

# Scale Up Event Extraction Training with Little Human Involvement

## Abstract

Supervised learning is an effective method for building event extraction systems. The effectiveness of the learning heavily depends on the quality of the training datasets. Existing event extractor systems are typically built upon expert-annotated datasets. However, due to drastic efforts involved in annotating text, these datasets are usually small and cover only a limited variety of event types. This limits the quality of learned event extractor, making it hard to generalize. In this paper, we investigate ways to generate training datasets that require little expert involvement but can cover a richer set of event types. We achieve this by employing distant supervision to automatically create event annotations from unlabelled text using existing structured knowledge bases or tables. We then develop a novel neural network model with post inference, to detect multi-typed event mentions **FIX:ZW: reviewers may not know what is a event mention** with corresponding arguments. Experiments on the datasets collected through Freebase and Wikipedia tables show that it is possible to learn to extract events of rich types without human-annotated training data. **FIX:ZW: Will come back to the experimental results later.**

## Introduction

Event extraction, as represented by the Automatic Content Extraction (ACE) task, is a key enabling technique for many natural language processing (NLP) applications. Current event extractors are typically built through applying supervised learning to learn over labelled datasets. This means that the performance of event extraction systems highly dependent on the quality and coverage of the training datasets.

Constructing high-quality training data for event extraction is, however, an expensive and error-prone process (Aguilar et al. 2014; Song et al. 2015). This requires the involvement of linguists to design annotation templates and rules for a set of predefined event types, and the employment of annotators to manually label large datasets. Because of the drastic efforts required for manual labelling, existing training datasets are usually small, covering only a small set of event types. As a result, current event extraction systems built upon of human-annotated datasets have two major drawbacks. They can only handle (1) single-token trig-

gers<sup>1</sup> and (2) scenarios where each event is merely associated with a single type. In our real-world, however, there are countless examples where trigger annotations are unavailable (**FIX:ZW: This has nothing to do with single-token triggers??**) or one event is associated with multiple types. To make event extractors more practical, we need to find ways to support these scenarios.

This paper presents a novel approach for automatic training data generation, specifically targeting event extraction. Our approach advances prior work in two aspects: (1) it does not rely on expert-annotated texts (hence it reduces expert involvement) and (2) supports multi-typed event extraction. The former is achieved by automatically collecting training data through indirect supervision from structured knowledge bases or knowledge tables to automatically label event structures from plain text. The later is achieved by proposing a novel neural network model with post inference to detect multi-typed events without relying on explicit trigger annotations.

Our first insight is that structured knowledge bases (KB) typically organize complex structured information in tables; and such knowledge tables often share similar structures with ACE event definitions, i.e., a particular entry of such tables usually implies the occurrence of certain events. Recent studies (Mintz et al. 2009; Zeng et al. 2015) have demonstrated the effectiveness of KB as distant supervision (DS) for binary relation extraction. In this paper, we aim to extend DS to extract events of n-ary relations and multiple arguments (rather than just binary relation). One of the hurdles for doing this is the lacking of explicit trigger information in existing KB. Our solution is to use a group of **key arguments** rather than explicit triggers to capture a particular event. For example, we can use *spouse* as the key argument to identify *marriage* events.

Our second innovation is that unlike previous studies which focus on tasks defined by the ACE evaluation framework (Ahn 2006; Li, Ji, and Huang 2013; Chen et al. 2015; Nguyen, Cho, and Grishman 2016), we propose a novel event extraction paradigm with key arguments to characterize an event type. We consider event extraction as two se-

<sup>1</sup>In the ACE task, a trigger is the word that most clearly expresses the occurrence of an event. For example, *ex* in *ex-husband* may trigger a *divorce* event.

