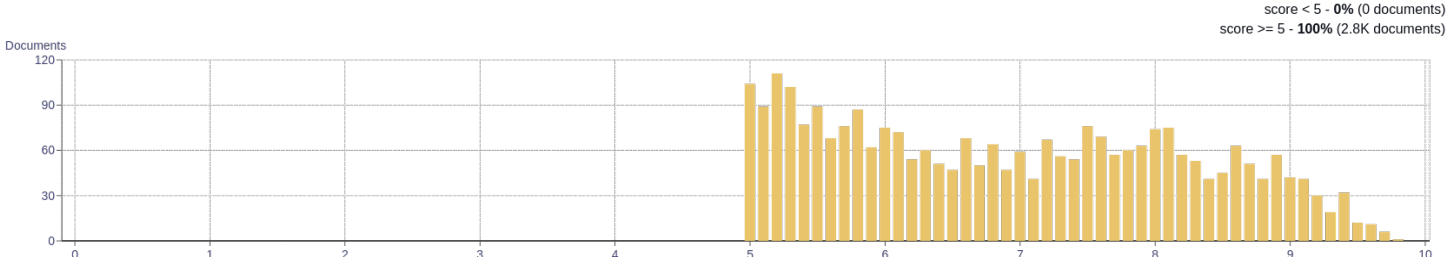
Number of segments



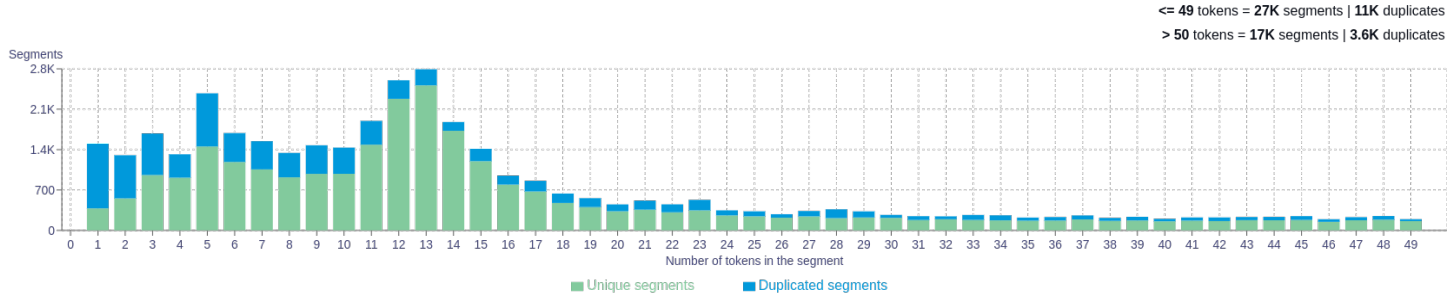
Percentage of segments in Chhattisgarhi (hne) inside documents



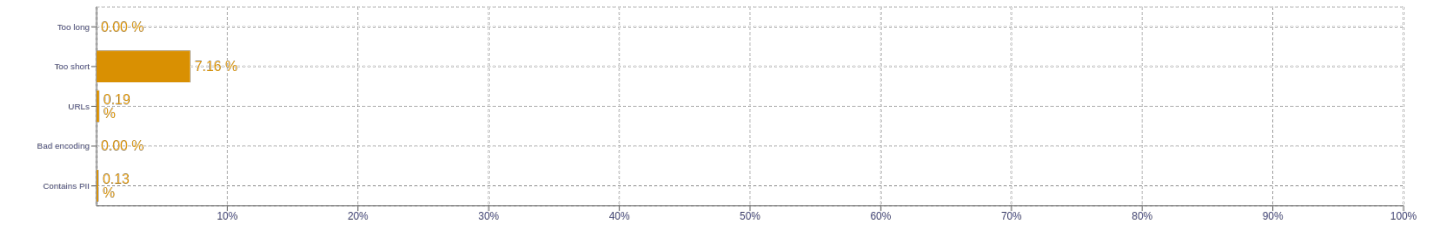
Distribution of documents by document score



Segment length distribution by token



Segment noise distribution



Frequent n-grams

| Size | n-grams  |
|------|--|
| 1    | मन   33111   ह   32170   ले   30393   हे   29654   अउ   20944  |
| 2    | देखे गए   6158   बार देखे   6154   डाइवर डाउनलोड   6154   वो ह   3033   मन ला   2046   |
| 3    | बार देखे गए   6154   keyword keyword keyword   1930   name name name   1629   insert column excel   714   column excel shortcut   693  |
| 4    | keyword keyword keyword keyword   1909   name name name name   1624   insert column excel shortcut   650   gk quiz in hindi   309   knowledge gk quiz in   306                                       |
| 5    | keyword keyword keyword keyword keyword   1888   name name name name name   1619   knowledge gk quiz in hindi   306   general knowledge gk quiz in   306   shortcut shortcut shortcut shortcut   255 |

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.stinfo.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/mbanoni/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>