

# HMEAE: Hierarchical Modular Event Argument Extraction

Xiaozhi Wang<sup>1\*</sup>, Ziqi Wang<sup>1\*</sup>, Xu Han<sup>1</sup>, Zhiyuan Liu<sup>1†</sup>,  
Juanzi Li<sup>1</sup>, Peng Li<sup>2</sup>, Maosong Sun<sup>1</sup>, Jie Zhou<sup>2</sup>, Xiang Ren<sup>3</sup>

<sup>1</sup>Department of Computer Science and Technology, Tsinghua University, Beijing, China

<sup>2</sup>Pattern Recognition Center, WeChat AI, Tencent Inc, China

<sup>3</sup>Department of Computer Science, University of Southern California, CA, USA

{xz-wang16, ziqi-wan16, hanxu17}@mails.tsinghua.edu.cn

## Abstract

Existing event extraction methods classify each argument role independently, ignoring the conceptual correlations between different argument roles. In this paper, we propose a Hierarchical Modular Event Argument Extraction (HMEAE) model, to provide effective inductive bias from the concept hierarchy of event argument roles. Specifically, we design a neural module network for each basic unit of the concept hierarchy, and then hierarchically compose relevant unit modules with logical operations into a role-oriented modular network to classify a specific argument role. As many argument roles share the same high-level unit module, their correlation can be utilized to extract specific event arguments better. Experiments on real-world datasets show that HMEAE can effectively leverage useful knowledge from the concept hierarchy and significantly outperform the state-of-the-art baselines. The source code can be obtained from <https://github.com/thunlp/HMEAE>.

## 1 Introduction

Event argument extraction (EAE) aims to identify the entities serving as event arguments and classify the roles they play in an event. For instance, given that the word “sold” triggers a Transfer-Ownership event in the sentence “*Steve Jobs sold Pixar to Disney*”, EAE aims to identify that “*Steve Jobs*” is an event argument and its argument role is “Seller”. Most event extraction (EE) methods treat EE as a two-stage problem, including event detection (ED, to identify the trigger word and determine the event type) and EAE. As ED is well-studied (Nguyen and Grishman, 2018; Zhao et al., 2018) in recent years, EAE becomes the bottleneck of EE.

\* indicates equal contribution

† Corresponding author: Z.Liu(liuzy@tsinghua.edu.cn)

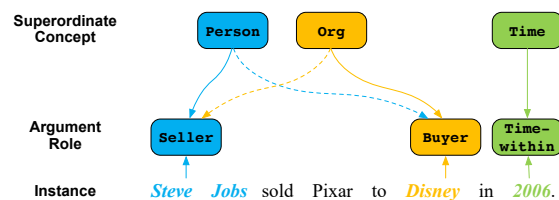


Figure 1: An example of the concept hierarchy.

Since EE benefits many NLP applications (Yang et al., 2003; Basile et al., 2014; Cheng and Erk, 2018), intensive efforts have been devoted to detecting events and extracting their event arguments. Traditional feature-based methods (Patwardhan and Riloff, 2009; Liao and Grishman, 2010b,a; Huang and Riloff, 2012; Li et al., 2013) rely on hand-crafted features and patterns. With the ongoing development of neural networks, various neural networks have been used to automatically represent textual semantics with low-dimensional vectors, and further extract event arguments based on those semantic vectors, including convolutional neural networks (Chen et al., 2015) and recurrent neural networks (Nguyen et al., 2016; Sha et al., 2018). Advanced techniques also have been adopted to further improve EE, such as zero-shot learning (Huang et al., 2018), multi-modal integration (Zhang et al., 2017), and weakly supervised methods (Chen et al., 2017; Wang et al., 2019).

