# HPLT Analytics report

**HPLT**Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| lao_Laoo.jsonl.tsv | 9/23/2024 | Lao (lo) |

### Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 29,504 | 319,953 | 221,388 (69.19 %) | 22M | 84,390,143 | 221.5 MB |

### Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| vnanet.vn | 5.2K | 17.55% |
| vovworld.vn | 2K | 6.66% |
| na.gov.la | 1.3K | 4.29% |
| lnr.org.la | 1.1K | 3.77% |
| lsr.com.la | 981 | 3.32% |
| thoidai.com.vn | 941 | 3.19% |
| wikipedia.org | 709 | 2.40% |
| tapchicongsan.o... | 534 | 1.81% |
| cri.cn | 473 | 1.60% |
| nuol.edu.la | 358 | 1.21% |

### Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| vn | 7.5K | 25.38% |
| com | 7.4K | 25.08% |
| gov.la | 5.6K | 19.11% |
| org | 2K | 6.67% |
| org.la | 1.5K | 5.04% |
| com.la | 1.4K | 4.66% |
| com.vn | 941 | 3.19% |
| edu.la | 620 | 2.10% |
| org.vn | 534 | 1.81% |
| cn | 487 | 1.65% |

## Register labels



- HI - 0.0%
- ID - 0.0%
- IN - 8.2%
- IP - 0.9%
- LY - 0.0%
- MIX - 0.3%
- NA - 72.3%
- OP - 0.5%
- SP - 0.3%
- UNK - 17.6%

**MT**:14.7% | 4.3K Documents

Documents

- HI_other - 0.0%
- HI_re - 0.0%
- ID_other - 0.0%
- IN_dtp - 1.8%
- IN_en - 1.8%
- IN_fi - 0.0%
- IN_lt - 0.8%
- IN_other - 3.7%
- IN_ra - 0.1%
- IP_ds - 0.6%
- IP_ed - 0.0%
- IP_other - 0.3%
- LY_other - 0.0%
- MIX - 0.3%
- NA_nb - 0.1%
- NA_ne - 68.9%
- NA_other - 3.3%
- NA_sr - 0.1%
- OP_av - 0.0%
- OP_ob - 0.0%
- OP_other - 0.1%
- OP_rs - 0.4%
- OP_rv - 0.0%
- SP_it - 0.0%
- SP_other - 0.3%
- UNK - 17.6%

## Documents size (in segments)

**<= 25** segments **92.57%** (27K documents)
**> 25** segments **7.43%** (2.2K documents)



## Documents by collection

**CC = 83.06%**
**IA = 16.94%**



cc22 (11K)
cc18 (7.4K)
cc21 (4K)
17 Others (6.8K)

## Language Distribution

### Number of segments in the Lao (lo) corpus



- Lao (lo) - 297K
- English (en) - 10K
- Italian (it) - 3.5K
- French (fr) - 1K
- German (de) - 749
- Thai (th) - 656
- Greek (el) - 565
- Slovenian (sl) - 540
- Ukrainian (uk) - 525
- Catalan (ca) - 473
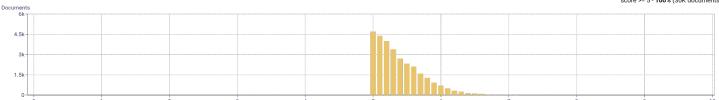- 99 Others - 4.4K

### Percentage of segments in Lao (lo) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (30K documents)

## Segment length distribution by token

≤ 49 tokens = **107K** segments | **70K** duplicates
> 50 tokens = **142K** segments | **28K** duplicates

Segments

14k

10.5k

7k

3.5k

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 | 3.15 % |
| Too short | 19.24 % |
| URLs | 0.65 % |
| Bad encoding | 0.01 % |
| Contains PII | 0.08 % |

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

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | ໆ \| 156290   ທ່ານ \| 129958   ລາວ \| 92003   ແນ່ນ \| 84606   ຜັກ \| 74319 |
| 2 | ກອງ ປະຊຸມ \| 53088   ສ້ ຳ \| 42316   ສະພາ ແຫ່ງຊາດ \| 25340   ບໍ ໆ \| 24536   ໆ ລົງ \| 20624 |
| 3 | ສ້ ຳ ລົງ \| 20609   ສ້ ຳ ດັບ \| 8441   ບໍ ໆ ໃຊ້ \| 8161   ຄະນະ ບໍລິຫານ ງານ \| 7607   ຈໍ ໆ ມວນ \| 7230 |
| 4 | ຕຽງ ວ່ງຈ ຊວນ ຜູກ \| 3904   ຄະນະ ບໍລິຫານ ງານ ສູນກາງ \| 3721   ແນວ ລາວ ຮັກ ຊາດ \| 3676   ບໍລິຫານ ງານ ສູນກາງ ພັກ \| 3156   ກິມ ການເມືອງ ສູນກາງ ພັກ \| 2971 |
| 5 | ຄະນະ ບໍລິຫານ ງານ ສູນກາງ ພັກ \| 3083   ກອງ ປະຊຸມ ໃຫຍ່ ຄັ້ງ ທີ \| 1614   ສູນກາງ ແນວ ລາວ ຮັກ ຊາດ \| 1590   ຄະນະ ໂຄສະນາ ອົບຮົມ ສູນກາງ ພັກ \| 1588   ກຳມະການ ກິມ ການເມືອງ ສູນກາງ ພັກ \| 1328 |

## 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 | | | | |