

A Convolution BiLSTM Neural Network Model for Chinese Event Extraction

No Author Given

No Institute Given

Abstract. Chinese event extraction is a challenging task in information extraction. Previous approaches highly depend on sophisticated feature engineering and complicated natural language processing (NLP) tools. In this paper, we first come up with the language specific issue in Chinese event extraction, and then propose a convolution bidirectional LSTM neural network that combines LSTM and CNN to capture both sentence-level and lexical information without any hand-craft features. Experiments on ACE 2005 dataset show that our approaches can achieve competitive performances in both trigger labeling and argument role labeling.

Keywords: event extraction, neural network, Chinese language processing

1 Introduction

Event extraction aims to extract events with specific types and their participants and attributes from unstructured data, which is an important and challenging task in information extraction. In this paper, we focus on the Chinese event extraction task proposed by the Automatic Content Extraction (ACE) program [7], which defines the following terminology for event extraction:

Trigger: the main word that most clearly expresses the occurrence of an event.

Argument: an entity, temporal expression or value that plays a certain role in the event.

There are two primary subtasks of an ACE event extraction system, namely *trigger labeling* and *argument labeling*. For example, consider the following Chinese sentence:

S1: Intel 在中国 成立了 研究中心。

Intel **founds** a research center in China.

In the trigger labeling task, “成立” (founds) should be labeled as the trigger of event type *Business*. Then in the argument labeling task, “Intel”, “中国” (China), and “研究中心” (research center) should be labeled as the roles of Agent, Place and Time respectively in this event.

Current state-of-the-art approaches [4, 14, 5] usually rely on a variety of elaborately features. In general, we can divide them into two categories: **lexical features** and **sentence-level features**. Take trigger labeling on the following sentences for example: the word “成立” (found) indicates a *Business* event in S1 but not in S2 or S3.

S2: 它成立于1994年，现在是一支深受欢迎的乐队。

It was **founded** in 1994, and now is a very popular band.

S3: 医院已成立救援中心。

The hospital has **founded** rescue centers.

Sentence-level features maintain important clues of the whole sentence. We can summarize S2 as “它是一支乐队” (it is a band) and encode it into a sentence-level feature, which indicates that the verb “成立” (founded) is not a trigger.

Lexical features contain semantic information of words and their surrounding context. In S3, given the next word “救援中心”(rescue centers) as a lexical feature, we can inference that “成立” (founded) is not a trigger of type *Business*.

Traditional approaches [2, 6, 13, 4] usually rely on a series of NLP tools to extract lexical features (e.g., part-of-speech tagging, named entity recognition) and sentence-level features (e.g., dependency parsing). Although they achieve high performance, they often suffer from hard feature engineering and error propagation from those external tools.

Recently, neural network models have been employed to produce competitive performance against traditional models for many NLP tasks. Chen et al. [5] propose a convolutional neural network to capture lexical-level clues, with a dynamic multi-pooling layer to capture sentence-level features, which yields state-of-art on English event extraction.

Inspired by the effectiveness of neural networks, in Section 2, we present a convolution bidirectional LSTM neural network that can learn both lexical and sentence-level features without any hand-engineered features in Chinese event extraction task. Specifically, we first use a bidirectional LSTM to encode the semantics of words in the whole sentence into **sentence-level features** without any parsing. Then, we can take advantage of a convolutional neural network to capture salient local **lexical features** for trigger disambiguation without any help from POS tags or NER.

We are also enlightened by the work of Chen and Ji [6], who are the first to report the language specific issue in Chinese trigger labeling. So we propose a character-based method committed to the problem in Section 2.3. Section 3 presents the model applied to argument labeling, and Section 4 discusses the experimental results. Section 5 concludes this paper.

2 Trigger Labeling

Trigger labeling, also called event detection, which aims to discover event triggers and assign them a predefined event type.

Unlike previous work [2, 6], which divide event detection into two subtasks: (1) trigger identification: to recognize the event trigger; (2) trigger classification: to assign an event type for an identified trigger. We jointly learn trigger identification and type classification by one network to reduce the error propagation problem of a pipeline model. Before we present our solution, we first come up with the language specific issue in Chinese trigger labeling.

2.1 Language Specific Issues

Unlike English, Chinese do not have delimiters between words. That makes word segmentation a fundamental step in Chinese event detection. However, we find that segmentation granularity does have an impact on the prediction. As showed in table 1, these triggers can not be recognized accurately if we simply predict whether a word is an event trigger or not. We summarize them as the following two types.

