

Dual-FOFE-net Neural Models for Entity Linking with PageRank

Feng Wei (✉), Uyen Trang Nguyen, and Hui Jiang

Department of Electrical Engineering and Computer Science,
York University, 4700 Keele St, Toronto, ON M3J 1P3, Canada
{fwei,utn,hj}@cse.yorku.ca

Abstract. This paper presents a simple and computationally efficient approach for entity linking (EL), compared with recurrent neural networks (RNNs) or convolutional neural networks (CNNs), by making use of feedforward neural networks (FFNNs) and the recent dual fixed-size ordinally forgetting encoding (dual-FOFE) method to fully encode the sentence fragment and its left/right contexts into a fixed-size representation. Furthermore, in this work, we propose to incorporate PageRank based distillation in our candidate generation module. Our neural linking models consist of three parts: a PageRank based candidate generation module, a dual-FOFE-net neural ranking model and a simple NIL entity clustering system. Experimental results have shown that our proposed neural linking models achieved higher EL accuracy than state-of-the-art models on the TAC2016 task dataset over the baseline system, without requiring any in-house data or complicated handcrafted features. Moreover, it achieves a competitive accuracy on the TAC2017 task dataset.

Keywords: Neural network · Entity linking · Knowledge base.

1 Introduction

Named entities (NEs) have received much attention over the last two decades [17], mostly focused on recognizing the boundaries of textual NE mentions and classifying them as, e.g., Person (PER), Organization (ORG), Facility (FAC), Geo-political Entity (GPE) or Location (LOC). In 2009, NIST proposed the shared task challenge of EL [15]. EL is a similar but broader task than named entity disambiguation (NED). NED is concerned with disambiguating a textual NE mention where the correct entity is known to be one of the Knowledge Base (KB) entries, while EL also requires systems to deal with the case where there is no entry for the NE in the reference KB.

In [10], the authors group and summarise different approaches to EL taken by participating systems. There is a vast body of research on NED, highlighted by [9]. The problem has been studied extensively by employing a variety of machine learning, and inference methods, including a pipeline of deterministic modules [13], simple classifiers [22], graphical models [3], classifiers augmented with ILP inference [1], and more recently, neural approaches [14,23,24,29].

In this paper, we propose to use the recent dual-FOFE method [25] to fully encode the left/right contexts for each target mention, and then a simple FFNNs can be trained to make a precise linking for each target mention based on the fixed-size presentation of the contextual information. Moreover, we propose to incorporate PageRank based distillation in our candidate generation system. Compared with [14,24,29], our proposed neural linking models, without requiring any in-house data or complicated handcrafted features, can efficiently achieve higher EL accuracy in terms of computing than the baseline system, and achieves a competitive accuracy on both TAC2016 and TAC2017 task datasets.

The remainder of this paper is organized as follows. Section 2 describes our proposed neural linking models. In Section 3, we discuss experimental results and compare the performance of our proposed models with that of existing state-of-the-art systems. Finally, Section 4 draws the conclusions and outlines our future work.

2 Our Proposed Neural Linking Models

In this section, we discuss our proposed neural linking models, which consist of three parts: PageRank based candidate generation module, dual-FOFE-net neural ranking model and NIL entity clustering system.

2.1 PageRank based Candidate Generation

Inspired by [19], we propose to extend the previous work in [14] to incorporate a PageRank based distillation in our candidate generation module to generate candidates for each detected mention. Candidates are generated based on KBs, including *Freebase* and *Wikipedia* [30]. Lucene fuzzy search strategy is applied in the implementation. The input to this module is a detected mention, and the output is a candidate list, which consists of a group of *Freebase* nodes potentially matching this mention, as shown in Figure 1.

Following are five types of mention extensions implemented in the candidate generation module:

- *Substring Extension*: For each mention, all the recognized named entities in its original context document containing that mention will be selected. For instance, given the mention “Trump” in document d , “Donald Trump” will be selected as its substring extension if the named entity “Donald Trump” is found in d .
- *Translation Extension*: If a mention is in Chinese or Spanish, we invoke Google Translation to obtain its English translation as Translation Extension.
- *Country Extension*: The abbreviation of a country name can be extended to a more concrete one. For example, the mention of the geo-political entity “UK” will be extended to “United Kingdom”.
- *Nominal Extension*: The nearest recognized entity with the same entity type as its nominal extension will be selected to be added to the query list.

