# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| guj_Gujr.jsonl.tsv | 9/16/2024 | Gujarati (gu) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 1,134,252 | 20,639,718 | 11,424,183 (55.35 %) | 667M | 3,366,421,654 | 7.99 GB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| divyabhaskar.co.in | 83K | 7.35% |
| wordpress.com | 48K | 4.25% |
| news18.com | 44K | 3.92% |
| oneindia.com | 29K | 2.58% |
| sandesh.com | 27K | 2.42% |
| wikipedia.org | 24K | 2.16% |
| gujjurocks.in | 23K | 2.02% |
| chitralekha.com | 18K | 1.58% |
| webdunia.com | 17K | 1.48% |
| vtvgujarati.com | 16K | 1.38% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 727K | 64.10% |
| in | 172K | 15.18% |
| co.in | 101K | 8.92% |
| org | 76K | 6.69% |
| net | 19K | 1.67% |
| news | 4.7K | 0.41% |
| app | 3K | 0.26% |
| online | 2.8K | 0.24% |
| gov.in | 2.4K | 0.21% |
| live | 1.4K | 0.12% |

## Register labels



- HI - 1.7%
- ID - 0.2%
- IN - 9.8%
- IP - 2.1%
- LY - 3.1%
- MIX - 1.2%
- NA - 53.9%
- OP - 11.2%
- SP - 0.2%
- UNK - 16.5%

**MT**:8.4% | 96K Documents



- HI_other - 1.0%
- HI_re - 0.7%
- ID_other - 0.2%
- IN_dtp - 2.9%
- IN_en - 2.2%
- IN_fi - 0.0%
- IN_lt - 0.1%
- IN_other - 4.6%
- IN_ra - 0.0%
- IP_ds - 1.2%
- IP_ed - 0.0%
- IP_other - 1.0%
- LY_other - 3.1%
- MIX - 1.2%
- NA_nb - 2.5%
- NA_ne - 44.7%
- NA_other - 4.7%
- NA_sr - 2.0%
- OP_av - 1.4%
- OP_ob - 1.4%
- OP_other - 3.6%
- OP_rs - 4.3%
- OP_rv - 0.5%
- SP_it - 0.1%
- SP_other - 0.1%
- UNK - 16.5%

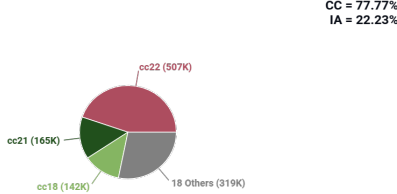## Documents size (in segments)

<= 25 segments **83.56%** (948K documents)
> 25 segments **16.44%** (186K documents)



## Documents by collection

CC = 77.77%
IA = 22.23%



- cc22 (507K)
- cc21 (165K)
- cc18 (142K)
- 18 Others (319K)

## Language Distribution

### Number of segments in the Gujarati (gu) corpus



- Gujarati (gu) - 18M
- English (en) - 1.4M
- Italian (it) - 326K
- Hindi (hi) - 69K
- French (fr) - 58K
- German (de) - 42K
- Sanskrit (sa) - 28K
- Greek (el) - 28K
- Spanish (es) - 18K
- Russian (ru) - 17K
- 161 Others - 260K

### Percentage of segments in Gujarati (gu) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (1.1M documents)

## Segment length distribution by token

Segments

1.2M
900k
600k
300k
0

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49

Number of tokens in the segment

■ Unique segments  ■ Duplicated segments

## Segment noise distribution



Too long — 0.98 %
Too short — 10.93 %
URLs — 1.07 %
Bad encoding — 0.01 %
Contains PII — 0.14 %

10%  20%  30%  40%  50%  60%  70%  80%  90%  100%

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | સાથે \| 2721183    હતી \| 1796456    કરવામાં \| 1275130    કરવા \| 1137108    દ્વારા \| 1131630 |
| 2 | ફેરફાર કરો \| 177866    સાબતે હતી \| 143408    ગોબા ભગે \| 130076    કરવામાં આવ્યું \| 122479    કરવામાં આવ્યો \| 119799 |
| 3 | all rights reserved \| 46726    db corp ltd \| 46044    code of ethics \| 46018    website follows the \| 46017    this website follows \| 46017 |
| 4 | website follows the dnpa \| 46017    this website follows the \| 46017    the dnpa code of \| 46017    follows the dnpa code \| 46017    dnpa code of ethics \| 46017 |
| 5 | website follows the dnpa code \| 46017    this website follows the dnpa \| 46017    the dnpa code of ethics \| 46017    follows the dnpa code of \| 46017    સહિત વધુ સમાચાર વાંચો «ન્યૂઝ18 \| 20648 |

## 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 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>, <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

**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 | | | | |