

# Leveraging Corpus Metadata to Detect Template-based Translation: An Exploratory Case Study of the Egyptian Arabic Wikipedia Edition

Saied Alshahrani<sup>1</sup> Hesham Haroon<sup>2</sup> Ali Elfilali<sup>3</sup> Mariama Njie<sup>4</sup> Jeanna Matthews<sup>1</sup>

<sup>1</sup>Clarkson University, USA <sup>2</sup>Sesame Labs, Egypt <sup>3</sup>Cadi Ayyad University, Morocco <sup>4</sup>M&T Bank, USA  
{saied, jnm}@clarkson.edu, hesham@smsm.ai, a.elfilali9805@uca.ac.ma, mnjie@mtb.com

## Abstract

Wikipedia articles (content pages) are commonly used corpora in Natural Language Processing (NLP) research, especially in low-resource languages other than English. Yet, a few research studies have studied the three Arabic Wikipedia editions, Arabic Wikipedia (AR), Egyptian Arabic Wikipedia (ARZ), and Moroccan Arabic Wikipedia (ARY), and documented issues in the Egyptian Arabic Wikipedia edition regarding the massive automatic creation of its articles using template-based translation from English to Arabic without human involvement, overwhelming the Egyptian Arabic Wikipedia with articles that do not only have low-quality content but also with articles that do not represent the Egyptian people, their culture, and their dialect. In this paper, we aim to mitigate the problem of template translation that occurred in the Egyptian Arabic Wikipedia by identifying these template-translated articles and their characteristics through exploratory analysis and building automatic detection systems. We first explore the content of the three Arabic Wikipedia editions in terms of density, quality, and human contributions and utilize the resulting insights to build multivariate machine learning classifiers leveraging articles' metadata to detect the template-translated articles automatically. We then publicly deploy and host the best-performing classifier, XGBoost, as an online application called EGYPTIAN WIKIPEDIA SCANNER\* and release the extracted, filtered, and labeled datasets to the research community to benefit from our datasets and the online, web-based detection system.

**Keywords:** Arabic, Egyptian, Moroccan, Wikipedia, Template Translation, Multivariate Classification

## 1. Introduction

Wikipedia articles are widely used as pre-training datasets for many Natural Language Processing (NLP) tasks like language modeling (language models) and word representation (word embedding models) tasks, especially for low-resource languages like Arabic, due to its large collection of multilingual content and its vast array of metadata that can be quantified and compared (Mittermeier et al., 2021). However, not all Wikipedia articles are organically produced by native speakers of those languages; while humans have naturally generated some articles in those languages, many others have been automatically generated using bots or automatically translated from high-resourced languages like English without human revision using off-the-shelf automatic translation tools like Google Translate<sup>1</sup> (Hautasaari, 2013; Nisioi et al., 2016; Baker, 2022; Alshahrani et al., 2022; Johnson and Lescak, 2022; Bhattacharjee and Giner, 2022; Wikipedia Foundation, 2022).

A few researchers have addressed this issue and highlighted its implications for NLP systems and tasks. For example, Alshahrani et al. (2022) have studied the three Arabic Wikipedia editions, Arabic Wikipedia (AR), Egyptian Arabic Wikipedia (ARZ), and Moroccan Arabic Wikipedia (ARY), and documented issues in the Egyptian Wikipedia with automatic creation/generation and translation of con-

tent pages from English without human supervision. They stressed that these issues could substantially affect the performance and accuracy of Large Language Models (LLMs) trained from these corpora, producing models that lack the cultural richness and meaningful representation of native speakers. In another research work by the same authors, they investigated the performance implications of using inorganic, unrepresentative corpora, mainly generated through automated techniques such as bot generation or automated template-based translation, to train a few masked language models and word embedding models. They found that models trained on bot-generated or template-translated articles underperformed the models trained on human-generated articles and underscored that, for good NLP performance, researchers need both large and organic corpora (Alshahrani et al., 2023a).

In this paper, we solely focus on the problem of template translation that took place in the Egyptian Arabic Wikipedia edition, where a few registered users employed simple templates to translate more than one million content pages (articles) from English to Arabic using Google Translate, all without translation error checking or culture misrepresentation verification, disregarding the consequences of using such poor articles (Baker, 2022; Das, 2020; Alshahrani et al., 2022; Agrawal et al., 2023; Al-Khalifa et al., 2024; Thompson et al., 2024). We first explore the three Arabic Wikipedia editions' content in terms of density, quality, and human contributions, highlighting how the template-based

\*<https://hf.co/spaces/Egyptian-Wikipedia-Scanner>.

<sup>1</sup>Google Translate: <https://translate.google.com>.

translation occurred on the Egyptian Wikipedia produces unrepresentative content. We second, attempt to build powerful multivariate machine learning classifiers leveraging corpus/articles' metadata to detect the template-translated articles automatically. We then deploy and publicly host the best-performing classifier, XGBoost, so researchers, practitioners, and other users can benefit from our online, web-based detection system. We lastly argue that practices such as template translations could not only impact the performance of models trained on these template-translated articles but also could misrepresent the native speakers and their culture and do not echo their views, beliefs, opinions, or perspectives.

## 2. Exploratory Analysis

We explore, in the following subsections, the three Arabic Wikipedia editions, Arabic Wikipedia (AR), Egyptian Arabic Wikipedia (ARZ), and Moroccan Arabic Wikipedia (ARY), regarding their articles' content in terms of density, quality, and human contributions.

### 2.1. Analysis Setup

We follow the same methodology [Alshahrani et al. \(2023a\)](#) used to quantify the bot-generated articles, but we, here, utilize the Wikimedia `XTools` API<sup>2</sup> to collect all Arabic Wikipedia editions' articles' metadata; specifically, we collect the total edits, total editors, top editors, total bytes, total characters, total words, creator name, and creation date for each article. We use the complete Wikipedia dumps of each Arabic Wikipedia edition, downloaded on the 1st of January 2024 ([Wikimedia, 2024](#)) and processed using the `Gensim` Python library ([Řehůřek and Sojka, 2010](#)). We also use Wikipedia's "List Users" service<sup>3</sup> to retrieve the full list of bots in each Arabic Wikipedia edition to measure the bot and human contributions to each article.

### 2.2. Shallow Content

We, in this subsection, study the density of the content of the three Arabic Wikipedia editions, highlighting general statistics and token/character length distributions per Arabic Wikipedia edition.

#### 2.2.1. Summary Statistics

We shed light on a few general statistics of the three Arabic Wikipedia editions regarding their total articles, total extracted articles, corpus size, total

bytes, total characters, and total tokens, highlighting the minimum, maximum, and mean values of the three articles' metadata: total bytes, total characters, and total tokens.<sup>4</sup> From [Table 1](#), it is notable that the Egyptian Arabic Wikipedia has a greater number of total articles than the Arabic Wikipedia (which is generally believed to be more organically generated), with almost 400K articles, yet as we will discuss later in [Table 3](#), this number of total articles does not reflect true measurements of organically generated contributions since all three Arabic Wikipedia editions include substantial bot generation and template translation activities ([Baker, 2022](#); [Alshahrani et al., 2022, 2023b](#)). We employ the `Gensim` Python library to parse and extract the textual content (articles) from each Wikipedia dump file. However, since the library discards any articles with less than 50 tokens/words, all three Arabic Wikipedia editions lost many articles. For example, the Egyptian Wikipedia lost nearly 741K (46%) of its articles, whereas the Moroccan Wikipedia and the Arabic Wikipedia lost 2.9K (30%) and 346K (28%) of their articles, respectively. This loss of articles exhibits how the Egyptian Arabic Wikipedia contains almost half of its total articles under 50 tokens per article, indicating that it has more limited and shallow content and reflecting the template translation that occurred on its articles.

#### 2.2.2. Token/Character Length Distribution

We visualize, in [Figure 1](#), the token and character distributions for each Arabic Wikipedia edition by plotting the tokens per article and characters per article with the mean lines for each Arabic Wikipedia edition. We observe that the Egyptian Wikipedia length distributions (token and character) are less dense than the Arabic Wikipedia and Moroccan Wikipedia, and a notable number of articles in the Egyptian Wikipedia are below the mean line/threshold, exhibiting that the Egyptian Wikipedia has unusually smaller and shorter articles than other Arabic Wikipedia editions. Surely, the Egyptian Wikipedia has more articles than the other Arabic Wikipedia editions, but it does have the lowest mean values of the total of characters and total of tokens/words, 610 and 100, respectively, compared to the mean values of the Arabic Wikipedia and the Moroccan Wikipedia, as shown in [Table 1](#). These observations signal that the template translation that happened on its articles does not produce rich, dense, and long content but only produces poor, limited, and shallow content.