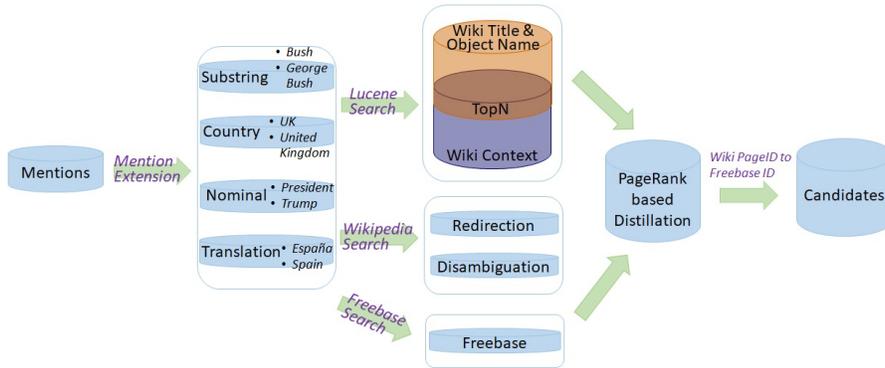


Fig. 1. The diagram of the PageRank based candidate generation module.

- *Traditional Chinese Extension*: If a mention is in the form of simplified Chinese, its traditional Chinese version will be obtained.

After the mention extensions, three parallel strategies are applied respectively:

- To invoke Lucene fuzzy search on *Wikipedia* titles, first paragraphs, document context and *Freebase* object names.
- To query a database with *Wikipedia* redirection and disambiguation information.
- To query a database with all *Freebase* entities.

PageRank based Distillation. As discussed above, although in general the mention extension step helps to enhance the candidate coverage, it also produces too many candidates for a single mention. This behavior leads to much noise, and slows down the whole system.

In this work, inspired by [19], we propose to incorporate PageRank based distillation at the last step. As depicted in Figure 2, a toy document graph includes three entity mentions and seven candidates: three candidates generated for Lincolnshire, and two candidates generated for United F.C. and Devon White each. Each graph node $e(m, c)$ is a pair of an entity mention m and a candidate c . An edge is drawn between two candidates of different entities whenever there is a link from the *Wikipedia* page for one candidate to the *Wikipedia* page for the other. There is no edge between candidates competing for the same entity.

It is worth noting that edges in our graph model represent relations between candidates. We insert an edge between two candidates if the *Wikipedia* entry corresponding to either of the two candidates contains a link to the other candidate. We assume that this relation is bidirectional and thus this edge is undirected.

We rank the candidates of each mention based on their outbound link counts to all the recognized mentions in the same document, and keep the top τ can-

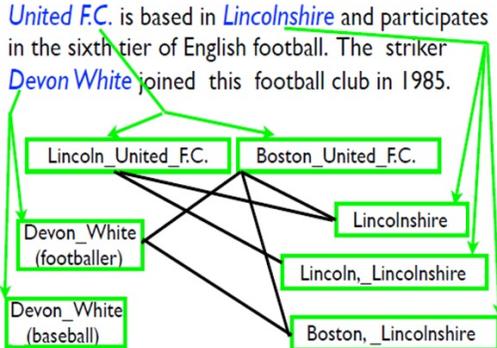


Fig. 2. A toy document graph for three entity mentions: United F.C., Lincolnshire, Devon White. Source: Adapted from Pershina et al. [21]

didates for each mention, where τ is a distillation factor. We name this step “distillation”. The score is calculated by the *Wikipedia*’s anchors as follows:

$$score(c) = \sum_{m'} count(c, m') \quad (1)$$

where c is the candidate entity, m' is the linked page (a mention identified in the same document) and $count(c, m')$ is the total co-occurrence count of c and m' .

Finally, for each detected mention m , the candidate generation module generates a list of K candidates as $C = \{c_1, \dots, c_K\}$, with $K \in [0, \tau]$.

2.2 Dual-FOFE-net Neural Ranking Models

Dual Fixed-Size Ordinally Forgetting Encoding (Dual-FOFE). FOFE [31] was proposed as an alternative to commonly used sequence embedding representations, and achieved competitive results in language modeling. There is a nice theoretical property to guarantee that FOFE codes can almost uniquely encode any variable-length sequence of words into a fixed-size representation without losing any information.

