# *ANETAC*: Arabic Named Entity Transliteration and Classification Dataset

**Mohamed Seghir Hadj Ameur**[*]
Department of Computer Science
USTHB University
Bab-Ezzouar, Algiers, Algeria
mhadjameur@usthb.dz

**Farid Meziane**
Informatics Research Centre
University of Salford
M5 4WT, United Kingdom
f.meziane@salford.ac.uk

**Ahmed Guessoum**
Department of Computer Science
USTHB University
Bab-Ezzouar, Algiers, Algeria
aguessoum@usthb.dz

July 9, 2019

## Abstract

In this paper, we make freely accessible *ANETAC*[1] our English-Arabic named entity transliteration and classification dataset that we built from freely available parallel translation corpora. The dataset contains $79,924$ instances, each instance is a triplet $(e, a, c)$, where $e$ is the English named entity, $a$ is its Arabic transliteration and $c$ is its class that can be either a Person, a Location, or an Organization. The *ANETAC* dataset is mainly aimed for the researchers that are working on Arabic named entity transliteration, but it can also be used for named entity classification purposes. This dataset was developed and used as part of a previous research study done by Hadj Ameur et al. [1].

***Keywords*** Natural Language Processing · Arabic Language · Arabic Transliteration · Named Entity Transliteration · Arabic Named Entity · Arabic Transliteration Dataset

## 1 Introduction

The task of transliteration is the process of converting words (e.g. named entities) that are written in one language alphabet to another language that has a different alphabet while still preserving the phonetics of the transliterated words. One of the main difficulties when attempting to transliterate named entities from a given source language to another is the lack of some phonetic character correspondences. For example, in the task of named entity transliteration between Arabic and English, several Arabic letters such as "ث"and "ظ" do not have direct single-letter correspondences in the English language alphabet. Table 1 presents some English named entities and their transliteration in the Arabic language.

Table 1: English named entities and their equivalent Arabic transliterations

| English | Arabic |
|---------|--------|
| Brandes | برانديس (Brandees) |
| Mayhawk | مايهوك (Mayhouk) |
| Cressner | كريسنير (Crissneer) |
| Husseini | حسيني (Husseini) |

Accurate transliteration of named entities is useful for several applications such as machine translation [2, 3], and cross-lingual information retrieval [4, 5]. Though a great deal of attention has been devoted to improving this task for

---

[*]Corresponding author. Feel free to contact me via my personal email mohamedhadjameur@gmail.com
[1]The ANETAC dataset is freely available on Github https://github.com/MohamedHadjAmeur/ANETAC.

many languages such as English, only limited studies have been made with regard to Arabic mainly due to the lack of transliteration datasets. In this paper, we make accessible *ANETAC*, an English-Arabic named entity transliteration and classification dataset that we built from freely available parallel translation corpora. It contains 79,924 English-Arabic named entities along with their respective classes that can be either a Person, a Location, or an Organization. Table 2 shows statistics about the *ANETAC* named entities classes.

Table 2: Statistics about the number of named entities belonging to each class [1]

| Named entity | Count |
|---|---|
| Person | 61,662 |
| Location | 12,679 |
| Organization | 5,583 |
| **All** | 79,924 |

To make it easier for other researchers to train and compare their own models, the *ANETAC* dataset is divided into training, development, and test sets as shown in Table 3.

Table 3: Instance counts in the train, development and test datasets of our transliteration corpus [1]

| Sets | Train | Dev | Test |
|---|---|---|---|
| Named entities count | 75,898 | 1004 | 3013 |

As pointed out by many recent studies [1, 6], there is a lack of Arabic machine transliteration datasets. To the best of our knowledge, there is only one freely available English-Arabic transliteration dataset that contains no more than 12,877 pairs [2], thus, we believe that our dataset will be a valuable addition. The importance of the *ANETAC* dataset can be summarized as follows:

- This dataset is useful for many applications such as (1) training state-of-the-art English-Arabic machine transliteration models, (2) training Arabic named entity classification models, (3) handling Out-Of-Vocabulary (OOV) words in machine translation, (4) dealing with proper names in Cross-lingual Information Retrieval.

- This dataset is mainly aimed for those researchers working on Arabic named entity transliteration, but it can also be used for named entity classification purposes.