<sup>4</sup>We use the Wikimedia Statistics service, <https://stats.wikimedia.org>, to retrieve the total articles (content pages) for each Arabic Wikipedia edition, whereas the other statistics are generated from the extracted articles from each Arabic Wikipedia edition.

<sup>2</sup>XTools API: <https://www.mediawiki.org/wiki/XTools>.

<sup>3</sup><https://{{WIKI}}.wikipedia.org/wiki/Special:ListUsers>.

Wikipedia	Total Articles	Extracted Articles	Corpus Size	Total Bytes			Total Characters			Total Tokens		
				Min	Max	Mean	Min	Max	Mean	Min★	Max	Mean
Arabic (AR)	1,226,784	880,334	2.6GB	6,424,572,842			1,564,243,778			264,761,062		
				488	1,419,547	7,297	200	334,464	1,776	50	56,395	300
Egyptian (ARZ)	1,621,745	736,158	766MB	1,525,938,072			449,449,693			74,277,188		
				515	1,217,036	2,072	233	399,641	610	50	74,009	100
Moroccan (ARY)	9,659	6,754	11MB	25,109,824			6,802,694			1,153,946		
				646	105,009	3,717	248	32,853	1,007	50	5,635	170

Table 1: General statistics of the three Arabic Wikipedia editions in terms of total articles, total extracted articles, corpus/articles size, total bytes, total characters, and total tokens. ★As a result of the `Gensim` Python library discarding articles with tokens/words less than 50, all minimum tokens of articles are 50.

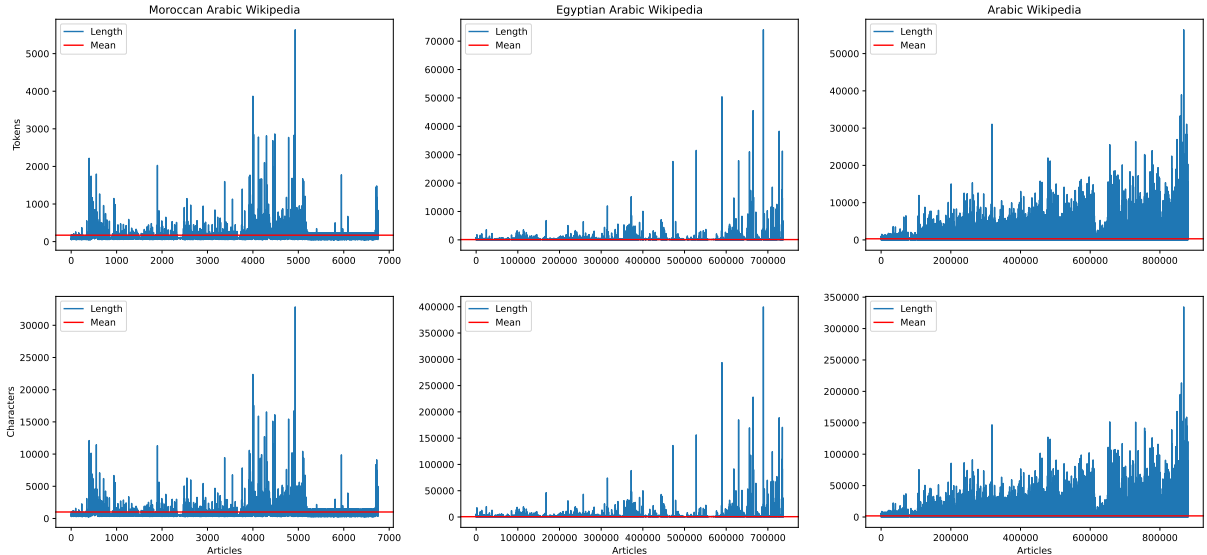


Figure 1: Visualizations of tokens and characters per article for each Arabic Wikipedia edition, displaying the total tokens and characters on the y-axes and articles on the x-axes, with plotting the mean lines.

## 2.3. Poor Quality Content

We study the quality of the Arabic Wikipedia editions’ content regarding lexical richness and diversity and the most common and duplicate n-grams.

### 2.3.1. Lexical Richness/Diversity

We use the terms lexical richness and lexical diversity equivalently and interchangeably in this study, as [Daller et al. \(2003\)](#) suggested. To measure the lexical richness and diversity, we first compute the total tokens and unique tokens per Arabic Wikipedia edition, and second, we utilize three simple but widely used lexical richness metrics: Type-Token Ratio (TTR) ([Chotlos, 1944](#); [Templin, 1957](#)), Root Type-Token Ratio (RTTR) ([Guiraud, 1954, 1959](#)), and Corrected Type-Token Ratio (CTTR) ([Carroll, 1964](#)). Yet, as many have emphasized, like [McCarthy \(2005\)](#), we find that these metrics are not often precise and sometimes erroneous and do not reflect the true lexical richness and diversity of a corpus. For example, we observe that the TTRs of Arabic Wikipedia and Egyptian Wikipedia are identical, and the RTTRs and CTTRs of Egyptian Wikipedia and Moroccan Wikipedia are similar, despite the massive difference between the Arabic Wikipedia editions’ corpora in terms of the lexicon size and vo-

cabulary size, as shown in [Table 2](#). Therefore, we adopt an advanced metric to measure the lexical richness and diversity called ‘Measure of Textual Lexical Diversity (MTLD)’, introduced by [Mccarthy and Jarvis \(2010\)](#). We utilize the `LexicalRichness` Python library’s implementation of the MTLD metric with a default factor size of 0.720 ([Shen, 2022](#)). We find that the results are consistent with the other metrics, as reported in [Table 2](#), in that the Moroccan Wikipedia has the best lexical richness and diversity among the three Arabic Wikipedia editions, where the Arabic Wikipedia comes second, and Egyptian Wikipedia comes in last, documenting the Egyptian Arabic Wikipedia corpus is not lexically rich and diverse, which we attribute to the template-based translation took place on its articles (content pages).

### 2.3.2. Most Common/Duplicate N-Grams

We generate n-grams from each Arabic Wikipedia corpus, where  $n=\{1, 2, 3, 5, 10, 50\}$ , to highlight the common and duplicate n-grams. We hypothesize that the higher the count of n-grams in an Arabic Wikipedia corpus, especially when  $n=\{5, 10, 50\}$ , the more we can detect templates used in the template translation activities in the Arabic Wikipedia

Wikipedia	Total Tokens	Unique Tokens	Type-Token Ratio (TTR)	Root Type Token Ratio (RTTR)	Corrected Type Token Ratio (CTTR)	Measure of Textual Lexical Diversity (MTLD)
Arabic (AR)	264,777,392	2,867,782	0.010	176.24	124.62	71.20
Egyptian (ARZ)	74,278,320	759,519	0.010	88.12	62.31	45.69
Moroccan (ARY)	1,154,058	94,827	0.082	88.27	62.41	89.77

Table 2: Calculations of four lexical richness and diversity metrics, TTR, RTTR, CTTR, and MTLD, accompanied with total tokens (lexicon) and unique tokens (vocabulary) for each Arabic Wikipedia edition.

editions, specifically in the Egyptian Wikipedia. We notice that n-grams in the Egyptian Wikipedia have very large counts compared to the Arabic and Moroccan Wikipedia editions, as shown in Tables 9 and 10 in Appendix A.<sup>5</sup> In Figure 2, we visualize the log values of the top K=1 counts of common and duplicate n-grams generated from each Arabic Wikipedia corpus, where  $n=\{1, 2, 3, \dots, 50\}$ , and we observe that all the n-grams in all the Arabic Wikipedia editions exhibit exponential decay, drastically (like Arabic Wikipedia) or gradually (like Egyptian Wikipedia and Moroccan Wikipedia). Yet, the large counts of Egyptian Wikipedia’s n-grams when  $n \geq \{5\}$  do not decline exponentially but linearly, suggesting that there is an anomaly in the Egyptian Wikipedia corpus, where the n-grams of the normally generated corpus by humans usually factorially decreases, as the n value increases. We believe the template-based translation on the Egyptian Wikipedia creates such an anomaly, as many parts/grams/phrases of templates used in the translation are duplicated repeatedly in its corpus.