However, the existing methods all treat argument roles as independent of each other, regardless of the fact that some argument roles are conceptually closer than others. Taking Figure 1 as an example, “Seller” is conceptually closer to “Buyer” than “Time-within”, because they share the same superordinate concepts “Person” and “Org” in the concept hierarchy. Intuitively, the concept hierarchy will provide extra informa-

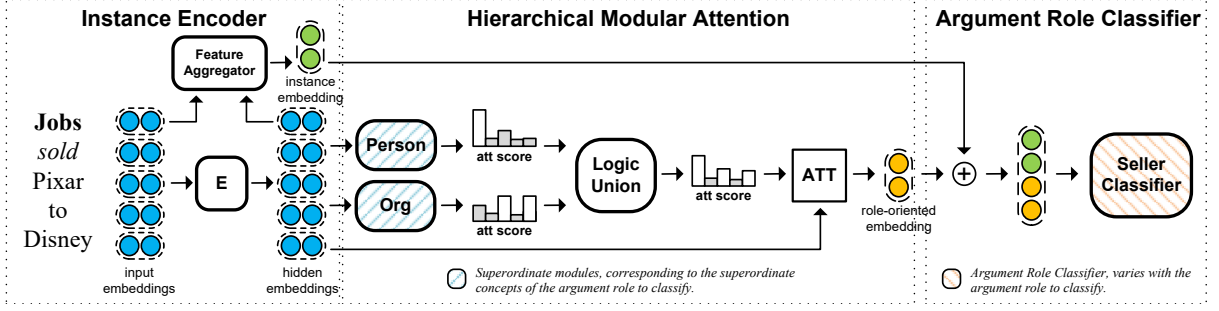


Figure 2: The overall architecture of HMEAE. Take the argument role “Seller” as an example.

tion about the correlation between argument roles and help the argument role classification.

To leverage the concept hierarchy information to improve EAE, we propose the Hierarchical Modular Event Argument Extraction (HMEAE) model. Inspired by the previous hierarchical classification works (Qiu et al., 2011; Shimura et al., 2018; Han et al., 2018) and the neural module networks (NMNs) (Andreas et al., 2016), HMEAE adopts the NMNs to enable a flexible network architecture imitating the concept hierarchical structure, which can provide effective inductive bias for better classification performance.

As Figure 1 shows, we divide the concepts into two types: the superordinate concepts representing more abstractive concepts, and the fine-grained argument roles. An argument role can belong to more than one superordinate concept, e.g., “Seller” belongs to both “Person” and “Org”. As shown in Figure 2, we set a neural module network for each concept, and hierarchically compose them as the structure of the concept hierarchy into a role-oriented modular network to predict the argument role for each entity: (1) First, for each superordinate concept, a superordinate concept module (SCM) is instantiated to highlight the textual information related to the concept; (2) Then for each argument role, the SCMs corresponding to its superordinate concepts are composed by role-specific logic union modules to obtain a unified high-level module; (3) Finally, an argument role classifier is set to predict whether the entity is of the given argument role relying on the output of high-level modules.

Intuitively, considering the concept hierarchy in our model brings the following benefits: (1) The high-level modules can significantly enhance the classifiers, e.g., it is easier to determine whether an entity is “Time-within” when paying more attention to the words about “Time”. (2) A super-

ordinate concept module is shared among different argument roles in the concept hierarchy. Hence, it can capture the concept features from the shared information in the data of its subordinate argument roles, and further provide the correlation information as an effective inductive bias to the argument role classifiers. We conduct experiments on two real-world datasets. The experimental results show that our methods can achieve state-of-the-art results.

## 2 Methodology

In this section, we will introduce the overall framework of HMEAE. As shown in Figure 2, HMEAE consists of three components: (1) The instance encoder represents a sentence into hidden embeddings and utilizes a feature aggregator to aggregate sentence information into a unified instance embedding. (2) The hierarchical modular attention component builds a role-oriented embedding to highlight the information about the superordinate concepts of the argument role in the predefined concept hierarchy. (3) The argument role classifier relies on the instance embedding and the role-oriented embedding to estimate the probability of a certain argument role for the instance.