Given a vocabulary V , where each word can be represented by a 1-of- $|V|$ one-hot vector. Let $S = \{w_1, \dots, w_N\}$ denote a sequence of N words from V , and e_n denote the one-hot vector of the n -th word in S , where $1 \leq n \leq N$. Assuming $z_0 = 0$, the FOFE code z_n of the sequence from word w_1 to w_n is shown as follows:

$$z_n = \alpha \cdot z_{n-1} + e_n \quad (2)$$

where α is a constant forgetting factor. Thus, z_n can be viewed as a fixed-size representation of the subsequence $\{w_1, \dots, w_n\}$. We can see that, according to the theoretical properties presented in [31], any sequence of variable length can be uniquely and losslessly encoded into a fixed-size representation by FOFE.

This simple ordinally-forgetting mechanism has been applied to some NLP tasks, e.g., [26,27,28] and have achieved very competitive results.

The main idea of dual-FOFE is to generate augmented FOFE encoding codes by concatenating two FOFE codes using two different forgetting factors. Each of these FOFE codes is still computed in the same way as the mathematical formulation shown in Equation (2). The difference between them is that we may select to use two different values for the forgetting factor for additional modeling benefits.

Our Proposed Dual-FOFE-net. As described in 2.1, we generate a candidate list C for each detected mention m . This list contains a special NIL candidate and some *Freebase* node IDs that match the mention in the candidate generation process. In this work, we propose to use a FFNNs probability ranking model to assign probabilities to all candidates in the list. The candidate with the highest probability is chosen as the final linking result. Each time, the FFNNs probability ranking model takes a mention m and a candidate c_k from the list C to compute a matching score, e_k . In order to do this, we make use of dual-FOFE to encode mention context features for the neural network.

As shown in Figure 3, the input feature vector to the FFNNs probability ranking model is a concatenation of all the following features:

- *Mention string embedding*: Each detected mention is represented as a bag-of-words vector. This bag-of-words vector is projected into a 128-dimension dense vector.
- *Document context*: The left and right contexts of each mention are encoded by dual-FOFE, and projected into a 256-dimension dense vector.
- *Knowledge base description*: The corresponding KB, *Freebase*, description of each candidate and target mention is individually represented as one bag-of-words vectors (weighted using the TFIDF schema), which is mapped to a 128-dimension dense vector. As for Chinese and Spanish, since the languages have fewer resources than English in *Freebase*, we invoke *Google* APIs, which extract the translation to expand their Chinese and Spanish descriptions separately.

In this paper, we use the rectified linear activation function, i.e., $f(x) = \max(0, x)$, to compute from activations to outputs in each hidden layer, which are in turn fed to the next layer as inputs. For the output layer, we make use of the softmax function to compute posterior probabilities between two nodes, standing for correct links or incorrect links, shown as follows:

$$P_r(c_k|m) = \frac{\exp(e_k)}{\sum_{k=1}^K \exp(e_k)}. \quad (3)$$

2.3 NIL Entity Clustering

For all mentions identified as NIL by the above dual-FOFE-net neural ranking model, we perform a simple rule-based algorithm to cluster them: Different

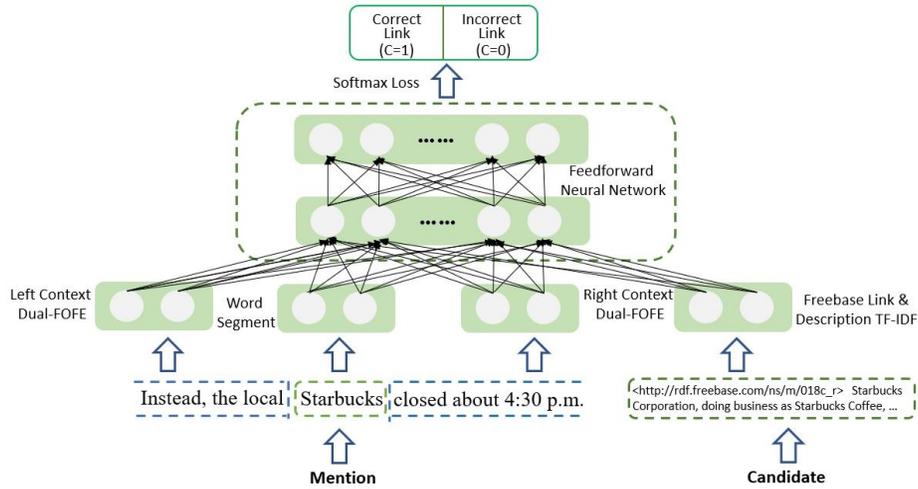


Fig. 3. Our proposed neural ranking model using dual-FOFE codes as input and a feed-forward neural network.

named NIL mentions are grouped into one cluster only if their mention strings are the same (case-insensitive).