- This dataset also contains a test set that can be used as a benchmark to compare the results of English-Arabic transliteration systems. First transliteration results have been already reported on this test set by Hadj Ameur et al. [1] and will be shown in Section 3.

In the remainder of this paper, section 2 presents the corpus construction methodology that we adopted in the development of this dataset. Section 3 presents the baseline transliteration results that have been obtained using the *ANETAC* dataset. Finally, section 4 provides a conclusion to this paper.

## 2    Building a Transliteration Corpus

As stated in the original work of Hadj Ameur et al. [1][3], the extraction system (see Fig. 1) uses freely available parallel corpora[4] in order to automatically extract bilingual named entities. The English-Arabic corpora that we have used are provided in Table 4.

Table 4: Statistics about the used English-Arabic parallel corpora [1]

| Corpus | Sentences (in millions) |
|---|---|
| United Nation | 10.6M |
| Open Subtitles | 24.4M |
| News Commentary | 0.2M |
| IWSLT2016 | 0.2M |
| **All** | 35.4M |

---

[2]https://github.com/google/transliteration

[3]We note that this description of the extraction system is mostly based on the original paper of Hadj Ameur et al. [1].

[4]The English-Arabic parallel corpora that we used are available on the opus website: http://opus.nlpl.eu.

As shown in Fig. 1, the system starts by a preprocessing phase in which the English and Arabic sentences are tokenized and normalized. Then, the English named entities are identified in each sentence belonging to the English-side of the parallel corpus. A set of Arabic transliteration candidates will then be associated with each English named entity. Finally, the best Arabic transliteration candidate will be selected for each English named entity. The detail of these step are provided in the remainder of this section.
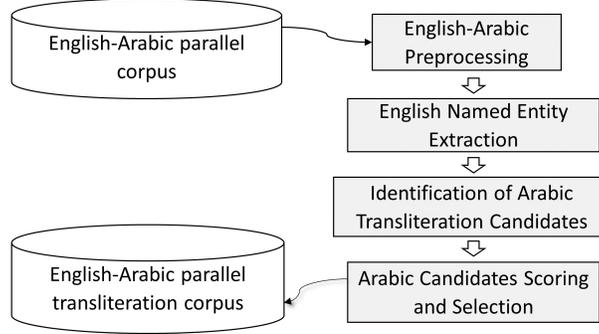


Figure 1: Architecture of our parallel English-Arabic Named entity extraction system [1]

## 2.1 Parallel Named Entity Extraction

The ultimate goal is to extract the correct Arabic transliteration of each English named entity. Given a corpus of English-Arabic parallel sentences $S = \{(e_1, a_1), ..., (e_m, a_m)\}$, we use the Stanford English Named Entity Recognizer [7] to find all the English named entities that are present in the parallel corpus $E_{ne} = \{n_1, n_2, ..., n_k\}$, where $k$ is the total number of named entities. Since each singleton word belonging to a multi-word English named entity can always be transliterated solely without needing its context, we decomposed all the English named entities containing multiple words to several singleton entities. For each English named entity $n_i$ belonging to an English sentence $e_j$, we end up with a list of pairs $(n_i, a_j)$ denoting that the $i^{th}$ English named entity (singleton word) is associated with the $j^{th}$ Arabic sentence.

## 2.2 Candidates Extraction and Scoring

The previous step leaves us with a set of pairs $(n_i, a_j)$, where $n_i$ is the English named entity (word) and $a_j$ is the Arabic sentence containing its transliteration. To find the correct transliterated word of $n_i$ in the Arabic sentence $a_j$, we first removed all the frequent Arabic words from it using a vocabulary containing the top $n$ most frequent Arabic words, with $n = 40000$, that we built automatically from our parallel corpus. This ensures that the remaining words in the Arabic sentence $a_j$ are mostly rare words. All the remaining words in $a_j$ are considered as transliteration candidates $C(a_j) = \{c_{j1}, c_{j2}, ..., c_{jt}\}$, where $c_{ji}$ denotes the $i^{th}$ candidate word found in the $j^{th}$ Arabic sentence, and $t$ is the total number of Arabic candidates in $C(a_j)$. We used the transliteration tool available in the polyglot multilingual NLP library[5] to obtain an approximate Arabic transliteration $t_i$ of each English named entity $n_i$. For each English named entity $n_i$ having the approximate transliteration $t_i$ and the list of Arabic candidates $C(a_j)$, the score of each Arabic candidate is estimated using the following three features:

1. The total number of shared characters: this feature takes into account the count of shared characters between each Arabic candidate in $C(a_j)$ and the approximate transliteration $t_i$.