S1: Remedy Corp was sold to BMC Software as the Service Management Business Unit in 2004.  
company acquired      acquiring company      divisions formed      date

(a) An example sentence from Wiki text

	company_acquired	acquiring_company	date	divisions_formed
m.07bh4j7	Remedy Corp	BMC Software	2004	Service Management Business Unit

(b) Entry of `business.acquisition` in Freebase

Figure 1: The event type of the sentence shown in (a) can be automatically inferred using the structural information given by Freebase in (b).

quence labeling subtasks, namely event detection and argument detection. Inspired by neural network models in sequence labeling tasks (Huang, Xu, and Yu 2015; Lample et al. 2016), we utilize BLSTM-CRF models to label key arguments and non-key arguments in each sentence separately. However, event structures are not simple sequences and there are strong dependencies among key arguments. We therefore reformulate the hypotheses as constraints, and apply linear integer programming to output multiple optimal label sequences to capture multi-typed events.

We evaluate our approach by applying it to automatically collect training data with both Freebase and Wikipedia tables. Our proposed event extractor is used to identify typed event mentions with typed roles. Our experimental results on both automatic and manual evaluations show that our approach can effectively extract a rich types of events without expert-annotated training data.

## Motivation

The task of event extraction is to (1) detect the occurrence of events with specific types and (2) extract arguments (i.e. typed participants or attributes) that are associated with an event. As a motivation example, consider a sentence from the Wiki text dataset, shown in Figure 1 (a). This sentence describes a business acquisition event. It consists of three **arguments**, *acquired company* (Remedy Corp), *acquiring company* (BMC Software), *formed division* (Service Management Business Unit) and *date* (2004). To learn an event extractor, existing supervision-based approaches all require manually identifying the event type by examining the event trigger from each individual sentence in the training data – for examples, trigger “sold” in this example indicate there is a business acquisition event.

The problem of manually labelling is that it takes a lot of time to construct a training dataset and in many cases the resulted dataset is too small to learn an effective event extractor. To make the learning process fast and cheap, we must take humans out of the loop of data labelling. If we can find a way to automatically label training data, we can then scale up the training data generation to millions of instances instead of thousands that current approaches are capable of. We discover that the structural information of

Freebase (FB) (Bollacker et al. 2008) already provide sufficient information to allow us to automatically label the event type of the particular sentence shown in Figure 1 (a). Figure 1 (b) shows a Compound Value Type (CVT) entry of FB. CVT organizes complex structured data with multiple *properties* as a table<sup>2</sup>. After examining this particular CVT schema, one can find that this schema implies there could be a *business.acquisition* event if the text contains properties (and arguments) like, `company_acquired` (Remedy Corp), `acquiring_company` (BMC Software) and `date` (2004) – all these are the identical arguments if they were extracted manually. If we can utilize such CVT information, we can then automatically label the sentence shown in Figure 1 (a) without needing to identify the trigger (which can only be achieved through human annotation if one wants to generate a high-quality training dataset).

A natural question to ask is that “How useful are CVTs?”. Later in this paper, we show that structured knowledge base, tables or lists that are used to sum up certain activities can be a useful source of supervision data generation for event extraction. For example, using just 24 CVTs, we are able to automatically generate 280 million instances (FIX:Section ??).

In this work, we develop a novel technique to automatically generate training data by exploiting the prior CVT knowledge offered by various datasets, including FB and Wikiedia pages. Our work shows that it is possible to scale up training data generation for event extraction with little human involvement. In the next section, we give an overview before presenting our approach in more details.

## Dataset Construction

Here, we employ the event-related entries of Freebase CVT tables to illustrate how to automatically annotate event mentions in Wikipedia’s articles, with the essence of distant supervision (**DS**): *A sentence that contains all key arguments of an entry in an event table (e.g., CVT table) is likely to express that event.* We will then label this sentence as a mention of this CVT event, and the words or phrases that match this entry’s properties as the involved arguments, with the roles specified by their corresponding property names.

We regard a sentence as *positive* when it mentions the occurrence of an event, or *negative* otherwise. For example, S1 and S2 are positive examples with their arguments in italics and underlined (also shown in Figure 2), while S3 and S4 are negative.