3 Experiments and Results

In this section, we evaluate the effectiveness of our proposed methods on the benchmark datasets: TAC2016 and TAC2017 Trilingual Entity Discovery and Linking (EDL) tasks, and compare the performance of our proposed neural linking model with state-of-the-art models [14,29] on both TAC2016 and TAC2017 task datasets.

In [14], the authors use some in-house data annotated by themselves, which consists of about 10,000 Chinese and English documents acquired through their web crawler. These documents are internally labelled using some annotation rules similar to the KBP guidelines. In [29], many complicated handcrafted features are created, including the mention level feature, entity level feature, mention-to-entity feature and entity-to-entity feature. It is worth noting that we have not used any in-house data or handcrafted features in our models, and all used features (either word or character level) are automatically derived from the data based on the simple FOFE formula.

3.1 Dataset

Given a document collection in three languages (English, Chinese and Spanish), the TAC trilingual EDL task [11,12] automatically identifies entities from a source collection of textual documents in multiple languages, as shown in Table

Table 1. Number of Documents in TAC2015-2017

	English	Chinese	Spanish	ALL
15 Train	168	147	129	444
15 Eval	167	166	167	500
16 Eval	168	167	168	503
17 Eval	167	167	166	500

1, classifies them into one of the following pre-defined five types: Person (PER), Geo-political Entity (GPE), Organization (ORG), Location (LOC) and Facility (FAC), links them to an existing KB (BaseKB)¹, and clusters mentions for those NIL entities that do not have corresponding BaseKB entries. The corpus consists of news articles and discussion forum posts published in recent years, related to but non-parallel across languages.

3.2 Neural Model Setup

Three models are trained and evaluated independently. Three sets of word embeddings of 128 dimensions are derived from English Gigaword [20], Chinese Gigaword [5] and Spanish Gigaword [16], respectively. Since Chinese segmentation is not reliable, based on a predefined set of all possible characters, we view the focus token as a character sequence and encode it using dual-FOFE. We then project the character encodings using a trainable character embedding matrix.

Relying on the development experiments, the set of hyper-parameters used in our experiments is summarized as follows: i) Learning rate: All models are trained using the stochastic gradient descent (SGD) algorithm while the learning rate is set to be 0.1; ii) Network structure: Three hidden layers and ReLUs [18] as the nonlinear activation function, randomly initialized based on a uniform distribution between $-\sqrt{\frac{6}{N_i+N_o}}$ and $\sqrt{\frac{6}{N_i+N_o}}$ [4]; iii) Dropout [6] is adopted during training; iv) Number of epochs: 30; v) Chinese character embeddings: 64 dimensions, randomly initialized; vi) Forgetting factor: $\alpha = (0.5, 0.9)$ ²; vii) Distillation factor: $\tau = 20$.

3.3 Evaluation Metrics

To evaluate the effectiveness of our proposed models, we use the standard NERLC and CEAFmC metrics, which are combined measures of linking and clustering performance.³

¹ <http://basekb.com/>

² The choice of the forgetting factors α is empirical. We’ve evaluated on a development set in some early experiments. It turns out that $\alpha = (0.5, 0.9)$ is the best. As a result, $(0.5, 0.9)$ is used for all EL tasks throughout this paper.

³ More details regarding data format and scoring metric can be found in <http://nlp.cs.rpi.edu/kbp/>

Table 2. Performance comparison with the best system on the TAC2016 datasets (in terms of NERLC F_1 and CEAFmC F_1).