### 2.1 Instance Encoder

We denote an instance as an  $n$ -word sequence  $x = \{w_1, \dots, t, \dots, a, \dots, w_n\}$ , where  $t, a$  denote the trigger word and the candidate argument respectively. The trigger word is detected by the previous event detection models (independent of our work) and each named entity in the sentence is a candidate argument.

**Sentence Encoder** is adopted to encode the word sequence into hidden embeddings,

$$\{h_1, h_2, \dots, h_n\} = E(w_1, \dots, t, \dots, a, \dots, w_n), \quad (1)$$

where  $E(\cdot)$  is the neural network to encode the sentence. In this paper, we select CNN (Chen et al., 2015) and BERT (Devlin et al., 2019) as encoders.

**Feature Aggregator** aggregates the hidden embeddings into an instance embedding. Our method is independent of the feature aggregator mechanism. Here, we follow Chen et al. (2015) and use dynamic multi-pooling as the feature aggregator:

$$\begin{aligned} [x_{1,p_t}]_i &= \max\{[h_1]_i, \dots, [h_{p_t}]_i\}, \\ [x_{p_t+1,p_a}]_i &= \max\{[h_{p_t+1}]_i, \dots, [h_{p_a}]_i\}, \\ [x_{p_a+1,n}]_i &= \max\{[h_{p_a+1}]_i, \dots, [h_n]_i\}, \\ \mathbf{x} &= [x_{1,p_t}; x_{p_t+1,p_a}; x_{p_a+1,n}] \end{aligned} \quad (2)$$

where  $[\cdot]_i$  is the  $i$ -th value of a vector,  $p_t, p_a$  are the positions of the trigger  $t$  and the candidate argument  $a$  respectively. We concatenate the piecewise max-pooling results as the instance embedding  $\mathbf{x}$ .

## 2.2 Hierarchical Modular Attention

As shown in Figure 2, given the hidden embeddings  $\{h_1, h_2, \dots, h_n\}$ , a superordinate concept module gives an attention score for each hidden embedding to model its correlation with the specific superordinate concept. As an argument role can belong to more than one superordinate concept, we set a logic union module to combine the scores from different superordinate modules together. For each argument role, we hierarchically compose its superordinate concept modules into the integrated hierarchical modular attention component to build its role-oriented embedding.

**Superordinate Concept Module** For a specific superordinate concept  $c$ , we represent its semantic features with a trainable vector  $\mathbf{u}_c$ . Following Luong et al. (2015), we adopt a multi-layer perceptron to calculate the attention scores. We first calculate the hidden state,

$$h_i^c = \tanh(\mathbf{W}_a[h_i; \mathbf{u}_c]). \quad (3)$$

Then, we apply a softmax operation to get the attention score for the hidden embedding  $h_i$ ,

$$s_i^c = \frac{\exp(\mathbf{W}_b h_i^c)}{\sum_{j=1}^n \exp(\mathbf{W}_b h_j^c)}, \quad (4)$$

where  $\mathbf{W}_a$  and  $\mathbf{W}_b$  are trainable matrices shared among different superordinate concept modules.

**Logic Union Module** Given an argument role  $r \in \mathcal{R}$ , we denote its  $k$  superordinate concepts as  $c_1, c_2, \dots, c_k$ , and the corresponding attention

scores for  $h_i$  are  $s_i^{c_1}, s_i^{c_2}, \dots, s_i^{c_k}$  computed by Eq. (4). As information about all the superordinate concepts should be retained in the role-oriented embedding, we calculate the mean of the attention scores as the role-oriented attention score,

$$s_i^r = \frac{1}{k} \sum_{j=1}^k s_i^{c_j}, \quad (5)$$

and then calculate the weighted sum of hidden embeddings as the role-oriented embedding,

$$\mathbf{e}^r = \sum_{i=1}^n s_i^r h_i. \quad (6)$$

## 2.3 Argument Role Classifier

We concatenate the instance embedding  $\mathbf{x}$  and the role-oriented embedding  $\mathbf{e}^r$  as the input feature for the argument role classifier, and estimate the probability of  $r \in \mathcal{R}$  of instance  $x$  as follows:

$$p(r|x) = \frac{\exp(\mathbf{r}^\top [\mathbf{x}; \mathbf{e}^r])}{\sum_{\tilde{r} \in \mathcal{R}} \exp(\tilde{\mathbf{r}}^\top [\mathbf{x}; \mathbf{e}^{\tilde{r}}])}, \quad (7)$$

where  $\mathbf{r}$  is the embedding of the argument role  $r$ .

