# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| scn_Latn.jsonl.tsv | 11/27/2024 | Sicilian (scn) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 81,970 | 1,650,375 | 735,503 (44.57 %) | 53M | 250,748,924 | 246.97 MB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| wikipedia.org | 45K | 54.61% |
| vsaduidoma.com | 1.6K | 1.90% |
| apiazzetta.com | 1.5K | 1.88% |
| tempicorsica.com | 1.1K | 1.34% |
| interromania.com | 700 | 0.85% |
| julinse.com | 679 | 0.83% |
| eodishasamachar... | 630 | 0.77% |
| blogspot.com | 541 | 0.66% |
| arritti.corsica | 526 | 0.64% |
| educationbro.com | 485 | 0.59% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| org | 47K | 57.82% |
| com | 24K | 28.97% |
| it | 2.8K | 3.44% |
| corsica | 1.7K | 2.12% |
| net | 1.3K | 1.60% |
| fr | 773 | 0.94% |
| pt | 510 | 0.62% |
| de | 330 | 0.40% |
| eu | 325 | 0.40% |
| zone | 294 | 0.36% |

## Register labels



- HI - 0.6%
- ID - 1.2%
- IN - 55.8%
- IP - 2.8%
- LY - 2.6%
- MIX - 0.2%
- NA - 5.2%
- OP - 2.5%
- SP - 0.1%
- UNK - 29.0%

🤖 **MT**:15.5% | 13K Documents

- HI_other - 0.3%
- HI_re - 0.3%
- ID_other - 1.2%
- IN_dtp - 3.5%
- IN_en - 39.3%
- IN_fi - 0.0%
- IN_lt - 0.1%
- IN_other - 13.0%
- IN_ra - 0.0%
- IP_ds - 2.2%
- IP_ed - 0.0%
- IP_other - 0.6%
- LY_other - 2.6%
- MIX - 0.2%
- NA_nb - 0.7%
- NA_ne - 1.5%
- NA_other - 2.6%
- NA_sr - 0.4%
- OP_av - 0.0%
- OP_ob - 0.3%
- OP_other - 0.5%
- OP_rs - 1.3%
- OP_rv - 0.4%
- SP_it - 0.0%
- SP_other - 0.0%
- UNK - 29.0%

## Documents size (in segments)

**<= 25** segments **81.53%** (67K documents)
**> 25** segments **18.47%** (15K documents)



## Documents by collection

**CC = 65.57%**
**IA = 34.43%**



cc18 (9.1K)  cc22 (20K)  19 Others (53K)

## Language Distribution

### Number of segments in the Sicilian (scn) corpus



- Italian (it) - 1.1M
- Sicilian (scn) - 116K
- English (en) - 104K
- Spanish (es) - 61K
- French (fr) - 56K
- Catalan (ca) - 25K
- Corsican (co) - 15K
- Lombard (lmo) - 14K
- Romanian (ro) - 13K
- German (de) - 13K
- 161 Others - 150K

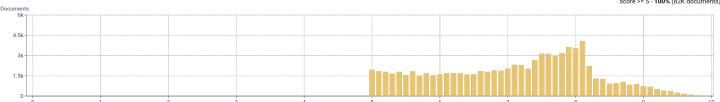### Percentage of segments in Sicilian (scn) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
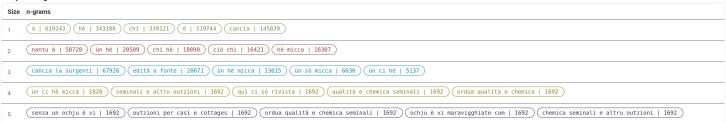score >= 5 - **100%** (82K documents)

## Segment length distribution by token

Segments / Number of tokens in the segment

■ Unique segments  ■ Duplicated segments

## Segment noise distribution



| Category | Value |
| --- | --- |
| Too long | 1.06 % |
| Too short | 11.57 % |
| URLs | 0.86 % |
| Bad encoding | 0.07 % |
| Contains PII | 0.07 % |

## Frequent n-grams

| Size | n-grams |
| --- | --- |
| 1 | à \| 619243   hè \| 343180   chì \| 330121   d \| 219744   cancia \| 145839 |
| 2 | nantu à \| 58728   ùn hè \| 20509   chì hè \| 18890   ciò chì \| 16421   hè micca \| 16307 |
| 3 | cancia la surgenti \| 67926   edità a fonte \| 20671   ùn hè micca \| 13615   ùn sò micca \| 6030   ùn ci hè \| 5137 |
| 4 | ùn ci hè micca \| 1820   seminali e altru outzioni \| 1692   quì ci sò rivista \| 1692   qualità e chemica seminali \| 1692   ordua qualità e chemica \| 1692 |
| 5 | senza un ochju è vi \| 1692   outzioni per casi e cottages \| 1692   ordua qualità e chemica seminali \| 1692   ochju è vi maravigghiate cum \| 1692   chemica seminali e altru outzioni \| 1692 |

## 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 | | | | |