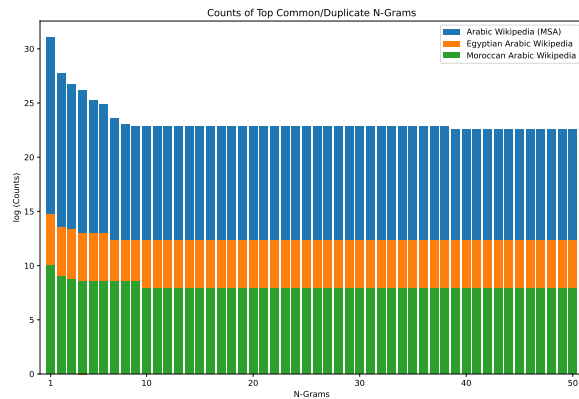


Figure 2: Counts of top common/duplicate n-grams of each Arabic Wikipedia edition; log values/counts are only for top K=1 common/duplicate n-grams.

## 2.4. Misleading Human Involvement

We shed light on the human involvement across the three Arabic Wikipedia editions, specifically the type of page creators and editors, debating how the template translation activities could produce misleading metadata regarding human involvement.

<sup>5</sup>We further analyze the 5-grams and 10-grams of each Arabic Wikipedia edition in Appendix A.

### 2.4.1. List of Contributors

We collect all the page creators for each article in the Arabic Wikipedia editions, count the number of their contributions (article creations), and categorize them into bots and humans. As shown in Table 3, it is clear that the Arabic Wikipedia and Moroccan Wikipedia suffer from mass auto-creation of articles by bots, especially by the ‘JarBot’, which has created nearly 260K articles (29.31%) in the Arabic Wikipedia, and the ‘DarijaBot’, which has created nearly 3.2K articles (34%) in the Moroccan Wikipedia.<sup>6</sup> However, the worst of all is the unguided, unreviewed, unsupervised template translation of articles from English in the Egyptian Wikipedia by registered users, largely by two registered users, ‘HitomiAkane’ and ‘Al-Dandoon’, who have created more than 1.4M articles (88.57%) and 113K articles (6.99%), respectively.<sup>7</sup>

### 2.4.2. Type of Contributors

We calculate the percentage of creators and editors of articles (bots and humans) in each Arabic Wikipedia edition. We use the absolute count of page creators and classify the creators based on their types, bots or humans, while with the page editors, we calculate the percentage using the total number of editors on each article and set a threshold of 50%, where if an article was edited by more than 50% by bots, we then consider this article a bot-edited, and vice versa. As shown in Figure 3, we see bots often create articles side-by-side with humans in the Arabic Wikipedia (31.5%) and Moroccan Wikipedia (22.30%) editions, which is normal and permitted to a certain degree according to Wikipedia’s bot policy (Wikipedia, 2024b). However, in the Egyptian Wikipedia edition, we observe that its articles are 100% created by humans, i.e., registered users, and this percentage is misleading given that 42.72% of its articles are

<sup>6</sup>These two bots, ‘JarBot’ and ‘DarijaBot’, have approval from Wikimedia to operate on the Arabic Wikipedia and the Moroccan Wikipedia (Wikidata, 2024b,a).

<sup>7</sup>These two registered users were local admins of the Egyptian Arabic Wikipedia edition until their permissions were revoked in May 2020 by the Stewards, the global admins of the Wikipedia project, for their abuse of admin permissions and their massive unsupervised and unauthorized creation of articles (Wikipedia, 2020).

<sup>8</sup>Wikiscan Statistics service: <https://wikiscan.org>.

Wikipedia \ Rank (percentage)	1st (%)	2nd (%)	3rd (%)	4th (%)	5th (%)	
Arabic (AR)	Name	JarBot	Mr. Ibrahim	جار الله	ElphiBot	Majed
	Count	359,677 (29.31%)	52,222 (4.25%)	43,691 (3.56%)	42,669 (3.47%)	26,228 (2.13%)
	Type	Bot	Human	Human	Bot	Human
Egyptian (ARZ)	Name	HitomiAkane	Al-Dandoon	Raafat	Ghaly	حمدى10
	Count	1,436,430 (88.57%)	113,468 (6.99%)	18,334 (1.13%)	7,212 (0.44%)	2,720 (0.16%)
	Type	Human	Human	Human	Human	Human
Moroccan (ARY)	Name	DarjaBot	Tifratin	Ideophagous	Sedrati	Rachidourkia
	Count	3,285 (34%)	1,302 (13.47%)	1,231 (12.74%)	765 (7.92%)	540 (5.59%)
	Type	Bot	Human	Human	Human	Human

Table 3: Top five page creators in the Arabic Wikipedia editions, highlighting their types (bots or humans) and how many articles they have created until March 1st, 2024, according to Wikiscan Statistics service.<sup>8</sup>

automatically template-translated from English to Arabic using templates without human supervision or intervention, as documented by Baker (2022) and Alshahrani et al. (2022).

### 3. Experimental Setup

We, here, attempt to build classifiers to identify and mitigate the impacts of the template-translated articles on the Egyptian Wikipedia edition since it particularly suffers from template translations, as documented by Alshahrani et al. (2022). We first extract all articles with their metadata, split the articles into two categories: before and after the template-based translation occurred, and lastly, label, preprocess, and encode these categorized articles using Arabic pre-trained models.

#### 3.1. Dataset Filtrating and Labeling

We follow a few heuristic rules to classify Egyptian Wikipedia into articles created before and after the massive template-based translation activities related to creation dates, total edits, and types of creators and editors. We take insights from our exploratory analysis, section 2, the Wikimedia Statistics service, and the previous research works that documented the template translation activities in the Egyptian Wikipedia (Baker, 2022; Alshahrani et al., 2022; Wikimedia Statistics, 2024), to craft these rules, specifically when selecting the dates.

Category	Total
<b>Total Articles (both categories)</b>	736,107
<b>Articles Before Template Translation</b>	11,126
<b>Articles After Template Translation</b>	155,275
<b>Uncategorized Articles</b>	569,706

Table 4: Statistics of filtered articles after applying our heuristic filtration rules, displaying the totals.

We list the heuristic rules for filtering the articles created *before* and *after* the translations in Appendix B, where we employ more rigorous heuristic rules to filter the articles created *after* the template translation appeared on the Egyptian Wikipedia. In Table 4, we show the statistics of our rule-based

filtration process. We then randomly select 10K articles from each category to train a multivariate machine learning classifier to detect the template-based translations automatically. We lastly label the articles *before* translation as 'human-generated' articles since all articles are created by registered users and label the articles *after* translation as 'template-translated' articles.

#### 3.2. Dataset Preprocessing

We lightly preprocess the filtered articles by replacing all non-alphanumeric and non-Arabic characters with white spaces and normalizing the extra unnecessary whitespaces to one whitespace. We do not apply stemming, lemmatization, or any Arabic text normalization on the articles to have organic content (articles) as much as possible.

#### 3.3. Dataset Encoding

We use two different types of embedding techniques to encode the randomly selected 20K articles separately: pre-trained Egyptian Arabic context-independent word embeddings (Word2Vec) of the size of 300 dimensions from Spark-NLP Python library<sup>9</sup> and context-dependent word embeddings (contextual) of the size of 768 dimensions produced by utilizing the pre-trained CAMeL-BERT-Mix POS-EGY model<sup>10</sup> (Inoue et al., 2021) as our feature extraction model. The goal is to test with different embedding techniques to maximize the performance of our multivariate machine learning classifiers and investigate how the type and size of the word embeddings would affect their performance.

### 4. Template Translation Detection

We experiment with a few supervised classification algorithms and unsupervised clustering algorithms

<sup>9</sup>Word2Vec Embeddings in Egyptian Arabic (300d): [https://sparknlp.org/2022/03/14/w2v\\_cc\\_300d\\_arz\\_3\\_0](https://sparknlp.org/2022/03/14/w2v_cc_300d_arz_3_0).

<sup>10</sup>CAMeL-BERT-Mix POS-EGY model: <https://huggingface.co/CAMeL-Lab/bert-base-arabic-camelbert-mix-pos-egy>.