Table 1. Examples of inconsistency problem between words and triggers. Words are segmented by spaces.

#	Sentence	Triggers
S4	犯罪嫌疑人都落入法网。 The suspects were arrested .	落入法网 (arrested)
S5	警察击毙了一名歹徒。 Policies shoot and kill a criminal.	击 (shoot) 毙 (kill)
S6	这是一件预谋的凶杀案。 It is a premeditated murder case.	凶杀 (murder)

Cross-word triggers: While many events anchor on a single word, multiple words could reasonably be called a trigger. In S4, the *arrest_jail* event should be triggered by neither “落入” nor “法网”, but by “落入法网” (arrested).

Inside-word triggers: Almost all Chinese characters have their own meanings, and some of which can be triggers themselves. There may be more than one trigger in a word like “击毙” (shoot and kill). Continuous characters of a word can also form a trigger such as “凶杀” (murder) in “凶杀案” (murder case).

Table 2 summarizes the number of problematic triggers we found in ACE 2005 Chinese corpus using different Chinese word segmentation tools. Even the minimum inconsistency rate is as high as 14%.

Table 2. Number of triggers inconsistent with the words.

NLP Tools	Cross-word	Inside-word	Total
Stanford NLP ¹	487	166	653
Jieba ²	554	85	639
NLPIR ³	314	172	486

To address the language specific issues, we treat event detection as a sequence labeling task rather than classification. Sentences are tagged in the *BIO* scheme, where each token is labeled as *B-type* if it is the beginning of an event trigger with type *type*, or

¹ <http://nlp.stanford.edu/software/segmenter.shtml>

² <https://github.com/fxsjy/jieba>

³ <https://github.com/NLPIR-team/NLPIR>

I-type if it is inside a trigger, or *O* otherwise. Our first labelling model is a word-based BiLSTM model with a CNN layer as shown in Figure 1.

2.2 Word-based Method

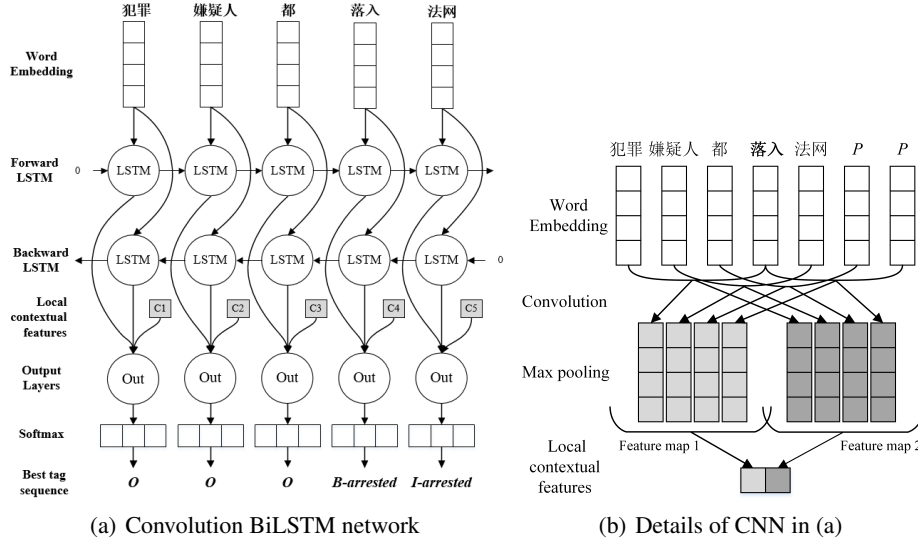


Fig. 1. The main architecture of our word-based model. The local contextual feature c_t (grey rectangle) in (a) for each word w_t is computed by the CNN as (b) illustrated. Our convolutional neural network learns an representation of local context information about the center word “落入”. Here the context size is 7 (3 words to the left and to the right of a center word), and we use a kernel of size 4 with two feature maps. The symbol P in sentence of (b) represents a padding word.

LSTM Network Recurrent neural networks (RNNs) maintain a memory based on history contextual information, which makes them a natural choice for processing sequential data. Unfortunately, it is difficult for standard RNNs to capture long range dependencies due to vanishing/exploding gradients [3]. Long Short-Term Memory [10] is explicitly designed to address the long-term dependency problem through purpose-built memory cells. They are composed of three multiplicative gates that control the proportion of information to forget and to store in the cell states.