The objective function is defined as follows:

$$\mathcal{L}(\theta) = - \sum_l \log p(r_l | x_l), \quad (8)$$

where  $\theta$  is all parameters of our model. We adopt Adam (Kingma and Ba, 2015) to minimize  $\mathcal{L}(\theta)$ .

## 3 Experiments

### 3.1 Experimental Settings

In the experiments, our model with CNN as the encoder is named **HMEAE (CNN)**, whose most hyperparameters are the same as Chen et al. (2015) for a fair comparison. Our model with BERT as the encoder is named **HMEAE (BERT)**. For a fair comparison, we set a vanilla BERT baseline named **DMBERT**, which is without the hierarchical modular attention module but with a dynamic multi-pooling layer (Chen et al., 2015) as the feature aggregator. As our work does not involve the event detection stage, we conduct the argument role classification based on the event detection models in Chen et al. (2015) and Wang et al. (2019) for **HMEAE (CNN)** and **HMEAE (BERT)** respectively.

## Datasets and Evaluation

We evaluate our models on two real-world datasets: the widely-used ACE 2005 (Walker et al., 2006) and the newly-developed TAC KBP 2016 (Ellis et al., 2015).

**ACE 2005** (LDC2006T06) is the most widely-used dataset in event extraction. It contains 599 documents, which are annotated with 8 event types, 33 event subtypes, and 35 argument roles. We evaluate our models by the performance of argument classification. An argument is correctly classified if its event subtype, offsets and argument role match the annotation results. Following the previous works (Liao and Grishman, 2010b; Chen et al., 2015), we use the same test set containing 40 newswire documents, a development set with 30 randomly selected documents and training set with the remaining 529 documents.

**TAC KBP 2016** The Text Analysis Conference (TAC) is a series of evaluation workshops organized to encourage research in NLP and related applications. In this paper, we use the data of the TAC KBP 2016 Event Argument Extraction track (LDC2017E05). This competition annotates difficult test data but no training data. They encourage participants to use training data from any other sources. Hence we use the ACE 2005 dataset as our training data, which is less than the data used by the baselines on TAC KBP 2016.

## Concept Hierarchy Design

Considering there is not an existing ontology in the datasets, we manually design a concept hierarchy with 8 different superordinate concepts for our models. The principle of designing the concept hierarchy is to induce superordinate concepts from the specific labels using human experience. For instance, people can easily summarize “Origin” and “Destination” into “Place”, which is the desired superordinate concept. Although the hierarchy used in this paper may not be generalized to other datasets with different label definitions, it is tractable to design an appropriate concept hierarchy with minimal human efforts and provide effective inductive bias via our method. For the details about the specific concept hierarchy used in this paper, please refer to the appendix.

## Hyperparameter Settings

**CNN** slides a convolution kernel over the input embedding sequence to get hidden embeddings.

Following previous work, the input embedding of each word consists of its word embedding, position embedding, and event type embedding. The hyperparameter settings for HMEAE (CNN) are shown in Table 1.

Learning Rate	1e-03
Batch Size	20
Word Embedding Dimension	100
Dropout Probability	0.5
Hidden Layer Dimension	300
Kernel Size	3
Position Embedding Dimension	5
Event Type Embedding Dimension	5
$u_c$ dimension	900
$W_b$ dimension	900

Table 1: Hyperparameter settings for CNN models.