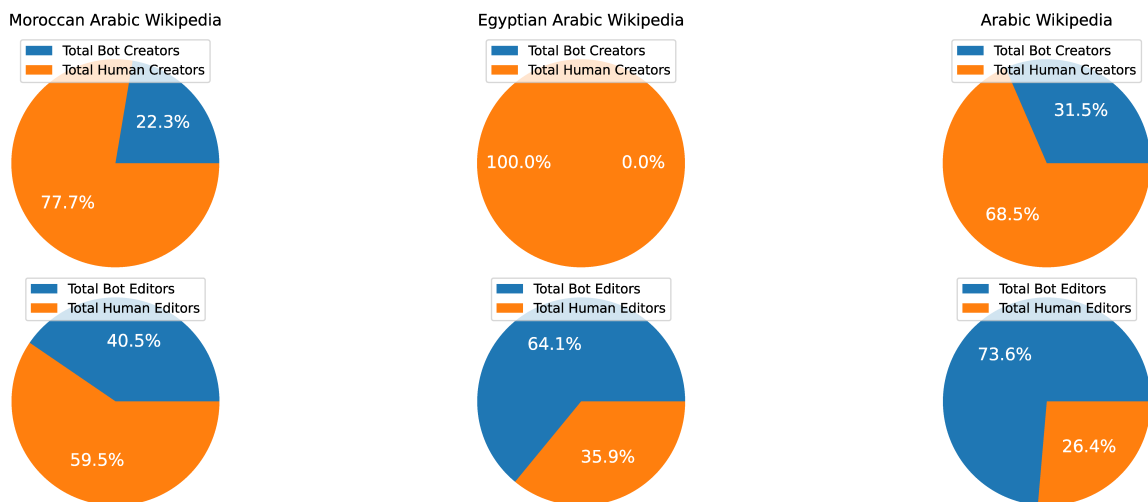


Figure 3: Visualizations displaying the percentage of article creators and editors in terms of their types, bots, and humans, and their number of contributions (article creations) in each Arabic Wikipedia edition.

to determine which approach and algorithm will best solve our template-based translation problem.

#### 4.1. Input Features Extraction

We aim to leverage the metadata of corpus, i.e., articles, collected using Wikimedia services to detect the template-translated articles in the Egyptian Wikipedia edition. Besides utilizing pre-trained Word2Vec and CAMELBERT word embeddings as input features, we also include the metadata we collect about every article: total edits, total editors, total bytes, total characters (charts), and total words. Overall, we test the machine learning algorithms' performance using three input features: only embeddings, only metadata, or both (metadata and embeddings), as illustrated in Figure 4.

#### 4.2. Metadata Ablation Studies

We perform two ablation studies for each machine learning algorithm (classification and clustering) to determine the best metadata features to include in the input features. We first test each metadata feature's performance individually and then combine two, three, and all metadata features consecutively.

#### 4.3. Classification Algorithms

We select five supervised classification algorithms to solve our multivariate classification problem: Logistic Regression (LR) (Fan et al., 2008), Support Vector Machine (SVM) (Chang and Lin, 2011), Gaussian Naive Bayes (GNB) (Pedregosa et al., 2011), Random Forests (RF) (Breiman, 2001), and XGBoost (eXtreme Gradient Boosting) (XGBoost, 2024). We, in the next subsections, discuss the experimental setups and the performance results of these supervised machine learning classifiers.

##### 4.3.1. Classification Experimental Setup

We split the randomly selected 20K articles into training (80%) and testing (20%) splits with data shuffling and stratification enabled to ensure that the training and test splits are randomized and have the same proportion of each class. We further evaluate our classifiers using the accuracy metric with the Stratified K-Folds Cross-Validation technique, where we set the number of folds  $K=5$ , ensuring every fold has a representative class distribution.

##### 4.3.2. Results of Classification Ablations

We report, in Table 5, the evaluation accuracy results on the testing splits of our metadata ablations. We can see that all machine learning classifiers achieve excellent (100%) to very good performance ( $100% > \text{accuracy} > 90%$ ) with the total edits and total editors separated or combined. In contrast, metadata features like the total bytes, total characters, and total words perform from fairly to poorly and, unfortunately, decrease the overall performance of all metadata features combined with some classifiers like SVM. Generally, we observe that the ensemble classifiers (RF and XGBoost) outperform the other classifiers even with the metadata features that contribute less to the classifiers' learning.

##### 4.3.3. Results of Classification Algorithms

We show, in Table 6, the evaluation accuracy scores on the testing splits of the multivariate machine learning classifiers studied, demonstrating how the classifiers would perform with three input features: two embedding styles (Word2Vec or CAMELBERT), corpus/articles metadata, and both embeddings and metadata combined. Here, we decided to include all the articles' metadata, not only



Figure 4: A basic process chart demonstrating the studied input features: embeddings (two word embeddings of sizes 300 or 768), metadata (five metadata of articles), or both (embeddings + metadata).

Classifier	Metadata							
	A	B	C	D	E	A+B	C+D+E	All
<b>Logistic Regression</b>	100	100	88.30	83.85	84.67	100	89.03	98.42
<b>Support Vector Machine</b>	90.30	100	87.95	83.60	83.95	99.78	87.62	87.75
<b>Naive Bayes</b>	100	100	82.00	74.28	78.00	100	80.50	99.60
<b>Random Forest</b>	100	100	86.17	82.23	84.80	100	91.25	100
<b>XGBoost</b>	100	100	88.60	84.52	84.70	100	90.53	100

Table 5: Accuracies of metadata ablations of the studied classifiers. Encoded columns denote metadata features as follows: A) total edits, B) total editors, C) total bytes, D) total characters, and E) total words.

Classifier	Embeddings		Metadata	Both (Embeddings + Metadata)	
	Word2Vec	CAMeLBERT		Word2Vec	CAMeLBERT
<b>Logistic Regression</b>	91.22	99.30	98.42	99.40	100
<b>Support Vector Machine</b>	99.02	98.45	87.75	87.90	87.90
<b>Naive Bayes</b>	88.90	95.17	99.60	99.60	99.52
<b>Random Forest</b>	98.08	98.17	100	100	99.95
<b>XGBoost</b>	98.28	98.78	100	100	100

Table 6: Accuracies of the machine learning classifiers studied, showing their performance with different input features: two embedding styles, corpus metadata, and both embeddings and metadata combined.

the features that performed well in our ablation studies, to diversify the classifiers’ learning and ensure that each category of the Egyptian Wikipedia articles (human-generated and template-translated) is well-represented. We report, again, that the SVM classification algorithm underperforms all the other algorithms and find that the metadata features present a bottleneck performance for it (i.e., highly variable features). We attribute the poor performance to the complex, multivariate nature of the dataset, specifically, the high variability of the metadata features like the total bytes, words, and characters, as seen in Table 5.<sup>11</sup> On the other bright side, we find that ensemble classification algorithms like RF and XGBoost excel and outperform the traditional, single classification algorithms due to their ability to overcome noise, bias, and variance; the RF algorithm uses the bagging technique, and XGBoost algorithm uses boosting technique to handle such technical challenges.<sup>12</sup>

<sup>11</sup>We handled the dataset noise through our filtration process and the bias by balancing the dataset classes, yet the dataset variance is challenging due to the high dispersion in metadata features collected.

<sup>12</sup>As an online application, we deploy our best classifier, XGBoost, with input features of metadata and CAMeLBERT embeddings. See Appendix C for details.

## 4.4. Clustering Algorithms

We explore three different unsupervised clustering algorithms to solve the template-based translation problem: K-Means (Wu, 2012), Hierarchical Agglomerative (Zepeda-Mendoza and Resendis-Antonio, 2013), and DBSCAN (Density-Based Spatial Clustering of Applications with Noise) (Ester et al., 1996). We, in the following, discuss the experimental setups and the performance results of these unsupervised machine learning clusters.

### 4.4.1. Clustering Experimental Setup

We feed the unsupervised clustering algorithms all the randomly selected 20K articles after removing the labels without splitting them due to their nature. We set the number of clusters to K=2 since our dataset only has two categories (human-generated and template-translated). We further evaluate our clusters using the Silhouette coefficient with the Euclidean distance, a widely used internal evaluation metric to measure how cohesive and separated the clusters are, based on the distances or similarities between the data points, i.g., articles.<sup>13</sup>

<sup>13</sup>Values of the Silhouette coefficient are always between 1 and -1. We apply a percentage normalization

#### 4.4.2. Results of Clustering Ablations

We report the Silhouette scores of our metadata ablation studies in Table 7. We can see that all machine learning clusterers achieve great performance with the total bytes, total characters, and total words, separated or combined, except for the DBSCAN algorithm. In contrast, metadata features like the total edits and total editors perform from fairly to poorly with K-Means and Hierarchical clustering algorithms, except for the DBSCAN algorithm. The results of these metadata ablations indicate an opposite behavior from those discussed in subsection 4.3.2, where the previously weak metadata features for the classification algorithms, like the total bytes, words, and characters, became strong metadata features for the clustering algorithms instead of the total edits and editors, which were previously strong. Generally, the K-Means and Hierarchical clustering algorithms outperform the DBSCAN algorithm even with the metadata features that contribute more to the clusterers' learning.