**S1:** *Remedy Corp* was sold to *BMC Software* as the *Service Management Business Unit* in *2004*.

**S2:** *Microsoft* spent \$6.3 billion buying online display advertising company *aQuantive* in *2007*.

**S3:** Microsoft hopes aQuantive’s Brian McAndrews can outfox Google.

**S4:** On April 29th, Elizabeth II and Prince Philip witnessed the marriage of Prince William.

<sup>2</sup>Therefore, we also use the term “argument to refer to a CVT property for the rest of the paper.

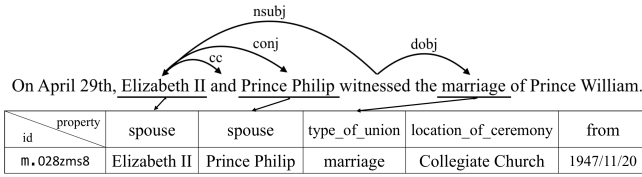


Figure 2: The dependency tree of S4, which partially matches an entry of *people.marriage*.

The selection strategy for key arguments of a given event type is based on two criteria: (1) *Key arguments should have high importance value*; (2) *Key arguments should include time-related arguments*.

The *importance value* of an argument *arg* (e.g., *date*) to its event type *cvt* (e.g., *business.acquisition*) can be defined as:

$$I_{cvt,arg} = \log \frac{\text{count}(cvt, arg)}{\text{count}(cvt) \times \text{count}(arg)} \quad (1)$$

where  $\text{count}(cvt)$  is the number of instances of type *cvt*,  $\text{count}(arg)$  is the number of times *arg* appearing in all CVT types, and  $\text{count}(cvt, arg)$  is the number of *cvt* instances that contain *arg*.

Although time-related arguments are often missing in the currently imperfect KBs, they are indeed crucial to indicate the actual occurrence of an event, e.g., S3, containing *Microsoft* as *acquiring\_company* and *aQuantive* as *company\_acquired* but without time-related arguments, will be mistakenly considered as a positive sample for event *business.acquisition*.

Intuitively, two arguments involving in the same event mention are likely to be closer within the syntactic structure. In Figure 3, both *Prince Philip* and *marriage* can be matched as key arguments in a *people.marriage* entry, but are far from each other on the dependency parse tree, thus S4 should be labeled as negative.

We conduct a series of manual evaluations on the quantity and quality of the datasets produced by different strategies (see Sec ), and our final strategy is: for each CVT, we first sort all its arguments in descending order by their importance values, and select the top half arguments as key arguments. We then include the time-related argument with highest importance value as a supplementary key argument. Finally, we eliminate sentences in which the dependency distances between any two key arguments are greater than 2.

## Our Approach

Previous event extraction systems mainly rely on explicit trigger identification to detect the occurrence of an event, which is then used to decide its event type and label its arguments. In our automatically collected dataset, where human-labeled event triggers are unavailable, we argue that **key arguments** can play the same role as explicit event triggers. We thus treat the event extraction as a pipeline of the following two steps:

- **Event detection:** to identify key arguments in a sentence. If a sentence contains **all key arguments** of a specific

event type, it will be considered to imply an event mention of this specified type.

- **Argument detection:** to identify other non-key arguments for each event in the sentence.

Take S1 as an example, in event detection, *Remedy Corp*, *BMC Software*, and *2004* could be identified as *company\_acquired*, *acquiring\_company*, and *date*, respectively, indicating that S1 may mention a *business.acquisition* event. During argument detection, *Service Management Business Unit* should be identified as *divisions\_formed*, which, together with the detected key arguments, form a full mention for a *business.acquisition* event.

## Event Detection

68% of arguments in our dataset consist of more than one word. We thus formulate each subtask in a sequence labeling paradigm rather than word-level classifications. Each word in the given sentence is tagged with the BIO scheme, where each token is labeled as B-role if it is the beginning of an event argument with its corresponding role *role*, or I-role if it is inside an argument, or O otherwise. We accordingly propose a novel BLSTM-CRF model with ILP-based inference to detect and organize key arguments.

