# HPLT Analytics report

**HPLT** Analytics

## General overview

| Corpus | Date | Language |
|---|---|---|
| hplt-v3-arb_Arab | 10/3/2025 | Standard Arabic |

## Volumes

| Docs | Segments | Unique segments | Duplication ratio | Tokens | Characters | Size |
|---|---|---|---|---|---|---|
| 50,071,127 | 755,927,954 | 564,065,401 (74.62 %) | 25.38% | 28B | 146,475,244,972 | 246 GB |

## Top 10 domains

| Domain | Docs | % of total |
|---|---|---|
| aljazeera.net | 406K | 0.81% |
| alittihad.ae | 360K | 0.72% |
| rt.com | 246K | 0.49% |
| netlify.app | 241K | 0.48% |
| web.app | 227K | 0.45% |
| blogspot.com | 196K | 0.39% |
| albayan.ae | 194K | 0.39% |
| alaraby.co.uk | 168K | 0.34% |
| sputniknews.com | 167K | 0.33% |
| alarabiya.net | 163K | 0.33% |

## Top 10 TLDs

| Domain | Docs | % of total |
|---|---|---|
| com | 28M | 55.39% |
| net | 7.4M | 14.74% |
| org | 2.7M | 5.30% |
| ae | 1.1M | 2.17% |
| ma | 903K | 1.80% |
| ps | 664K | 1.33% |
| news | 655K | 1.31% |
| info | 651K | 1.30% |
| app | 479K | 0.96% |
| tv | 421K | 0.84% |

## Register labels



- HI - 0.8%
- ID - 0.3%
- IN - 8.0%
- IP - 4.6%
- LY - 0.0%
- MIX - 1.1%
- NA - 70.0%
- OP - 4.1%
- SP - 0.4%
- UNK - 10.8%

- HI_other - 0.8%
- HI_re - 0.0%
- ID_other - 0.3%
- IN_dtp - 2.2%
- IN_en - 0.2%
- IN_fi - 0.0%
- IN_lt - 0.5%
- IN_other - 4.9%
- IN_ra - 0.1%
- IP_ds - 3.5%
- IP_other - 1.1%
- LY_other - 0.0%

- MIX - 1.1%
- NA_nb - 0.1%
- NA_ne - 66.5%
- NA_other - 2.2%
- NA_sr - 1.1%
- OP_av - 0.4%
- OP_ob - 2.0%
- OP_other - 1.4%
- OP_rs - 0.2%
- OP_rv - 0.1%
- SP_it - 0.2%
- SP_other - 0.2%
- UNK - 10.8%

🤖 **MT**:5.7% | 2.9M Documents

## Documents size (in segments) ⓘ

≤ 25 segments **87.59%** (44M documents)
> 25 segments **12.41%** (6.2M documents)



## Document collections

**CC = 89.55%**
**IA = 10.45%**



67 Others (50M)
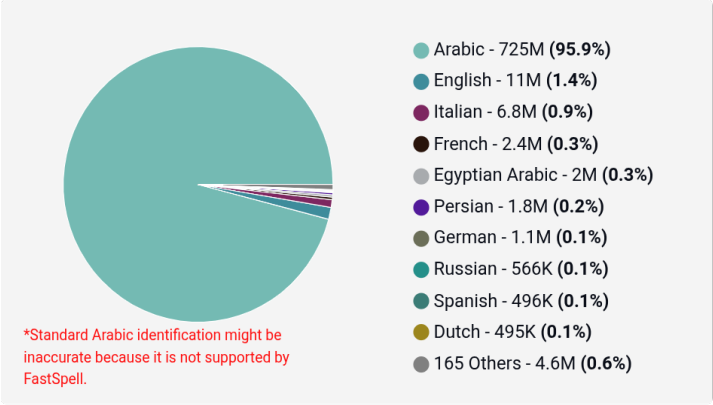
## Language Distribution

### Number of segments in the Standard Arabic corpus

- Arabic - 725M **(95.9%)**
- English - 11M **(1.4%)**
- Italian - 6.8M **(0.9%)**
- French - 2.4M **(0.3%)**
- Egyptian Arabic - 2M **(0.3%)**
- Persian - 1.8M **(0.2%)**
- German - 1.1M **(0.1%)**
- Russian - 566K **(0.1%)**
- Spanish - 496K **(0.1%)**
- Dutch - 495K **(0.1%)**
- 165 Others - 4.6M **(0.6%)**

*Standard Arabic identification might be inaccurate because it is not supported by FastSpell.

### Percentage of segments in Standard Arabic inside documents

segments < 50% - **7.91%** (4M documents)
segments ≥ 50% - **92.09%** (46M documents)

Documents

| Segments (Percentage) | Documents |
|---|---|
| 0% | 6.4K |
| 10% | 134K |
| 20% | 594K |
| 30% | 1.2M |
| 40% | 2.1M |
| 50% | 4.5M |
| 60% | 6M |
| 70% | 7.6M |
| 80% | 11M |
| 90% | 6.2M |
| 100% | 11M |

### Distribution of documents by document score

score < 5 - **0%** (0 documents)
score ≥ 5 - **100%** (50M documents)

Documents

### Segment length distribution by token

**≤ 49** tokens = **582M** segments | **176M** duplicates
**> 50** tokens = **174M** segments | **16M** duplicates

Segments

### Segment noise distribution

| | % of corpus |
|---|---|
| Too long | 0.83% |
| Too short | 8.57% |
| URLs | 0.72% |
| Bad encoding | 0.00% |
| Contains PII | 0.11% |

■ % of corpus

## Frequent n-grams

| SIZE | N-GRAMS |
|------|---------|
| 1 | خلال \| 64,431,878   اليوم \| 36,834,705   يمكن \| 33,624,991   عام \| 31,858,950   العام \| 31,534,283 |
| 2 | الولايات المتحدة \| 11,770,444   الخيارات الثنائية \| 7,491,096   أكمل القراءة \| 5,819,230   الأمم المتحدة \| 5,179,880   عبر الإنترنت \| 4,888,034 |
| 3 | المملكة العربية السعودية \| 3,845,615   تداول العملات الأجنبية \| 1,944,553   الإمارات العربية المتحدة \| 1,857,390   الولايات المتحدة الأمريكية \| 1,724,755   تداول الخيارات الثنائية \| 1,667,447 |
| 4 | دولة الإمارات العربية المتحدة \| 678,105   صلى الله عليه وسلم \| 479,909   بن زايد آل نهيان \| 421,279   الأمين العام للأمم المتحدة \| 417,975   الرئيس عبد الفتاح السيسي \| 373,271 |
| 5 | محمد بن راشد آل مكتوم \| 308,125   الثنائية هي نوع من الأدوات \| 291,635   الأدوات المالية التي تسمح للمستثمر \| 290,403   المالية التي تسمح للمستثمر تحقيق \| 290,062   تسمح للمستثمر تحقيق مكاسب مالية \| 289,340 |

## 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 | | | | |