#### 4.4.3. Results of Clustering Algorithms

We show, in Table 8, the Silhouette scores of the machine learning clusterers studied, demonstrating how the unsupervised clusterers would perform with three input features: two embedding styles (Word2Vec or CAMELBERT), corpus/articles metadata, and both embeddings and metadata combined. We, here, fit all the articles' metadata, not only the features that performed well in our ablation studies, to diversify the clusterers' learning and ensure that each class of the Egyptian Arabic Wikipedia articles (human-generated and template-translated) is included. We report that all the clustering algorithms perform poorly with the word embeddings as features, whereas the metadata features present a performance improvement. We assume clustering the word embeddings is challenging, especially with their large dimensionality; Word2Vec's size is 300, and CAMELBERT's is 768. Overall, the unsupervised clustering algorithms underperform the supervised classification algorithms, yet we can confirm that the clustering algorithms do better with low-dimensionality features like articles' metadata, even though they introduce high-variable and dispersed features.

## 5. Discussion

We discuss three negative implications of the unguided, unreviewed, unsupervised template-based translation from English to Arabic on the Egyptian Wikipedia articles: societal, representation, and

performance implications. On the societal implications, we argue that using off-the-shelf-translation tools like Google Translate, which is widely known for its social problems like gender, cultural, and religious biases and stereotypes, could not only cause linguistic and grammatical errors but also amplify these social risks like biases and stereotypes (Prates et al., 2020; Ullmann and Saunders, 2021; Lopez-Medel, 2021; Naik et al., 2023; Al-Khalifa et al., 2024). Many researchers have emphasized how unsupervised translations are prone to serious gender bias issues, like producing translations with inaccurate gender, that could impact native speakers. For example, Stanovsky et al. (2019) have automatically evaluated the gender bias for eight highly-gendered languages like Arabic and found that a few popular industrial and academic machine translation systems (like Google Translate and Microsoft Translator<sup>14</sup>) were significantly prone to gender-biased translation errors for all tested target languages. We believe those machine translation systems are greatly beneficial tools, yet they should not be used to naively, directly, or automatically translate content without human review, especially if the content is related to the societal representation of Arabic native speakers.

On the representation implications, we argue that such automatic template-based translations without humans in the loop could misrepresent the Egyptian Arabic native speakers, where instead of the Egyptian people enriching the content of Wikipedia by sharing their voices, opinions, knowledge, perspectives, and experiences, a couple of registered users automated the creation and translation of more than a million and a half million articles (95.56%) from English on their behalf without supervision or revision of the translated articles, disregarding that the main goal of Wikipedia is to be written by the people to the people (Cohen, 2008). Another troubling drawback of such a practice is the cultural misrepresentation of the Egyptian people and their culture, where the unfiltered and unsupervised translation from English could introduce content that is not representative of the culture of native speakers. Lastly, we argue that including culturally unrepresentative articles from the Egyptian Arabic Wikipedia in pre-training corpora for language models could present cultural implications and generate culturally misaligned outputs from these models, where the majority of Arabic and multilingual language models have been fundamentally pre-trained on Wikipedia dumps like Jais and Jais-chat (Sengupta et al., 2023), AraMUS (Alghamdi et al., 2023), and JASMINE (Nagoudi et al., 2023). We believe research works, like ours, that automatically identify these template-translated articles could promote data transparency and help

---

(multiply values by 100) when reporting the values to draw a head-to-head comparison between algorithms.

<sup>14</sup>Microsoft Bing: <https://www.bing.com/translator>.



Clusterer	Metadata							
	A	B	C	D	E	A+B	C+D+E	All
<b>K-Means</b>	82.68	78.32	97.10	96.46	96.39	81.77	96.89	96.89
<b>Hierarchical</b>	86.85	81.42	97.10	97.37	97.32	82.28	96.08	97.52
<b>DBSCAN</b>	97.80	99.62	37.20	67.58	89.79	77.11	68.35	68.33

Table 7: Silhouette scores of the metadata ablations of the studied clusterers. Encoded columns denote metadata features: A) total edits, B) total editors, C) total bytes, D) total characters, and E) total words.

Clusterer	Embeddings		Metadata	Both (Embeddings + Metadata)	
	Word2Vec	CAMeLBERT		Word2Vec	CAMeLBERT
<b>K-Means</b>	12.50	14.95	96.89	96.89	96.89
<b>Hierarchica</b>	11.79	10.82	97.52	96.77	96.77
<b>DBSCAN</b>	61.64	8.43	68.33	68.34	68.68

Table 8: Silhouette scores of the machine learning clusterers studied, showing their performance with different features: two embedding styles, corpus/articles metadata, and both embeddings and metadata.

researchers make an informed decision about what to include in their pre-training datasets/corpora.

On the performance implications, we argue that the template-based translations that occurred on the Egyptian Wikipedia produce not only short and shallow articles, where we have reported that nearly 46% of the Egyptian Wikipedia articles are less than 50 tokens/words and recognized a large number of duplicate n-grams due to the templates used in translations, but also articles that lack lexical richness and diversity, where we have found that the Egyptian Wikipedia scored the worst among other Arabic Wikipedia editions in the MTLT metric. These poorly translated articles could negatively impact the performance of language models and NLP tasks that are trained on them. One research that supports our claim is the recent work of [Alshahrani et al. \(2023a\)](#), where they documented that models trained on the template-translated articles of the Egyptian Wikipedia performed the worst when compared with the models trained on the Arabic Wikipedia articles. Finally, we recommend excluding the unfiltered template-translated articles from Egyptian Wikipedia from training datasets to mitigate their negative societal, representation, and performance implications and encourage using automatic detection systems, like ours, to identify such articles that are not only mispicturing the Egyptian people and their culture but also affecting the performance of language models and NLP tasks.

## 6. Limitations

We leverage five metadata of articles of different sizes (total edits, total editors, total bytes, total characters, and total words) and then append them to two types of word embeddings (Word2Vec and CAMeLBERT) of sizes of 300 or 768 vectors to build powerful classifiers, yet concatenating all these different features could produce highly vari-

able features due to the high dispersion between the extracted input features, which could present a performance challenge for our proposed automatic detection system and could increase the non-deterministic behavior of its classifiers.

## 7. Conclusion

We attempt to mitigate the template translations on the Egyptian Arabic Wikipedia by identifying these template-translated articles and their characteristics through exploratory analysis and developing automatic detection systems. We first investigate the content of the three Arabic Wikipedia editions in terms of density, quality, and human contributions and use such insights to build powerful multivariate machine learning classifiers leveraging articles' metadata to detect template-translated articles automatically; we find that the supervised classification algorithms are better than the unsupervised clustering algorithms. We then publicly deploy the best-performing classifier, XGBoost, as an application and release the extracted, filtered, labeled, and preprocessed datasets to the community to benefit from our datasets and the online detection system.

## Reproducibility

We share our labeled datasets, code and scripts of the exploratory analysis, and the multivariate machine learning classifiers on GitHub at <https://github.com/SaiedAlshahrani/leveraging-corpus-metadata>.

## Acknowledgments

We would like to thank Clarkson University and the Office of Information Technology (OIT) for providing computational resources. We also would like to thank Norah Alshahrani for her valuable feedback.