**BLSTM** Long Short-Term Memory Network (LSTM) (Hochreiter and Schmidhuber 1997) is a natural fit for sequence labeling, which maintains a memory based on historical contextual information. Formally, given a sentence  $w = \{w_1, w_2, \dots, w_n\}$  of length  $n$ , we use  $\mathbf{x}_t$  to represent feature vector, e.g., word embeddings, corresponding to the  $t$ -th word  $w_t$ . At each time  $t$ , a forward LSTM layer takes  $\mathbf{x}_t$  as input and computes the output vector  $\vec{\mathbf{h}}_t$  of the past context, while a backward LSTM layer reads the same sentence in reverse and outputs  $\overleftarrow{\mathbf{h}}_t$  given the future context. We concatenate these two vectors to form the output vector of BLSTM, which is fed into a softmax layer to estimate a probability distribution over all possible labels.

**CRF** A straightforward way to find the label sequence for given sentence is to choose the best label for each word individually according to the BLSTM output. However, this greedy strategy ignores the dependencies between labels, thus can not guarantee the best sequence. Therefore, we introduce a CRF layer over the BLSTM output, which is admittedly effective in various sequence labeling tasks (Collobert et al. 2011; Huang, Xu, and Yu 2015).

We consider  $\mathbf{P}$  to be a matrix of confidence scores output by BLSTM, and the element  $\mathbf{P}_{i,j}$  of the matrix denotes the probability of the label  $j$  for the  $i$ -th word in a sentence. The CRF layer takes a transition matrix  $\mathbf{A}$  as parameter, where  $\mathbf{A}_{i,j}$  represents the score of a transition from label  $i$  to label  $j$ . The score of a sentence  $w$  along with a path of labels  $y = \{y_1, y_2, \dots, y_n\}$  is measured by the sum of BLSTM outputs and transition scores:

$$\text{score}(w, y) = \sum_{i=0}^n \mathbf{P}_{i,y_i} + \sum_{i=1}^n \mathbf{A}_{y_{i-1},y_i}, \quad (2)$$

During test, given a sentence  $w$ , we adopt the Viterbi algorithm (Rabiner 1989) to find the optimal label sequence with the maximum score among all possible label sequences.

**ILP-based Post Inference** Basically, event detection is a structure prediction problem, while the output sequences of BLSTM-CRF do not necessarily satisfy the structural constraints. For instance, regardless of how many key arguments are correctly identified by BLSTM-CRF, if there is one key argument missing, this detection should be considered as failed.

We thus propose to apply Integer Linear Programming (ILP) to further globally optimize the BLSTM-CRF output to produce the best label sequence. Formally, let  $\mathcal{L}$  be the set of possible argument labels. For each word  $w_i$  in the sentence  $w$  and a pair of labels  $\langle l, l' \rangle \in \mathcal{L} \times \mathcal{L}$ , we create a binary variable  $v_{i,l,l'} \in \{0, 1\}$ , denoting whether or not the  $i$ -th word  $w_i$  is tagged as label  $l$  and its following word  $w_{i+1}$  is tagged as label  $l'$  at the same time. The objective of ILP is to maximize the overall score of the variables as:

$$\sum_{i,l,l'} v_{i,l,l'} * (\mathbf{P}_{i,l} + \mathbf{A}_{l,l'}).$$

where we consider the following four constraints:

**C1:** Each word should be and only be annotated with one label, i.e.:

$$\sum_{l,l'} v_{i,l,l'} = 1 \quad (3)$$

**C2:** If the value of  $v_{i,l,l'}$  is 1, then there has to be a label  $l^*$  that will make  $v_{i+1,l',l^*}$  equal to 1, i.e.:

$$v_{i,l,l'} = \sum_{l^*} v_{i+1,l',l^*} \quad (4)$$

**C3:** If the current label is I-arg, then its previous label must be B-arg or I-arg, i.e.:

$$v_{i,\text{I-arg},l'} = v_{i-1,\text{B-arg},\text{I-arg}} + v_{i-1,\text{I-arg},\text{I-arg}} \quad (5)$$

**C4:** For a specific event type, all its key arguments should co-occur in the sentence, or none of them appears in the resulting sequence. For any pair of key arguments  $arg_1$  and  $arg_2$  with respect to the same event type, the variables related to them are subject to:

$$\sum_{i,l'} v_{i,\text{B-arg}_1,l'} \leq n * \sum_{j,l^*} v_{j,\text{B-arg}_2,l^*} \quad (6)$$

where  $n$  is the length of the sentence.