	[14] TAC Rank 1		Our proposed models	
	NERLC	CEAFmC	NERLC	CEAFmC
Trilingual	64.7	66.0	65.9	67.1
English	66.6	67.6	67.7	69.0
Chinese	65.0	70.2	66.4	70.7
Spanish	61.6	63.5	62.5	64.4

3.4 Results and Discussion

Table 2 shows the performance of our proposed model on the TAC2016 task dataset along with the TAC Rank 1 system [14]. Our model outperforms the other system by 1.2% in terms of NERLC, and by 1.1% in terms of CEAFmC with the overall trilingual EL performance. Furthermore, each of the three individual models is better than its counterpart in the TAC Rank 1 system in terms of both NERLC and CEAFmC. As for the TAC2017 task dataset, shown in Table 3, encouragingly, the NERLC performance of English, Chinese and trilingual overall, outperforms the best system [29], by 0.4%, 0.6% and 0.2% separately, and the CEAFmC performance of English, Spanish and trilingual overall, is slightly better than the best system [29], by 0.5%, 0.4% and 0.4% separately.

Following are the advantages of our proposed models over the state-of-the-art on both TAC2016 and TAC2017 task datasets. First, unlike the systems in [14,29], our models do not rely on any in-house data or complicated handcrafted features. It is very time consuming and labour intensive to prepare clean annotated in-house data, or to collect and select good handcrafted features. More importantly, we have not used any handcrafted features in our models, and all used features (either word or character level) are automatically derived from the data based on the simple FOFE formula. Secondly, we present a simple and computationally efficient approach compared with recurrent neural networks (RNNs) or convolutional neural networks (CNNs), by making use of feedforward neural networks (FFNNs) and the recent dual fixed-size ordinal forgetting encoding (dual-FOFE) method to fully encode the sentence fragment and its left/right contexts into a fixed-size representation. Feedforward neural networks (FFNNs) use rather simple structures consisting of several fully-connected layers. These neural networks are known to be powerful as universal approximators [8], and they are simpler and faster to train and inference than the more recent variants such as long short-term memory (LSTM) [7], Gated Recurrent Unit (GRU) [2] or CNNs. Thus, our proposed dual-FOFE-net neural model is light and highly efficient compared with RNNs or CNNs. Last but not least, our proposed PageRank based distillation not only enhances the candidate coverage, but also speeds up the whole models.

Table 3. Performance comparison with the best system on the TAC2017 datasets (in terms of NERLC F_1 and CEAFmC F_1).

	[29] TAC Rank 1		Our proposed models	
	NERLC	CEAFmC	NERLC	CEAFmC
Trilingual	67.8	70.5	68.0	70.9
English	66.8	68.8	67.2	69.3
Chinese	71.0	73.2	71.6	72.4
Spanish	65.0	68.9	64.8	69.3

4 Conclusion

This paper presents a simple and computationally efficient approach to EL by applying FFNNs on top of dual-FOFE features. Furthermore, we propose to incorporate PageRank based distillation in our candidate generation module. Our experiments have shown that, without requiring any in-house data or complicated handcrafted features, it achieves higher EL accuracy than state-of-the-art systems on the TAC2016 task dataset, and offers a competitive accuracy on the TAC2017 task dataset.

In our future work, we will evaluate our neural linking models on more datasets and conduct more experiments to measure the sensitivity of the system to the values of some hyperparameters (e.g., number of hidden layers). In addition, we will explore more architectures (e.g., convolutional layers), to quantify the contribution of some modules.

References

1. Cheng, X., Roth, D.: Relational inference for wikification. In: Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing. pp. 1787–1796 (2013)
2. Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078 (2014)
3. Durrett, G., Klein, D.: A joint model for entity analysis: Coreference, typing, and linking. Transactions of the Association for Computational Linguistics **2**, 477–490 (2014)
4. Glorot, X., Bordes, A., Bengio, Y.: Deep sparse rectifier neural networks. In: Proceedings of the fourteenth international conference on artificial intelligence and statistics. pp. 315–323 (2011)
5. Graff, D., Chen, K.: Chinese gigaword. LDC Catalog No.: LDC2003T09, ISBN **1**, 58563–58230 (2005)
6. Hinton, G.E., Srivastava, N., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.R.: Improving neural networks by preventing co-adaptation of feature detectors. arXiv preprint arXiv:1207.0580 (2012)
7. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural computation **9**(8), 1735–1780 (1997)