2. The longest shared sequence: this feature takes into account the length of the longest common sequence of characters between each Arabic candidate in $C(a_j)$ and the approximate transliteration $t_i$.

3. Length difference penalty: this feature is used to penalize the $C(a_j)$ candidates according to their level of dissimilarity with the approximate transliteration $t_i$.

The final score of each candidate is then estimated by averaging the score of all the three features. The candidate having the highest score is then selected if its corresponding final score surpasses a certain confidence threshold. Some examples of the extracted English-Arabic named entities are provided in Table 5. The reader should recall that the Arabic language has no letters for the English sound "v", "p" and "g".

---

[5] https://github.com/aboSamoor/polyglot

Table 5: Some examples of the extracted English-Arabic named entities [1]

| Entity class | English | Arabic |
|---|---|---|
| PERSON | Villalon | فيلالون (filaloun) |
| LOCATION | Nampa | نامبا (namba) |
| ORGANIZATION | Soogrim | سوغريم (soughrim) |

## 3   Baseline Results

This section provides the English-to-Arabic and Arabic-to-English baselines' transliteration results that we have obtained when using the *ANETAC* dataset for both the training and testing of our models [1]. The baseline results (Table 6) are reported in terms of both Word Error Rate (WER) and Character Error Rate (CER) on the *ANETAC* test set[6].

Table 6: Baseline transliteration results in terms of WER and CER reported on the *ANETAC* test set

| Directions | WER | CER |
|---|---|---|
| English-to-Arabic | **5.40** | **0.95** |
| Arabic-to-English | **65,16** | **16.35** |

As shown in Table 6, the results of the Arabic-to-English transliteration are still poor, thus much work is still needed to improve them. We note the baseline models that we have used are based on the attention-based encoder-decoder architecture [8] and trained at the character level.

## 4   Conclusion

In this work, we have made accessible the *ANETAC* dataset, that we developed as part of our previous work [1]. We have shown how this dataset is built from parallel translation corpora by relying on several features and tools. We also presented the baseline results that we have achieved on the tasks of English-to-Arabic and Arabic-to-English machine transliteration. We encourage all researchers that are interested in this task to try and achieve better results. Finally, we hope that this dataset will have a positive impact on the current state of Arabic-English named entity transliteration.

## References

[1] Mohamed Seghir Hadj Ameur, Farid Meziane, and Ahmed Guessoum. Arabic machine transliteration using an attention-based encoder-decoder model. *Procedia Computer Science*, 117:287–297, 2017.

[2] Ulf Hermjakob, Kevin Knight, and Hal Daumé III. Name translation in statistical machine translation-learning when to transliterate. In *ACL*, pages 389–397, 2008.

[3] Nizar Habash. Four techniques for online handling of out-of-vocabulary words in arabic-english statistical machine translation. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics on Human Language Technologies: Short Papers*, pages 57–60. Association for Computational Linguistics, 2008.

[4] Paola Virga and Sanjeev Khudanpur. Transliteration of proper names in cross-lingual information retrieval. In *Proceedings of the ACL 2003 workshop on Multilingual and mixed-language named entity recognition-Volume 15*, pages 57–64. Association for Computational Linguistics, 2003.

[5] Atsushi Fujii and Tetsuya Ishikawa. Japanese/english cross-language information retrieval: Exploration of query translation and transliteration. *Computers and the Humanities*, 35(4):389–420, 2001.

[6] Mihaela Rosca and Thomas Breuel. Sequence-to-sequence neural network models for transliteration. *arXiv preprint arXiv:1610.09565*, 2016.

[7] Jenny Rose Finkel, Trond Grenager, and Christopher Manning. Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43rd annual meeting on association for computational linguistics*, pages 363–370. Association for Computational Linguistics, 2005.

[8] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.

---

[6]`https://github.com/MohamedHadjAmeur/ANETAC`