# HPLT Analytics report


HPLTAnalytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| hplt-v3-jpn_Jpan-SAMPLED | 10/3/2025 | Japanese |

## Volumes

| Docs | Segments | Unique segments | Duplication ratio | Tokens | Characters | Size |
|---|---|---|---|---|---|---|
| 667,404,438 | 35,767,735,744 | 33,130,761 (0.09 %) | 99.91% | 818B | 1,461,409,987,151 | 3.63 TB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| fc2.com | 32M | 4.76% |
| hatenablog.com | 11M | 1.62% |
| cocolog-nifty.com | 10M | 1.50% |
| exblog.jp | 7.1M | 1.07% |
| fortune-uranai.... | 5.4M | 0.80% |
| yahoo.co.jp | 4.4M | 0.66% |
| seesaa.net | 3.8M | 0.57% |
| goo.ne.jp | 3.6M | 0.54% |
| tabelog.com | 3.1M | 0.47% |
| biglobe.ne.jp | 2.6M | 0.39% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 304M | 45.58% |
| jp | 146M | 21.86% |
| net | 50M | 7.53% |
| co.jp | 46M | 6.83% |
| info | 16M | 2.45% |
| org | 16M | 2.33% |
| ne.jp | 14M | 2.03% |
| biz | 6M | 0.90% |
| or.jp | 4.8M | 0.72% |
| xyz | 4.4M | 0.66% |

## Register labels



- HI - 2.0%
- ID - 4.2%
- IN - 14.4%
- IP - 32.0%
- LY - 0.1%
- MIX - 5.9%
- NA - 21.4%
- OP - 10.3%
- SP - 1.0%
- UNK - 8.7%

- HI_other - 1.6%
- HI_re - 0.4%
- ID_other - 4.2%
- IN_dtp - 3.6%
- IN_en - 1.0%
- IN_fi - 0.0%
- IN_lt - 0.7%
- IN_other - 9.1%
- IN_ra - 0.0%
- IP_ds - 28.6%
- IP_ed - 0.0%
- IP_other - 3.3%
- LY_other - 0.1%

- MIX - 5.9%
- NA_nb - 14.6%
- NA_ne - 3.4%
- NA_other - 2.3%
- NA_sr - 1.1%
- OP_av - 2.3%
- OP_ob - 1.1%
- OP_other - 2.0%
- OP_rs - 0.2%
- OP_rv - 4.8%
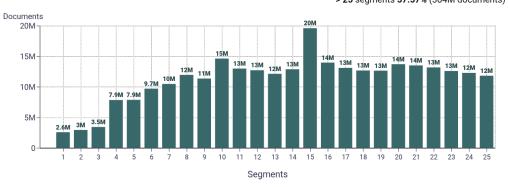- SP_it - 0.6%
- SP_other - 0.3%
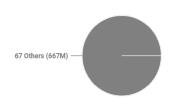- UNK - 8.7%

♛ **MT**:3.3% | 22M Documents

## Documents size (in segments) ⓘ

≤ **25** segments **42.43%** (283M documents)
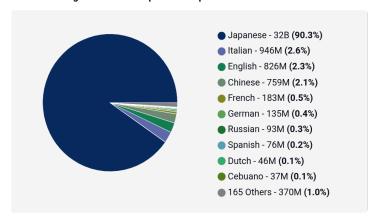> **25** segments **57.57%** (384M documents)
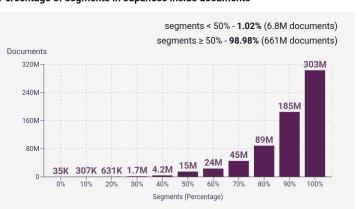


## Document collections

**CC = 84.89%**
**IA = 15.11%**



67 Others (667M)

## Language Distribution

### Number of segments in the Japanese corpus

- Japanese - 32B **(90.3%)**
- Italian - 946M **(2.6%)**
- English - 826M **(2.3%)**
- Chinese - 759M **(2.1%)**
- French - 183M **(0.5%)**
- German - 135M **(0.4%)**
- Russian - 93M **(0.3%)**
- Spanish - 76M **(0.2%)**
- Dutch - 46M **(0.1%)**
- Cebuano - 37M **(0.1%)**
- 165 Others - 370M **(1.0%)**

### Percentage of segments in Japanese inside documents

segments < 50% - **1.02%** (6.8M documents)
segments ≥ 50% - **98.98%** (661M documents)

Documents

| 0% | 10% | 20% | 30% | 40% | 50% | 60% | 70% | 80% | 90% | 100% |
|---|---|---|---|---|---|---|---|---|---|---|
| 35K | 307K | 631K | 1.7M | 4.2M | 15M | 24M | 45M | 89M | 185M | 303M |

Segments (Percentage)

## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score ≥ 5 - **100%** (667M documents)

Documents

## Segment length distribution by token

≤ **49** tokens = **32B** segments | **18B** duplicates
> **50** tokens = **3.4B** segments | **1.2B** duplicates

Segments

Number of tokens in the segment

## Segment noise distribution

| | % |
|---|---|
| Too long | 0.08% |
| Too short | 3.78% |
| URLs | 1.01% |
| Bad encoding | 0.00% |
| Contains PII | 0.06% |

■ % of corpus

## Frequent n-grams

| SIZE | N-GRAMS | | | | |
|---|---|---|---|---|---|
| 1 | て \| 21,012,489 | で \| 16,793,720 | た \| 13,473,279 | ます \| 8,191,653 | な \| 6,729,083 |
| 2 | て い \| 2,862,824 | て いる \| 2,623,445 | まし た \| 2,427,244 | い ます \| 1,506,673 | され \| 1,505,429 |
| 3 | て います \| 1,409,695 | され て \| 642,107 | て いた \| 639,080 | され た \| 417,165 | れ てい \| 409,248 |
| 4 | され て いる \| 240,605 | され てい \| 236,053 | れ て います \| 235,438 | て いました \| 191,964 | こと が できます \| 185,241 |
| 5 | かも しれ ません \| 176,734 | され て います \| 135,902 | では あり ません \| 105,411 | する こと が できます \| 63,342 | では ない でしょう か \| 51,223 |

# 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 | | | | |