## Bibliographical References

- Ameeta Agrawal, Lisa Singh, Elizabeth Jacobs, Yaguang Liu, Gwyneth Dunlevy, Rhitabrat Pokharel, and Varun Uppala. 2023. [All Translation Tools Are Not Equal: Investigating the Quality of Language Translation for Forced Migration](#). In *2023 IEEE 10th International Conference on Data Science and Advanced Analytics (DSAA)*, pages 1–10.
- Hend Al-Khalifa, Khaloud Al-Khalefah, and Hesham Haroon. 2024. [Error Analysis of Pretrained Language Models \(PLMs\) in English-to-Arabic Machine Translation](#). *Human-Centric Intelligent Systems*.
- Asaad Alghamdi, Xinyu Duan, Wei Jiang, Zhenhai Wang, Yimeng Wu, Qingrong Xia, Zhefeng Wang, Yi Zheng, Mehdi Rezagholizadeh, Baoxing Huai, Peilun Cheng, and Abbas Ghaddar. 2023. [ArAMUS: Pushing the Limits of Data and Model Scale for Arabic Natural Language Processing](#). *arXiv preprint arXiv:2306.06800*.
- Saied Alshahrani, Norah Alshahrani, Soumyabrata Dey, and Jeanna Matthews. 2023a. [Performance Implications of Using Unrepresentative Corpora in Arabic Natural Language Processing](#). In *Proceedings of ArabicNLP 2023*, pages 218–231, Singapore (Hybrid). Association for Computational Linguistics.
- Saied Alshahrani, Norah Alshahrani, and Jeanna Matthews. 2023b. [DEPTH+: An Enhanced Depth Metric for Wikipedia Corpora Quality](#). In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, pages 175–189, Toronto, Canada. Association for Computational Linguistics.
- Saied Alshahrani, Esmā Wali, and Jeanna Matthews. 2022. [Learning From Arabic Corpora But Not Always From Arabic Speakers: A Case Study of the Arabic Wikipedia Editions](#). In *Proceedings of the The Seventh Arabic Natural Language Processing Workshop (WANLP)*, pages 361–371, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Maher Asaad Baker. 2022. *How I Wrote a Million Wikipedia Articles*, 2 edition. BookRix GmbH & Co. KG., Munich, Germany.
- Runa Bhattacharjee and Pau Giner. 2022. [You Can Now Use Google Translate to Translate Articles on Wikipedia](#). Last accessed on 2024-03-01.
- Leo Breiman. 2001. Random Forests. *Machine Learning*, 45:5–32.
- John Bissell Carroll. 1964. *Language and Thought*. Prentice-Hall.
- Chih-Chung Chang and Chih-Jen Lin. 2011. [LIB-SVM: A Library for Support Vector Machines](#). In *ACM Transactions on Intelligent Systems and Technology*, volume 2, New York, NY, USA. Association for Computing Machinery.
- John W Chotlos. 1944. IV. A Statistical and Comparative Analysis of Individual Written Language Samples. *Psychological Monographs*, 56(2):75.
- Noam Cohen. 2008. [Open-Source Troubles in Wiki World](#). The New York Times. Last accessed on 2024-03-01.
- Helmut Daller, Roeland van Hout, and Jeanine Treffers-Daller. 2003. [Lexical Richness in the Spontaneous Speech of Bilinguals](#). *Applied Linguistics*, 24(2):197–222.
- Alok Das. 2020. [Neural Machine Translation \(NMT\): Inherent Inadequacy, Misrepresentation, and Cultural Bias](#). *International Journal of Translation*, 32:115–145.
- M Ester, H P Kriegel, J Sander, and Xu Xiaowei. 1996. [A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases With Noise](#). *U.S. Department of Energy Office of Scientific and Technical Information*.
- Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. 2008. [LIBLINEAR: A Library for Large Linear Classification](#). *the Journal of Machine Learning Research*, 9:1871–1874.
- Pierre Guiraud. 1954. *Les Caractères Statistiques du Vocabulaire: Essai de Méthodologie*. Presses universitaires de France, Paris, France.
- Pierre Guiraud. 1959. *Problèmes et Méthodes de la Statistique Linguistique*. D. Reidel, Dordrecht, Holland.
- Ari Hautasaari. 2013. [“Could Someone Please Translate This?”: Activity Analysis of Wikipedia Article Translation by Non-experts](#). In *Proceedings of the 2013 Conference on Computer Supported Cooperative Work, CSCW '13*, page 945–954, New York, NY, USA. Association for Computing Machinery.
- Go Inoue, Bashar Alhafni, Nurpeiis Baimukan, Houda Bouamor, and Nizar Habash. 2021. [The Interplay of Variant, Size, and Task Type in Arabic Pre-trained Language Models](#). In *Proceedings of the Sixth Arabic Natural Language Processing Workshop*, pages 92–104, Kyiv, Ukraine (Virtual). Association for Computational Linguistics.

- Isaac Johnson and Emily Lescak. 2022. [Considerations for Multilingual Wikipedia Research](#). *arXiv preprint arXiv:2204.02483*.
- Maria Lopez-Medel. 2021. [Gender bias in machine translation: an analysis of Google Translate in English and Spanish](#). *Academia.edu*.
- Philip Mccarthy and Scott Jarvis. 2010. [MTLD, vocd-D, and HD-D: A validation study of sophisticated approaches to lexical diversity assessment](#). *Behavior research methods*, 42:381–92.
- Philip M McCarthy. 2005. [An Assessment of the Range and Usefulness of Lexical Diversity Measures and the Potential of the Measure of Textual, Lexical Diversity \(MTLD\)](#). Ph.D. thesis, The University of Memphis.
- John Mittermeier, Ricardo Correia, Rich Grenyer, Tuuli Toivonen, and Uri Roll. 2021. [Using Wikipedia to Measure Public Interest in Biodiversity and Conservation](#). *Conservation Biology*, 35.
- El Moatez Billah Nagoudi, Muhammad Abdul-Mageed, AbdelRahim Elmadany, Alcides Alcoba Inciarte, and Md Tawkat Islam Khondaker. 2023. [JASMINE: Arabic GPT Models for Few-Shot Learning](#). *arXiv preprint arXiv:2212.10755*.
- Ranjita Naik, Spencer Rarrick, and Vishal Chowdhary. 2023. [Reducing Gender Bias in Machine Translation through Counterfactual Data Generation](#). *arXiv preprint arXiv:2311.16362*.
- Sergiu Nisioi, Ella Rabinovich, Liviu P. Dinu, and Shuly Wintner. 2016. [A Corpus of Native, Non-native and Translated Texts](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 4197–4201, Portorož, Slovenia. European Language Resources Association (ELRA).
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. [Scikit-learn: Machine Learning in Python](#). *Journal of Machine Learning Research*, 12:2825–2830.
- Marcelo Prates, Pedro Avelar, and Luís Lamb. 2020. [Assessing Gender Bias in Machine Translation: A Case Study With Google Translate](#). *Neural Computing and Applications*, 32.
- Motaz Saad and Basem Alijla. 2017. [WikiDocsAligner: An Off-the-Shelf Wikipedia Documents Alignment Tool](#). In *Palestinian International Conference on Information and Communication Technology (PICICT)*.
- Neha Sengupta, Sunil Kumar Sahu, Bokang Jia, Satheesh Katipomu, Haonan Li, Fajri Koto, William Marshall, Gurpreet Gosal, Cynthia Liu, Zhiming Chen, Osama Mohammed Afzal, Samta Kamboj, Onkar Pandit, Rahul Pal, Lalit Pradhan, Zain Muhammad Mujahid, Massa Baali, Xudong Han, Sondos Mahmoud Bsharat, Alham Fikri Aji, Zhiqiang Shen, Zhengzhong Liu, Natalia Vassilieva, Joel Hestness, Andy Hock, Andrew Feldman, Jonathan Lee, Andrew Jackson, Hector Xuguang Ren, Preslav Nakov, Timothy Baldwin, and Eric Xing. 2023. [Jais and Jais-chat: Arabic-Centric Foundation and Instruction-Tuned Open Generative Large Language Models](#). *arXiv preprint arXiv:2308.16149*.
- Lucas Shen. 2022. [LexicalRichness: A Small Module to Compute Textual Lexical Richness](#). Last accessed on 2024-03-01.
- Gabriel Stanovsky, Noah A. Smith, and Luke Zettlemoyer. 2019. [Evaluating Gender Bias in Machine Translation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1679–1684, Florence, Italy. Association for Computational Linguistics.
- Mildred C. Templin. 1957. [Certain Language Skills in Children: Their Development and Interrelationships](#), NED—New edition, volume 26. University of Minnesota Press.
- Brian Thompson, Mehak Preet Dhaliwal, Peter Frisch, Tobias Domhan, and Marcello Federico. 2024. [A Shocking Amount of the Web is Machine Translated: Insights from Multi-Way Parallelism](#). *arXiv preprint arXiv:2401.05749*.
- Stefanie Ullmann and Danielle Saunders. 2021. [Google Translate is sexist. What it needs is a little gender-sensitivity training](#). Last accessed on 2024-03-01.
- Wikidata. 2024a. [Wikidata: Requests For Permissions/Bot/DarijaBot](#). Last accessed on 2024-03-01.
- Wikidata. 2024b. [Wikidata: Requests For Permissions/Bot/JarBot](#). Last accessed on 2024-03-01.
- Wikimedia. 2024. [Wikimedia Downloads](#). Last accessed on 2024-03-01.
- Wikimedia Foundation. 2022. [Content Translation – Mediawiki](#). Last accessed on 2024-03-01.
- Wikimedia Statistics. 2024. [New Pages: Egyptian Arabic Wikipedia](#). Last accessed on 2024-03-01.
- Wikipedia. 2020. [Steward Removal of Flags on ARZWiki](#). Last accessed on 2024-03-01.

