# HPLT Analytics report

**HPLT** Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| hau_Latn.jsonl.tsv | 9/20/2024 | Hausa (ha) |

## Volumes

| Docs | Segments | Unique segments | Tokens | Characters | Size |
|---|---|---|---|---|---|
| 315,870 | 5,688,420 | 3,787,469 (66.58 %) | 180M | 848,139,069 | 820.35 MB |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| voahausa.com | 48K | 15.16% |
| legit.ng | 33K | 10.45% |
| leadership.ng | 23K | 7.27% |
| premiumtimesng.com | 18K | 5.74% |
| rfi.fr | 8.5K | 2.68% |
| bbc.com | 8.4K | 2.65% |
| cri.cn | 6.4K | 2.01% |
| dw.com | 5K | 1.60% |
| isyaku.com | 4.5K | 1.43% |
| wondershare.com | 4.5K | 1.43% |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 190K | 60.15% |
| ng | 61K | 19.34% |
| com.ng | 19K | 6.14% |
| org | 11K | 3.46% |
| fr | 9.3K | 2.94% |
| cn | 6.5K | 2.07% |
| net | 4K | 1.26% |
| ir | 3.6K | 1.14% |
| zone | 1.8K | 0.56% |
| co.uk | 1.3K | 0.40% |

## Register labels



- HI - 0.6%
- ID - 0.4%
- IN - 3.2%
- IP - 1.6%
- LY - 0.2%
- MIX - 0.4%
- NA - 69.1%
- OP - 5.3%
- SP - 0.6%
- UNK - 18.8%

**MT**:16.3% | 52K Documents

Documents

- HI_other - 0.5%
- HI_re - 0.1%
- ID_other - 0.4%
- IN_dtp - 0.9%
- IN_en - 0.8%
- IN_fi - 0.0%
- IN_lt - 0.1%
- IN_other - 1.5%
- IN_ra - 0.0%
- IP_ds - 1.2%
- IP_ed - 0.0%
- IP_other - 0.3%
- LY_other - 0.2%
- MIX - 0.4%
- NA_nb - 0.2%
- NA_ne - 63.5%
- NA_other - 2.9%
- NA_sr - 2.5%
- OP_av - 0.2%
- OP_ob - 0.3%
- OP_other - 1.5%
- OP_rs - 3.1%
- OP_rv - 0.1%
- SP_it - 0.2%
- SP_other - 0.3%
- UNK - 18.8%

## Documents size (in segments)

<= 25 segments **86.05%** (272K documents)
> 25 segments **13.95%** (44K documents)



## Documents by collection

CC = 90.27%
IA = 9.73%



cc22 (133K)
cc21 (74K)
cc18 (50K)
18 Others (59K)

## Language Distribution

### Number of segments in the Hausa (ha) corpus



- English (en) - 2.2M
- Filipino (tl) - 605K
- Swahili (sw) - 408K
- Italian (it) - 406K
- Indonesian (id) - 316K
- German (de) - 267K
- Esperanto (eo) - 177K
- Spanish (es) - 128K
- Ido (io) - 91K
- Portuguese (pt) - 85K
- 161 Others - 1M

*Hausa (ha) identification might be inaccurate because it is not supported by FastSpell
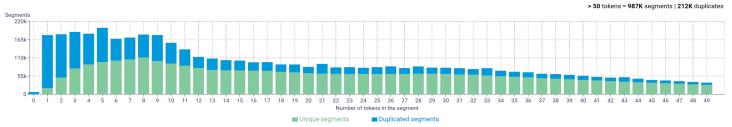
### Percentage of segments in Hausa (ha) inside documents



## Distribution of documents by document score

score < 5 - **0%** (0 documents)
score >= 5 - **100%** (316K documents)

## Segment length distribution by token

≤ 49 tokens = **3M** segments | **1.7M** duplicates
> 50 tokens = **987K** segments | **212K** duplicates



Legend: ■ Unique segments  ■ Duplicated segments

X-axis: Number of tokens in the segment
Y-axis: Segments

## Segment noise distribution



| Category | Value |
|---|---|
| Too long | 0.56 % |
| Too short | 9.87 % |
| URLs | 3.39 % |
| Bad encoding | 0.03 % |
| Contains PII | 0.47 % |

## Frequent n-grams

| Size | n-grams |
|---|---|
| 1 | tare \| 443713   haka \| 427231   kasar \| 414901   allah \| 355204   jihar \| 339309 |
| 2 | ci gaba \| 139703   shugaban kasa \| 79477   boko haram \| 41826   kayan aiki \| 40359   gwamnan jihar \| 39630 |
| 3 | majalisar dinkin duniya \| 17795   kasa da kasa \| 15999   dandalin sada zumunta \| 13811   kasa muhammadu buhari \| 13176   shugaban kasa muhammadu \| 13173 |
| 4 | shugaban kasa muhammadu buhari \| 12988   wayar ku ta hannu \| 12395   shawara ko bukatar bamu \| 11826   shafukanmu na dandalin sada \| 11675   latsa wannan domin samun \| 9175 |
| 5 | shawara ko bukatar bamu labari \| 11825   shafukanmu na dandalin sada zumunta \| 11675   latsa wannan domin samun sabuwar \| 8055   sabuwar manhajar labarai ta legit \| 7418   ng a shafinka na facebook \| 5274 |

## 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 | Encyclopedia article | en | Denominational religious blog or sermon | rs |
| Sports report | sr | Research article | ra | Advice | av |
| Narrative blog | nb | | | | |