For the event extraction task, if we access to both past and future contexts for a given time, we can take advantage of more sentence-level information and make better prediction. This can be achieved by bidirectional LSTM networks [8, 9]. Figure 1(a) shows the layers of a BiLSTM trigger identification model.

A forward LSTM network computes the hidden state \vec{h}_t of the past (left) context of the sentence at word w_t , while a backward LSTM network reads the same sentence in reverse and outputs \overleftarrow{h}_t given the future (right) context. In our implementation, we concatenate these two vectors to form the hidden state of a BiLSTM network, i.e. $h_t = [\vec{h}_t; \overleftarrow{h}_t]$.

Convolutional Neural Network Convolutional neural networks are originally applied to computer vision to capture local features [12]. CNN architectures have gradually shown effectiveness in various NLP tasks, and have been used for event extraction in previous studies [5]. We employ a convolutional neural network as illustrated in Figure 1(b) to extract local contextual information for each word in a sentence.

Given a sentence containing n words $\{w_1, w_2, \dots, w_n\}$, and the current center word w_t , a convolution operation involves a kernel, which is applied to words surrounding w_t in a window to generate the feature map. We can utilize multiple kernels with different widths to extract local features of various granularities. Then max pooling is performed over each map so that only the largest number of each feature map is recorded. One property of pooling is that it produces a fixed size output vector, which allows us to apply variable kernel sizes. And by performing the max operation, we are keeping the most salient information. Finally, we take the fixed length output vector c_{w_t} as a representation of local contextual information of center word w_t .

In our implementation, the sliding window size is 7 (3 words to the left and to the right of a center word), and we use several sizes of kernels to capture context information of various granularities.

The Output Layer We concatenate the hidden state h_t of BiLSTM with contextual feature c_{w_t} extracted by CNN at each time step t . Then $[h_t; c_{w_t}]$ is fed into a softmax layer to produce the log-probabilities of each label for w_t .

However, word-based method still can not solve the inconsistency problem caused by inside-word triggers. Like Chen and Ji [6], we construct a global errata table to record the most frequent triggers in the training set. During testing, if a word has an entry in the errata table, we replace its label with its corresponding trigger type directly.

2.3 Character-based Method

Despite of the effectiveness of the errata table, word-based method is not a flawless solution because it only recognizes triggers across words or frequent inside-word triggers appearing in training data.

Ideally, character-based method may solve both inconsistency problem. It uses the same tagging scheme as word-based method to label each character. As shown in Figure 2, the only difference between them is the input layers of their networks: character-based method uses character embedding while word-base method uses word embedding.

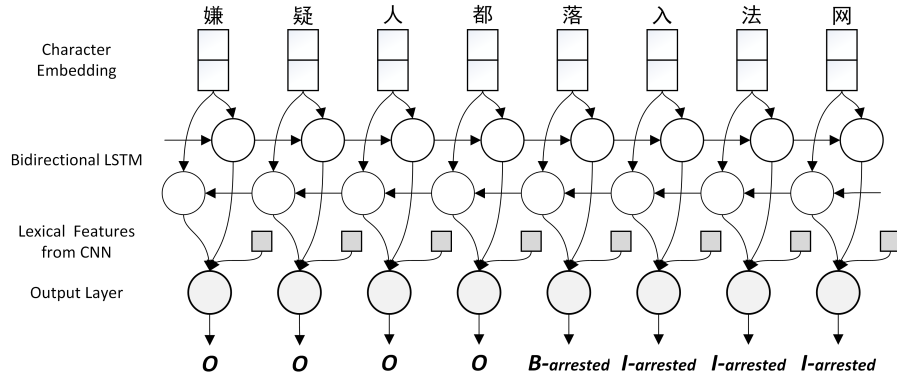


Fig. 2. Character-based Convolution BiLSTM network.

3 Argument Labeling

In the above section, we present our convolution BiLSTM model for trigger labeling. The idea of this neural network architecture is also suitable for argument labeling: we use a bidirectional LSTM to encode its sentence-level information, concatenated with a CNN-extracted local lexical feature, to predict whether an entity serves as an argument in a sentence. Next, we will introduce the main differences between the models used in trigger labeling and argument labeling.

3.1 Input Layer

As a pipeline system, besides word embeddings, we can use the information extracted from upstream trigger labeling task. Therefore, we propose four additional types of feature embeddings to form the input layer of BiLSTM and CNN.

- Trigger position feature: whether a word is in a trigger
- Trigger type feature: classified trigger type of a word, and *NONE* type for non-trigger words
- Entity position feature: whether a word is in an entity
- Entity type feature: entity type of a word, and *NONE* type for non-entity word. The ACE dataset provides ground truth of entity recognition, so we can generate entity features directly without external NLP tools.