Wikipedia. 2024a. [Wiki Markup](#). Last accessed on 2024-03-01.

Wikipedia. 2024b. [Wikipedia: Bot Policy](#). Last accessed on 2024-03-01.

Junjie Wu. 2012. *Advances in K-means Clustering: A Data Mining Thinking*. Springer Science & Business Media.

XGBoost. 2024. [XGBoost Documentation](#). Last accessed on 2024-03-01.

Marie Lisandra Zepeda-Mendoza and Osbaldo Resendis-Antonio. 2013. *Hierarchical Agglomerative Clustering*, pages 886–887. Springer New York, New York, NY.

Radim Řehůřek and Petr Sojka. 2010. [Software Framework for Topic Modelling with Large Corpora](#). In *Proceedings of LREC 2010 workshop New Challenges for NLP Frameworks*, pages 46–50, Valletta, Malta. University of Malta.

## A. Analysis of N-Grams

We analyze the 5-grams and 10-grams closely since they are suitable, not long or short. The n-grams in the Egyptian Wikipedia are very large compared to the Arabic and Moroccan Wikipedia editions, as indicated in Tables 9 and 10. Plus, it is noticeable that these counts do not decay exponentially as they normally should (the larger the n-gram size, the smaller the n-grams' count) but linearly and slowly (all near 222K even with different sizes of n-grams), suggesting this abnormal decay is a symptom of the template translations that Egyptian Wikipedia suffered from, where some grams/parts/phrases from the used templates are frequently and constantly repeated.

We additionally observe that most of the top ten 5-grams and 10-grams of the Moroccan Wikipedia edition are predominantly non-Arabic grams, which seems in a format of the Wikitext Markup Language (Wikipedia, 2024a), as exhibited in Tables 9, 10, and 11. We further investigate this issue by testing our parsing code scripts and find that it does not occur when parsing articles from the other two Arabic Wikipedia editions, Arabic (AR) and Egyptian (ARZ), using the same code scripts; it only surfaces when parsing the Moroccan Wikipedia articles. We attribute this issue to either leaking Wikipedia templates used to create articles or insert images into articles or an issue with the method used to dump and compress Moroccan Wikipedia articles. We urge the global and local admins of the Moroccan Wikipedia edition to investigate this issue, which could affect not only the Moroccan Wikipedia content but also the performance of perspective NLP models and tasks trained on such content.

## B. Heuristic Filtration Rules

We list the heuristic filtration rules used to filter the articles *before* and *after* the template-based translation in the Egyptian Wikipedia edition and further shed light on the effectiveness of each enforced rule. We demonstrate, in Figures 5 and 6, the effectiveness of the implemented rule-based filtration. We can see that our heuristic filtration rules are practical, as each rule consecutively and rigorously filters out unfit articles that do not meet the heuristic filtration rules.

\* Heuristic filtration rules for *before* the translation:

1. Include articles created before 2019-12-01.
2. Include articles with more than five edits.
3. Include articles with more than three editors.
4. Include articles with greater than or equal to 50% human editors.

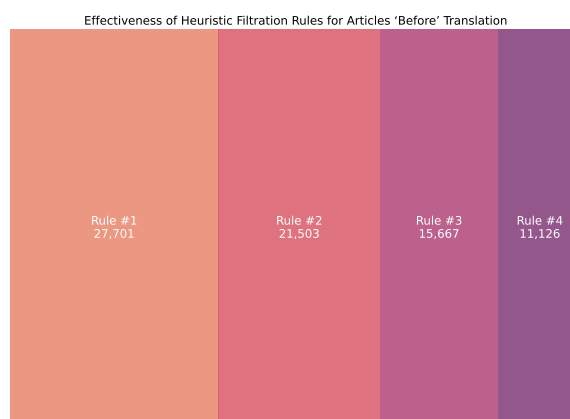


Figure 5: A treemap showing the effectiveness of the heuristic rules for articles *before* the template-based translation in Egyptian Wikipedia, highlighting the number of articles filtered out by each rule.

\* Heuristic filtration rules for *after* the translation:

1. Include articles created between 2019-12-1 and 2023-12-01 and discard young articles with an age of less than 30 days (2023-12-1 and 2024-1-1).
2. Include articles with less than five edits.
3. Include articles with less than three editors.
4. Include articles with greater than or equal to 50% bot editors.
5. Include articles created by these registered users, 'HitomiAkane' and 'Al-Dandoon', who overwhelmed the Egyptian Arabic Wikipedia with massive auto-generated and template-translated articles without human supervision.

Wikipedia	Count	5-Gram
Arabic (AR)	141,880	تصنيف أشخاص على قيد الحياة Classification of surviving people
	80,460	لاعب كرة قدم رجالية مغتربون Men's football players expatriate
	38,793	تعداد عام وبلغ عدد الأسر A general census, and the number of families reached
Egyptian (ARZ)	222,964	صوره هيا مجال الكره السماويه It is a picture of the celestial sphere
	222,961	الكره السماويه اللي المجره جزء The celestial sphere, of which the galaxy is a part
	222,939	مجموعه من النجوم اللي بتكون A collection of stars that forms
Moroccan (ARY)	5,172	width text textcolor black fontsize
	2,057	لعاطلين اللي سبق ليهوم خدمو For unemployed people, who have previously served
	1,483	على حساب لإحصاء الرسمي عام According to the official census of the year

Table 9: Selected top three 5-grams from each Arabic Wikipedia edition with their English translations.

Wikipedia	Count	10-Gram
Arabic (AR)	38,790	تعداد عام وبلغ عدد الأسر أسرة وعدد العائلات عائلة مقيمة A general census, the number of families was one family, and the number of families was one resident family
	38,710	وبلغت نسبة الأزواج القاطنين مع بعضهم البعض من أصل المجموع The percentage of couples living together was out of the total
	38,524	نسبة منها لديها أطفال تحت سن الثامنة عشر تعيش معهم A percentage of them have children under the age of eighteen living with them
Egyptian (ARZ)	222,935	صوره هيا مجال الكره السماويه اللي المجره جزء منها الانزياح A picture of the celestial sphere, of which the galaxy is a part of the displacement
	222,935	المطلع المستقيم هو الزاويه المحصوره بين الدائره الساعيه لجرم سماوي The right ascension is the angle enclosed between the hourly circle of a celestial body
	221,251	أو صوره هيا مجال الكره السماويه اللي المجره جزء منها Or a picture of the celestial sphere, of which the galaxy is a part
Moroccan (ARY)	2,586	imagesize width height plotarea left right top bottom timeaxis orientation
	1,483	لعاداد كان ديالو واصل شخص على حساب لإحصاء الرسمي عام The number of people was counted up to according to the official census of the year
	1,348	ما كايعرفوش يقرأو ولا يكتبو نسبة كان قارين فوق أنوي They did not know how to read or write the percentage of literate was above

Table 10: Selected top three 10-grams from each Arabic Wikipedia edition with their English translations.