8. Hornik, K.: Approximation capabilities of multilayer feedforward networks. *Neural networks* **4**(2), 251–257 (1991)
9. Ji, H.: Entity discovery and linking reading list. <http://nlp.cs.rpi.edu/kbp/2014/elreading.html> (2016)
10. Ji, H., Grishman, R.: Knowledge base population: Successful approaches and challenges. In: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*. pp. 1148–1158. Association for Computational Linguistics (2011)
11. Ji, H., Nothman, J., Dang, H.T., Hub, S.I.: Overview of tac-kbp2016 tri-lingual edl and its impact on end-to-end cold-start kbp. *Proceedings of TAC* (2016)
12. Ji, H., Pan, X., Zhang, B., Nothman, J., Mayfield, J., McNamee, P., Costello, C., Hub, S.I.: Overview of tac-kbp2017 13 languages entity discovery and linking. In: *Proceedings of the Tenth Text Analysis Conference (TAC2017)* (2017)
13. Ling, X., Singh, S., Weld, D.S.: Design challenges for entity linking. *Transactions of the Association for Computational Linguistics* **3**, 315–328 (2015)
14. Liu, D., Lin, W., Wei, S., Zhang, S., Jiang, H.: The ustc nelslip systems for trilingual entity detection and linking tasks at tac kbp 2016. In: *TAC* (2016)
15. McNamee, P., Dang, H.T.: Overview of the tac 2009 knowledge base population track. In: *Text Analysis Conference (TAC)*. vol. 17, pp. 111–113 (2009)
16. Mendonca, A., Graff, D.A., DiPersio, D.: *Spanish gigaword second edition*. Linguistic Data Consortium (2009)
17. Nadeau, D., Sekine, S.: A survey of named entity recognition and classification. *Linguisticae Investigationes* **30**(1), 3–26 (2007)
18. Nair, V., Hinton, G.E.: Rectified linear units improve restricted boltzmann machines. In: *Proceedings of the 27th international conference on machine learning (ICML-10)*. pp. 807–814 (2010)
19. Page, L., Brin, S., Motwani, R., Winograd, T.: The pagerank citation ranking: Bringing order to the web. *Tech. rep.*, Stanford InfoLab (1999)
20. Parker, R., Graff, D., Kong, J., Chen, K., Maeda, K.: *English gigaword*. Linguistic Data Consortium (2011)
21. Pershina, M., He, Y., Grishman, R.: Personalized page rank for named entity disambiguation. In: *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. pp. 238–243 (2015)
22. Ratnoff, L., Roth, D., Downey, D., Anderson, M.: Local and global algorithms for disambiguation to wikipedia. In: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*. pp. 1375–1384. Association for Computational Linguistics (2011)
23. Sil, A., Dinu, G., Kundu, G., Florian, R.: The ibm systems for entity discovery and linking at tac 2017. *Proc. TAC2017* (2017)
24. Tang, S., Ren, Y., Yang, X., Liu, D., Hu, G., Wu, F., Zhuang, Y.: The zhi-edl system for entity discovery and linking at tac kbp 2017. In: *TAC* (2017)
25. Watcharawittayakul, S., Xu, M., Jiang, H.: Dual fixed-size ordinal forgetting encoding (fofe) for competitive neural language models. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. pp. 4725–4730 (2018)
26. Xu, M., Jiang, H., Watcharawittayakul, S.: A local detection approach for named entity recognition and mention detection. In: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. vol. 1, pp. 1237–1247 (2017)

27. Xu, M., Nosirova, N., Jiang, K., Wei, F., Jiang, H.: Fofe-based deep neural networks for entity discovery and linking. In: TAC (2017)
28. Xu, M., Wei, F., Watcharawittayakul, S., Kang, Y., Jiang, H.: The yorknrm systems for trilingual edl tasks at tac kbp 2016. In: TAC (2016)
29. Yang, T., Du, D., Zhang, F.: The tai system for trilingual entity discovery and linking track in tac kbp 2017. In: TAC (2017)
30. Zesch, T., Müller, C., Gurevych, I.: Extracting lexical semantic knowledge from wikipedia and wiktionary. In: LREC. vol. 8, pp. 1646–1652 (2008)
31. Zhang, S., Jiang, H., Xu, M., Hou, J., Dai, L.: The fixed-size ordinally-forgetting encoding method for neural network language models. In: Proceedings of ACL (2015)