We then transform these features into vectors by their lookup tables, and concatenating them with the origin word embeddings, as the final input layer of BiLSTM and CNN.

3.2 Output Layer

It is worth mentioning that argument labeling is no longer a sequence tagging task, but a classification task. ACE dataset provides ground truth of entity recognition, and it guarantees that arguments can only be labeled from those entities. As a result, we only need to predict the role of a tagged entity instead of every word in the whole sentence. For example, there are three triggers (bold words), and three entities (italic words) in S7, which together makes up nine pairs of trigger and argument candidate to be classified.

S7: 六起 **谋杀案** 发生在 *France*, 包括 *Bob* 的 **暗杀** 和 *Joe* 的 **杀害**。

Six **murders** occurred in *France*, including the **assassination** of *Bob* and the **kill**ing of *Joe*.

We modify the output layers of both CNN and BiLSTM network to adjust to the new task. For BiLSTM, we still want to make use of its ability to memory long sequences, so we regard the hidden state of the last word h_N as sentence-level information. And for CNN, we take all words of the whole sentence as the context, rather than a shallow window for each center word. Finally, we feed the concatenation of output vectors from two networks into a softmax classifier just like trigger labeling.

4 Experiments

4.1 Experimental Setup

We used the standard ACE 2005 corpus for our experiments, which contains 633 Chinese documents. Following the setup of Chen and Ji [6], we also randomly selected 509 documents for training and 64 documents as test set, and the reserved 60 documents for validation.

Evaluation Metric: Similar to previous work, we evaluated our models in terms of *precision* (P), *recall* (R), and *F-measure* (F) for each subtask. These performance metrics are computed following the standards of correctness for these subtasks:

- A trigger is correctly *identified* if its offsets exactly match a reference trigger;
- A trigger is correctly *classified* if its trigger type and offsets exactly match a reference trigger;
- An argument is correctly *identified* if its offset, related trigger type and trigger's offsets exactly match a reference argument;
- An argument is correctly *classified* if its offsets, role, related trigger type and trigger's offsets exactly match a reference argument.

4.2 Network Training

We implement the neural network using the Tensorflow library[1]. During training, we keep checking performance on the validation set and pick the highest F-score parameters for final evaluation.

Parameter initialization: Weight matrix parameters are randomly initialized with uniform samples from $[-0.01, 0.01]$. Bias vectors are initialized to zero.

Pre-trained embeddings: Word and character embeddings are pre-trained on over 261 thousand articles crawled from Chinese news website. All embeddings are fine-tuned during training.

Optimization algorithm: For all models presented, parameter optimization is performed using Adam[11] with gradient clipping[15]. We also apply the dropout method [16] on both the input and output vectors of all models to mitigate overfitting.

Hyper-parameters: In different stages of event extraction, we adopted different parameters. Table 3 summarizes the chosen hyper-parameters for all experiments.

Table 3. Hyper-parameters for all experiments.

Layer	Hyper-parameter	Trigger identification and classification	Argument identification and classification
Input	word embedding	100	100
	character embedding	50	100
	entity feature embedding	–	32
	trigger feature embedding	–	32
LSTM	state size	100	120
CNN	context size	7	sentence length
	kernel size	[3, 5, 7]	[2, 3, 4, 5]
	number of filters	[32, 32, 32]	[64, 64, 64, 64]
Dropout	dropout rate	0.5	0.5
	batch size	32	20
	gradient clipping	1.0	2.0

4.3 Trigger Labeling

Table 4 lists the results of previous work [6, 4] and our models. The performances of *Char-MEMM* and *Rich-L* are reported in their paper.

- *Char-MEMM* [6] is the first character-based method to handle the language specific issue, which trains a Maximum Entropy Markov Model to label each character with BIO tagging scheme.
- *Rich-L* [4] is a joint-learning, knowledge-rich approach that extends the union of the features employed by *Char-MEMM* and Li et al. [13] with six groups of linguistic features, including character-based features and discourse consistency features, which is the feature-based state-of-art system.

Compared with the feature-based approaches, all neural network based models outperform *Char-MEMM*, because they can capture semantic and syntactic information in the absence of feature engineering and avoid the errors propagated from other NLP tasks like NER and POS tagging.