**BERT** adopts multi-layer bidirectional transformers to encode the input embedding sequence into hidden embeddings. The hyperparameters of DMBERT and HMEAE (BERT) are the same as the BERT<sub>BASE</sub> model. To utilize the event type information in our model, we append a special token into each input sequence for BERT to indicate the event type. Additional hyperparameters used in our experiments are shown in Table 2.

Learning Rate	6e-05
Batch Size	50
Kernel Size	3
Warmup Rate	0.1
$u_c$ dimension	900
$W_b$ dimension	900

Table 2: Hyperparameter settings for BERT models.

## 3.2 Overall Evaluation Results

We compare our models with various state-of-the-art baselines on ACE 2005: (1) Feature-based methods, including **Li’s joint** (Li et al., 2013) and **RBPB** (Sha et al., 2016). (2) Vanilla neural network methods, including **DMCNN** (Chen et al., 2015) and **JRNN** (Nguyen et al., 2016). (3) Neural network with syntax information, like **dbRNN** (Sha et al., 2018) enhancing the recurrent neural network with dependency bridges to consider syntactically related information.

On TAC KBP 2016, we compare our models with the top systems (Dubbin et al., 2016; Hsi et al., 2016; Ferguson et al., 2016) of the competition as well as DMCNN and DMBERT.



Method	Argument Role Classification		
	P	R	F1
Li’s Joint (Li et al., 2013)	64.7	44.4	52.7
DMCNN (Chen et al., 2015)	62.2	46.9	53.5
RBPB (Sha et al., 2016)	54.1	53.5	53.8
JRNN (Nguyen et al., 2016)	54.2	<b>56.7</b>	55.4
dbRNN (Sha et al., 2018)	<b>66.2</b>	52.8	58.7
HMEAE (CNN)	57.3	54.2	55.7
DMBERT	58.8	55.8	57.2
HMEAE (BERT)	62.2	56.6	<b>59.3</b>

Table 3: The overall results (%) on ACE 2005.

Method	Argument Role Classification		
	P	R	F1
DISCERN-R (Dubbin et al., 2016)	7.9	7.4	7.7
Washington4 (Ferguson et al., 2016)	32.1	5.0	8.7
CMU CS Event1 (Hsi et al., 2016)	<b>31.2</b>	4.9	8.4
Washington1 (Ferguson et al., 2016)	26.5	6.8	10.8
DMCNN (Chen et al., 2015)	17.9	16.0	16.9
HMEAE (CNN)	15.3	22.5	18.2
DMBERT	22.6	24.7	23.6
HMEAE (BERT)	24.8	<b>25.4</b>	<b>25.1</b>

Table 4: The overall results (%) on TAC KBP 2016.

The results are shown in Table 3 and Table 4. We have the following observations from the results: (1) HMEAE (CNN) and HMEAE (BERT) achieve improvements (about 2% in F1) as compared with DMCNN and DMBERT respectively, which have almost the same network framework with our models except the hierarchical modular attention. It indicates that our hierarchical modular method works well to enhance the EAE models with the inductive bias from the concept hierarchy. (2) HMEAE (BERT) is comparable to dbRNN and achieves the state-of-the-art performance among existing methods. Our methods adopt modular networks to consider the hierarchical concept knowledge in EAE, which is compatible with the sophisticated linguistic knowledge adopted in dbRNN. We will explore to integrate this two kinds of external knowledge together to further improve EAE in the future.