<p>توريرت (سيدي أحمد وعبدالله): أرمد هو دوار مجمع كاين جماعة أسني دائرة أسني إقليم لحوز جهة مراكش أسفي لمغرب هاد وار كينتامي مشيخة إلمليل لعاداد كان ديالو واصل شخص على حساب لإحصاء الرسمي عام هو دوار لي كاين الجبل السلسلة ديال لأطلس الكبير الغربي الجغرافيا دوار أرمد بعيد كلم على مدينة مراكش على ارتفاع حوالي ميمرو على البحر فالجبال ديال الأطلس الكبير هاد الدوار مشهور بلفلاحة خصوصا التفاح والكرز فيه بزاف لوبيرجات والمحلات ولعشاش والتجارة السياحية والبيع والمنتجات التقليدية والماكلة السكان إحصائيات عامة عدد السكان ديال أرمد تزداد عدد لفاميلات تزداد ما بين عدد لبالعين كان واحد منهوم دكور</p> <p>imagesize width height plotarea left right top bottom timeaxis orientation vertical alignbars justify colors id gray value gray dateformat yyyy period from till scalemajor unit year increment start gridcolor gray plotdata bar color green width from till width text textcolor black fontsize px bar color red width from till width text textcolor black fontsize px imagesize width height plotarea left right top bottom timeaxis orientation vertical alignbars justify colors id gray value gray dateformat yyyy period from till scalemajor unit year increment start gridcolor gray plotdata bar color green width from till width text textcolor black fontsize px bar color red width from till width text textcolor black fontsize px green</p> <p>الجواج أرمد واصله لومعد ال لعمر عند الجواج اللواني هو عام عند الرجال عند لعيلات لخصوبة عند لعيلات واصله لخصوبة لكاملة التسكويل نسبة التسكويل واصله نسبة لأمية واصله لخدمة نسبة الناس النشيطين دوار أرمد واصله نسبة الشوماج واصله نوطات عيون لكلام تصنيف جهة مراكش أسفي تصنيف دوار لمغرب تصنيف دوار إقليم لحوز تصنيف مقالات فيها مصدر بايت.</p>
--

Table 11: A sample of a parsed article from Moroccan Wikipedia, showing the embedded Wiki markups.

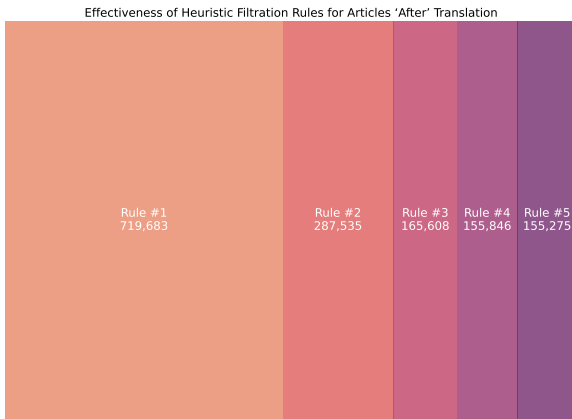


Figure 6: A treemap showing the effectiveness of the heuristic rules for articles *after* the template-based translation in Egyptian Wikipedia, highlighting the number of articles filtered out by each rule.

### C. EGYPTIAN WIKIPEDIA SCANNER

We evaluate our multivariate supervised machine learning classifiers using metrics like accuracy and ROC-AUC (Receiver Operating Characteristic Area Under Curve). We then publicly deploy and host our best classifier, XGBoost, which takes input features of articles' metadata and CAMELBERT embeddings, as illustrated in Figures 7 and 8. We include the articles' metadata because we find that, from our two ablation studies, metadata could be practical and encode features useful for the classifier's learning. We also choose CAMELBERT over Word2Vec word embeddings because CAMELBERT's embeddings take the context into account, and Word2Vec's embeddings are context-free and need to be retrieved word by word and then averaged for the whole article; this is not ideal.

We call this online application EGYPTIAN WIKIPEDIA SCANNER, where users can search for an article directly or select a suggested article retrieved using fuzzy search from the Egyptian Arabic Wikipedia edition. The application automatically fetches the article's metadata (using the Wikimedia XTOOLS API), displays the fetched metadata in a table, and automatically classifies the article as 'human-generated' or 'template-translated'. The application also dynamically displays the full summary of the article and provides the URL to the article to read the full text, as shown in Figure 9.

We utilize the Streamlit Framework<sup>15</sup> to design, host, and deploy the application on the free Streamlit Community Cloud<sup>16</sup> service, making it publicly accessible to everyone at <https://egyptian-wikipedia-scanner.streamlit.app>. We also host the application on Hugging Face Spaces to avoid run-

ning out of Streamlit Cloud free, limited resources: <https://huggingface.co/spaces/SaiedAlshahrani/Egyptian-Wikipedia-Scanner>. This online application, EGYPTIAN WIKIPEDIA SCANNER, is open-sourced on GitHub with an MIT license, here: <https://github.com/SaiedAlshahrani/Egyptian-Wikipedia-Scanner>.

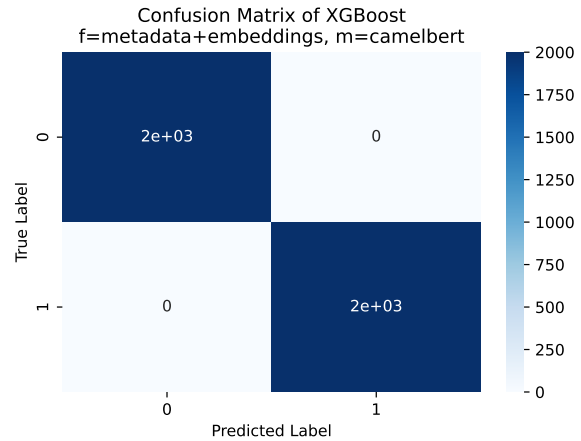


Figure 7: Confusion matrix of the best, deployed classifier, XGBoost, which takes input features of articles' metadata combined with CAMELBERT's embeddings, showing the excellent performance of this multivariate ensemble classifier.

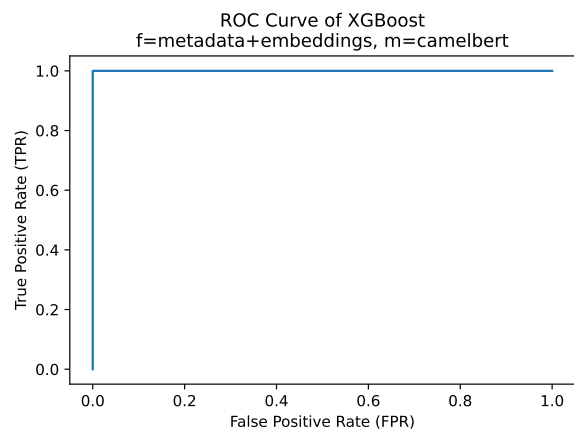


Figure 8: ROC curve of the best, deployed classifier, XGBoost, which takes input features of articles' metadata combined with CAMELBERT's embeddings, showing the excellent performance of this multivariate ensemble classifier.

<sup>15</sup>Streamlit Framework: <https://streamlit.io>.

<sup>16</sup>Streamlit Cloud: <https://streamlit.io/cloud>.

# Egyptian Arabic Wikipedia Scanner

## Automatic Detection of Template-translated Articles in the Egyptian Wikipedia

Search for an article in Egyptian Arabic Wikipedia:

ويكيبيديا مصرى

### ■ Collected Metadata of ويكيبيديا مصرى

Total Edits	Total Editors	Total Bytes	Total Characters	Total Words	Creator Name	Creation Date
242	37	3,929	2,222	388	Ghaly	2008-05-01


### ■ Automatic Classification of ويكيبيديا مصرى

Human-generated Article

### ■ Full Summary of ويكيبيديا مصرى

ويكيبيديا مصرى

ويكيبيديا مصرى مشروع موسوعه حره اى حد ممكن يساهم فى كتابتها و مكتوبه بالمصرى بطريقه اى مصرى يعرف يقرأها. ويكيبيديا مصرى هى النسخه المصرى بتاعه ويكيبيديا, الموسوعه الحره. ويكيبيديا مصرى فيها 1,622,097 مقاله دلوقتى. ف يونيه 2020, ويكيبيديا مصرى كانت تالت لغه ف ويكيبيديا بيزورها يوزرز من مصر 961,000 قراية صفحه

 Read Full Text of ويكيبيديا مصرى:

[https://arz.wikipedia.org/wiki/%D9%88%D9%8A%D9%83%D9%8A%D8%A8%D9%8A%D8%AF%D9%8A%D8%A7\\_%D9%85%D8%B5%D8%B1%D9%89](https://arz.wikipedia.org/wiki/%D9%88%D9%8A%D9%83%D9%8A%D8%A8%D9%8A%D8%AF%D9%8A%D8%A7_%D9%85%D8%B5%D8%B1%D9%89)

Figure 9: A screenshot of the EGYPTIAN WIKIPEDIA SCANNER, illustrating its capabilities and features.