Rich-L performs 10-fold cross-validation experiments, so that the results reported by them obtain more accurate estimation of system performance. Therefore, it is unfair

to directly compare our results with them. But we can still see the models evaluated in the bucket achieve competitive F-score on event identification without any human-designed features and discourse level knowledge.

Table 4. Comparison of different models on Chinese event detection (trigger identification and classification) (%).

Model	Trigger Identification			Trigger Classification		
	P	R	F	P	R	F
<i>Char-MEMM</i> [6]	82.4	50.6	62.7	78.8	48.3	59.9
<i>Rich-L</i> [4]	62.2	71.9	66.7	58.9	68.1	63.2
Word-based CNN	71.7	58.3	64.3	67.7	55.1	60.7
Word-based BiLSTM	68.4	61.2	64.6	63.9	57.4	60.5
Word-based C-BiLSTM	75.8	59.0	66.4	69.8	54.2	61.0
+ Errata table	76.0	63.8	69.3	69.8	59.9	64.5
Character-based C-BiLSTM	65.6	66.7	66.1	60.0	60.9	60.4
+ Errata table	68.1	69.2	68.7	61.6	64.7	63.1

Word-based Models vs. Character-based Models As we can summarize from table 4, when applying the same network architecture, word-based methods always have higher precisions while character-based methods always have higher recalls.

We then take a further step to see their impacts on different kinds of triggers. Table 5 shows that: (1) Word-based methods can not label inside-word triggers, while character-based methods can handle this issue nicely, which brings them higher overall recall; (2) Two methods achieve similar F-measure in regular trigger identification; (3) It is harder for character-based method to correctly identify cross-word triggers. As there are more cross-word triggers than inside-word triggers in dataset, the overall F-measure of word-based method is slightly higher.

Table 5. Results of different types of triggers with different models on trigger identification. Regular triggers mean triggers composed of exactly one word.

Model	Regular Triggers			Inside-word Triggers			Cross-word Triggers		
	P	R	F	P	R	F	P	R	F
Word-based C-BiLSTM	0.78	0.65	0.71	–	0	–	0.64	0.39	0.48
Character-based C-BiLSTM	0.72	0.70	0.71	0.30	0.69	0.41	0.57	0.22	0.32

There are several reasons causing the low precision of character-based method:

(1) Character-based method has to learn the extra word segmentation by themselves. 7.3% of triggers identified by it are partially mislabeled, like triggers in S8 and S9.

(2) Word embedding brings richer semantic information than character embedding. Take S10 as an example, characters “胡” and “同” do not have any meaning related to

the formed word “胡同” (the end of a road). But this word strongly suggests that “死” (dead) is not a trigger. Given the more accurate embedding of surrounding context, word-based networks can understand the meaning of the center word better and do better disambiguation.

(3) RNN in character-based method needs to maintain information for longer sequence, as 1.7 times longer than the average length of word sequences. Evaluating on sentences containing more than 150 characters, F-measure of character-based method is 70%, while word-based method can achieve 72.8%.

Table 6. Error analysis: examples of triggers mislabeled by character-based C-BiLSTM, but can be identified correctly by word-based C-BiLSTM.

#	Sentence	Correct Labels	Wrong Labels
S8	走访 相关 人员 以后, ... After visiting to relevant staff, ...	走 (go) 访 (visit) BI	BO
S9	贺电 全文 如下: There is the full congratulatory message	贺 (congratulatory) 电 (message) BI	OB
S10	小偷 被 逼 进 死 胡同。 The thief was chased into a dead end .	死 (dead) 胡同 (end) O OO	B OO

Neural Network Architectures Our convolution BiLSTM model has two main components that we could take them apart to understand their impacts on the overall performance. As table 4 shows, BiLSTM is slightly more efficient than CNN, and the convolution BiLSTM model outperforms other models. Constructing an errata table is an effective method that increases both precisions and recalls.

We also evaluate the capacity of each network on trigger disambiguation. Table 7 provides suggestive evidence that CNN-extracted local features, together with LSTM-extracted sentence-level information can help reduce some errors caused by ambiguous triggers as we expected.

Table 7. Percentages of ambiguous words whose all occurrences in the test set are classified correctly. Ambiguous words can have different labels according to their meaning and context, like the word “成立” (found) in S1 ~ S3. All networks listed in this table are word-based.