### 3.3 Case Study

To verify whether the superordinate concept modules work as we designed, we conduct a case study. We visualize the attention score  $s_i^c$  of HMEAE (BERT) on a sentence randomly sampled from the ACE 2005 dataset in Figure 3. We

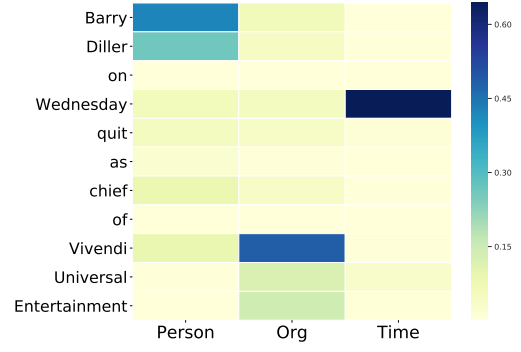


Figure 3: Heatmap for attention scores of three superordinate concept modules of the left sentence.

observe that the attention scores for the hidden embeddings of words related to the superordinate concept are much higher than the others. It indicates that as the superordinate concept modules are shared among their subordinate argument roles, the superordinate concept modules can capture the concept features well without being specially trained with exclusive data.

## 4 Conclusion and Future work

In this paper, we propose a hierarchical modular event argument extraction model (HMEAE), which adopts flexible modular networks to utilize the hierarchical concept correlation among argument roles as effective inductive bias. Experimental results show that HMEAE achieves the state-of-the-art performance. In the future, we will further explore to leverage other kinds of inductive bias from human experience to improve extensive tasks with our modular networks.

## Acknowledgement

This work is supported by the National Natural Science Foundation of China (NSFC No. 61572273, 61661146007), the NSFC key project (U1736204) and Tsinghua University Initiative Scientific Research Program (20151080406). This work is also supported by the Pattern Recognition Center, WeChat AI, Tencent Inc. Xu Han and Xiaozhi Wang are also supported by 2018 and 2019 Tencent Rhino-Bird Elite Training Program respectively. Xiaozhi Wang is also supported by Tsinghua University Initiative Scientific Research Program. Xiang Ren is supported in part by National Science Foundation SMA 18-29268, DARPA MCS, GAILA, IARPA BETTER, and Schmidt Family Foundation.

