

Neural and Symbolic Arabic Paraphrasing with Automatic Evaluation

First Author

Affiliation / Address line 1
Affiliation / Address line 2
Affiliation / Address line 3
email@domain

Second Author

Affiliation / Address line 1
Affiliation / Address line 2
Affiliation / Address line 3
email@domain

Weijian Lin

Carnegie Mellon University
School of Computer Science
5000 Forbes Avenue, Pittsburgh, PA 15213
wlin1@cs.cmu.edu

Abstract

We present a tool for Arabic paraphrasing that yields good paraphrasing accuracy. We present and compare several methods for paraphrasing and obtaining monolingual parallel data. We also present first results on Arabic paraphrasing using neural methods. Additionally, we propose a new evaluation metric for paraphrasing that is shown to correlate highly with human judgement.

1 Introduction

2 Related Work

3 Obtaining Parallel Monolingual Data

4 Extracting Paraphrases

In order to extract paraphrases, we first obtained a parallel bilingual corpus using the EUROPARL dataset (Koehn, 2005). (Bannard and Callison-Burch, 2005)

5 Generating Paraphrased Sentences

To address the shortage problem for Arabic language data, here we propose a novel method for generating Arabic parallel language sentences using other two languages pairs as resource. The idea is to use paired sentences data from two other languages, translate them into Arabic correspondingly, then use them as the training data for sequence to sequence machine translation model to train a translator to generate Arabic to Arabic paraphrases.

The first advantage of this approach is its scalability. After preparing enough training data for the

paraphraser model, the generation step for Arabic paraphrases is easily scalable. The second advantage is that paraphraser model may contain valuable insight for building Arabic paraphrases database at words and phrases level, since state of the art neural machine translation techniques are capable of capturing word and phrase level similarity by projecting word embeddings into vector space.

We use europarl-v7.fr-en data which contains two millions sentence pairs, and then use Google translate API to generate French-Arabic and English-Arabic sentence pairs correspondingly. Then we use them as training data to train a single layer sequence to sequence attentional encoder-decoder model, with embedding size and hidden size set to be 512 while attention size set to be 128. The whole process is demonstrated in the following ??

5.1 Phrase Substitution Method

5.2 Neural Seq-to-Seq Method

6 An Automatic Evaluation Metric

6.1 Semantic Similarity

6.2 Surface Variation

7 Analysis

8 Evaluation

9 Future Work

We plan to explore other options for obtaining monolingual parallel data. One possible approach is to retrieve headlines of news articles from different agencies covering the same event. We expect headlines describing the same event to have some degree

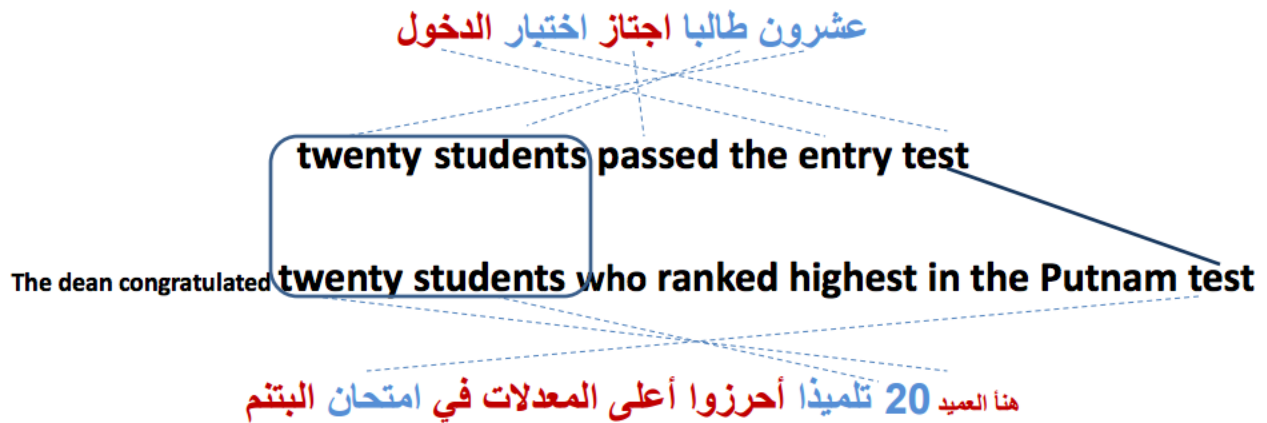


Figure 1: An example paraphrased sentence produced using the pivot method.

of semantic similarity yet different surface realizations.

The sequence-to-sequence models required a relatively long time to run which limited the testing of other architectures and model options. We plan to conduct more experiments on different architectures and compare results.

Acknowledgments

References

- Colin Bannard and Chris Callison-Burch. 2005. Paraphrasing with bilingual parallel corpora. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics, ACL '05*, pages 597–604, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. <http://www.statmt.org/europarl/>.