	BiLSTM	CNN	C-BiLSTM
ambiguous word classification (%)	58.4	60.0	62.5

4.4 Argument Labeling

Table 8 shows results for argument labeling after trigger labeling. As we can observe from our evaluation standards that once a trigger has not been labeled correctly, neither

of its arguments will be labeled correctly. In the stage of trigger labeling, the character-based methods have higher recalls and can extract more golden triggers. As a result, character-based methods perform much better than word-based methods in argument labeling. And we can draw the same conclusion that word-based methods have higher precision while character-based methods achieve higher recalls, as trigger labeling.

Table 8. Comparison of different models on Chinese argument labeling (%). +e means that the input (result of trigger labeling) has been modified by an errata table.

Model	Argument Identification			Argument Classification		
	P	R	F	P	R	F
<i>Char-MEMM</i>	64.4	36.4	46.5	60.6	34.3	43.8
<i>Rich-L</i>	43.6	57.3	49.5	39.2	51.6	44.6
word-based C-BiLSTM	56.6	43.6	49.3	49.7	38.3	43.2
word-based C-BiLSTM +(e)	56.5	47.0	51.3	49.6	41.3	45.0
character-based C-BiLSTM	53.2	51.6	52.4	47.0	45.6	46.3
character-based C-BiLSTM +(e)	53.0	52.2	52.6	47.3	46.6	46.9

Errata table is not such effective as in trigger labeling, especially for character-based C-BiLSTM. Adding an errata table even drops a little precision in argument identification.

Char-MEMM concludes that *neighbor word* features are fairly effective. They use the left word and right word of an entity to reduce spurious argument, which is a similar objective with our CNN-extracted lexical features. But we can achieve much better results in argument identification and classification.

It is worth noting that some of the arguments are not in the same sentence with their triggers. It is a bottleneck of our C-BiLSTM model, while *Rich-L* uses discourse-level features to handle this problem. Under this unfavorable circumstance, our C-BiLSTM can still achieve a comparable result against sophisticated human designed features.

5 Conclusion

In this paper, we propose a novel convolution bidirectional LSTM model on Chinese event extraction task. Our model departs from the inherent characteristic of Chinese, formulates the event detection task as a sequence labeling fashion, and features both bidirectional LSTM and CNN to capture both sentence-level and lexical features from raw text. Experimental results show that without human-designed features and external resources, our neural network method can achieve comparable performances on ACE 2005 datasets with traditional feature based methods.

References

1. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., et al.: Tensorflow: Large-scale machine learning on heterogeneous distributed systems. arXiv preprint arXiv:1603.04467 (2016)
2. Ahn, D.: The stages of event extraction. In: Proceedings of the Workshop on Annotating and Reasoning about Time and Events. pp. 1–8. Association for Computational Linguistics (2006)
3. Bengio, Y., Simard, P., Frasconi, P.: Learning long-term dependencies with gradient descent is difficult. *Neural Networks, IEEE Transactions on* 5(2), 157–166 (1994)
4. Chen, C., Ng, V.: Joint modeling for chinese event extraction with rich linguistic features. In: COLING. pp. 529–544. Citeseer (2012)
5. Chen, Y., Xu, L., Liu, K., Zeng, D., Zhao, J.: Event extraction via dynamic multi-pooling convolutional neural networks. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing. vol. 1, pp. 167–176 (2015)
6. Chen, Z., Ji, H.: Language specific issue and feature exploration in chinese event extraction. In: Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Short Papers. pp. 209–212. Association for Computational Linguistics (2009)
7. Doddington, G.R., Mitchell, A., Przybocki, M.A., Ramshaw, L.A., Strassel, S., Weischedel, R.M.: The automatic content extraction (ace) program-tasks, data, and evaluation. In: LREC. vol. 2, p. 1 (2004)
8. Graves, A., Schmidhuber, J.: Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks* 18(5), 602–610 (2005)
9. Graves, A., Mohamed, A.r., Hinton, G.: Speech recognition with deep recurrent neural networks. In: Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on. pp. 6645–6649. IEEE (2013)
10. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Computation* 9(8), 1735–1780 (1997)
11. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
12. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86(11), 2278–2324 (1998)
13. Li, P., Zhou, G., Zhu, Q., Hou, L.: Employing compositional semantics and discourse consistency in chinese event extraction. In: Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning. pp. 1006–1016. Association for Computational Linguistics (2012)
14. Li, Q., Ji, H., Huang, L.: Joint event extraction via structured prediction with global features. In: ACL (1). pp. 73–82 (2013)
15. Pascanu, R., Mikolov, T., Bengio, Y.: On the difficulty of training recurrent neural networks. *ICML (3)* 28, 1310–1318 (2013)
16. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* 15(1), 1929–1958 (2014)