## References

- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. [Neural module networks](#). In *Proceedings of CVPR*, pages 39–48.
- P Basile, A Caputo, G Semeraro, and L Siciliani. 2014. [Extending an information retrieval system through time event extraction](#). In *Proceedings of DART*, pages 36–47.
- Yubo Chen, Shulin Liu, Xiang Zhang, Kang Liu, and Jun Zhao. 2017. [Automatically labeled data generation for large scale event extraction](#). In *Proceedings of ACL*, pages 409–419.
- Yubo Chen, Liheng Xu, Kang Liu, Daojian Zeng, and Jun Zhao. 2015. [Event extraction via dynamic multi-pooling convolutional neural networks](#). In *Proceedings of ACL-IJCNLP*, pages 167–176.
- Pengxiang Cheng and Katrin Erk. 2018. [Implicit argument prediction with event knowledge](#). In *Proceedings of ACL*, pages 831–840.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of NAACL-HLT*, pages 4171–4186.
- Greg Dubbin, Archana Bhatia, Bonnie Dorr, Adam Dalton, Kristy Hollingshead, Suriya Kandaswamy, Ian Perera, and Jena D Hwang. 2016. [Improving discern with deep learning](#). In *TAC*.
- Joe Ellis, Jeremy Getman, Dana Fore, Neil Kuster, Zhiyi Song, Ann Bies, and Stephanie M Strassel. 2015. [Overview of linguistic resources for the TAC KBP 2016 evaluations: Methodologies and results](#). In *TAC*.
- James Ferguson, Colin Lockard, Natalie Hawkins, Stephen Soderland, Hannaneh Hajishirzi, and Daniel S Weld. 2016. [University of Washington TAC-KBP 2016 system description](#). In *TAC*.
- Xu Han, Pengfei Yu, Zhiyuan Liu, Maosong Sun, and Peng Li. 2018. [Hierarchical relation extraction with coarse-to-fine grained attention](#). In *Proceedings of EMNLP*, pages 2236–2245.
- Andrew Hsi, Jaime G Carbonell, and Yiming Yang. 2016. [CMU CS event TAC-KBP2016 event argument extraction system](#). In *TAC*.
- Lifu Huang, Heng Ji, Kyunghyun Cho, Ido Dagan, Sebastian Riedel, and Clare Voss. 2018. [Zero-shot transfer learning for event extraction](#). In *Proceedings of ACL*, pages 2160–2170.
- Ruihong Huang and Ellen Riloff. 2012. [Modeling textual cohesion for event extraction](#). In *Proceedings of AAAI*, pages 1664–1670.
- Diederik P Kingma and Jimmy Ba. 2015. [Adam: a method for stochastic optimization](#). In *Proceedings of ICLR*.
- Qi Li, Heng Ji, and Liang Huang. 2013. [Joint event extraction via structured prediction with global features](#). In *Proceedings of ACL*, pages 73–82.
- Shasha Liao and Ralph Grishman. 2010a. [Filtered ranking for bootstrapping in event extraction](#). In *Proceedings of COLING*, pages 680–688.
- Shasha Liao and Ralph Grishman. 2010b. [Using document level cross-event inference to improve event extraction](#). In *Proceedings of ACL*, pages 789–797.
- Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. [Effective approaches to attention-based neural machine translation](#). In *Proceedings of EMNLP*, pages 1412–1421.
- Thien Nguyen and Ralph Grishman. 2018. [Graph convolutional networks with argument-aware pooling for event detection](#). In *Proceedings of AAAI*, pages 5900–5907.
- Thien Huu Nguyen, Kyunghyun Cho, and Ralph Grishman. 2016. [Joint event extraction via recurrent neural networks](#). In *Proceedings of NAACL-HLT*, pages 300–309.
- Siddharth Patwardhan and Ellen Riloff. 2009. [A unified model of phrasal and sentential evidence for information extraction](#). In *Proceedings of EMNLP*, pages 151–160.
- Xipeng Qiu, Xuanjing Huang, Zhao Liu, and Jinlong Zhou. 2011. [Hierarchical text classification with latent concepts](#). In *Proceedings of ACL-HLT*, pages 598–602.
- Lei Sha, Jing Liu, Chin-Yew Lin, Sujian Li, Baobao Chang, and Zhifang Sui. 2016. [RBPB: Regularization-based pattern balancing method for event extraction](#). In *Proceedings of ACL*, pages 1224–1234.
- Lei Sha, Feng Qian, Baobao Chang, and Zhifang Sui. 2018. [Jointly extracting event triggers and arguments by dependency-bridge RNN and tensor-based argument interaction](#). In *Proceedings of AAAI*, pages 5916–5923.
- Kazuya Shimura, Jiyi Li, and Fumiyo Fukumoto. 2018. [HFT-CNN: Learning hierarchical category structure for multi-label short text categorization](#). In *Proceedings of EMNLP*, pages 811–816.
- Christopher Walker, Stephanie Strassel, Julie Medero, and Kazuaki Maeda. 2006. [ACE 2005 multilingual training corpus](#). *LDC, Philadelphia*.
- Xiaozhi Wang, Xu Han, Zhiyuan Liu, Maosong Sun, and Peng Li. 2019. [Adversarial training for weakly supervised event detection](#). In *Proceedings of NAACL-HLT*, pages 998–1008.

- Hui Yang, Tat-Seng Chua, Shuguang Wang, and Chun-Keat Koh. 2003. [Structured use of external knowledge for event-based open domain question answering](#). In *Proceedings of SIGIR*, pages 33–40.
- Tongtao Zhang, Spencer Whitehead, Hanwang Zhang, Hongzhi Li, Joseph Ellis, Lifu Huang, Wei Liu, Heng Ji, and Shih-Fu Chang. 2017. [Improving event extraction via multimodal integration](#). In *Proceedings of MM*, pages 270–278.
- Yue Zhao, Xiaolong Jin, Yuanzhuo Wang, and Xueqi Cheng. 2018. [Document embedding enhanced event detection with hierarchical and supervised attention](#). In *Proceedings of ACL*, pages 414–419.