In order to address the multi-type event mention issue, we allow our ILP solver to output multiple optimal sequences. Specifically, after our model outputs the best sequence  $s^t$  at time  $t$ , we remove the previously best solutions  $\{s^1, \dots, s^t\}$  from the solution space, and re-run our solver to obtain the next optimal sequences  $s^{t+1}$ . We repeat the optimization procedure until the difference between the scores of  $s^1$  and  $s^T$  is greater than a threshold  $\lambda$ , and consider all solutions  $\{s^1, s^2, \dots, s^{T-1}\}$  as the optimal label sequences. We use Gurobi (Gurobi Optimization 2016) as our ILP solver and set  $\lambda = 0.05 \times n$ , which averagely produce 1.07 optimal sequences for each sentence.

	Train	Dev	Test
#PosSent.	29912	7477	9346
#NegSent.	50904	12726	15906
#Eve.			
#Arg.			
%Multi_Eve.			

Table 1: Statistics for the generated dataset. #PosSent. is the number of positive sentences, #NegSent. is the number of positive sentences, #Eve. is the number of event mentions, and #Arg. is the number of event arguments. %Multi\_Eve. is the ratio of multi-type events.

## Argument Detection

After event detection, a sentence will be classified into different event types, and labeled with its corresponding key arguments. Next, we will identify the remaining non-key arguments in the sentence.

We adopt the same BLSTM-CRF architecture (in Sec ) for argument detection, where we encode the event label (output of event detection) of each word into a key-argument feature vector through a look-up table, and concatenate it with the original word embedding as the input to the new BLSTM-CRF. Note that we do not need post inference here.

## Experiments

We use Freebase and the English Wikipedia dump of 2016-11-20, to construct our dataset. Statistics for the generated dataset from Freebase is shown in Table 1. Note that we only include positive sentences in this dataset.

We first manually evaluate the quality of our test set and then regard the automatically generated data as gold standard and evaluate our model accordingly. Next, we manually evaluate a subset of events detected by our model and analyze the differences with regards to the automatic evaluation. Finally, we conduct evaluation on a smaller dataset constructed according to Wikipedia tables and articles.

**Evaluation Metrics:** We evaluate our models in terms of precision (P), recall (R), and F-measure (F) for each subtask. These metrics are computed according to the following standards of correctness: For *event type classification*, an event is correctly classified if its reference sentence contains **all key arguments** of this event type. For *event detection*, an event is correctly detected if its type and all its key arguments match a reference event within the same sentence. For *argument detection*, an argument is correctly detected if its offsets, role, and related event type exactly match the reference argument within the same sentence.

**Training:** All hyper-parameters are tuned on the development set. During event detection, we set the size of word embeddings to 200, the size of LSTM layer to 100. In argument detection, we use the same size of word embedding, while the size of LSTM layer is 150, and the size of key argument embedding is 50. Word embeddings are pre-trained using skip-gram word2vec (Mikolov et al. 2013) on English

Wikipedia and fine tuned during training. We apply dropout (0.5) on both input and output layers.

### Negative samples:

### Dataset Evaluation

To investigate the possibility of automatically constructing training data for event extraction, we evaluate five datasets that utilize following strategies to determine key arguments and collect positive instances: (1) *ALL* means regarding all arguments as key arguments; (2) *IMP* means selecting the top half arguments with high importance value as key arguments; (3) *IMP&TIME* means adding a time-related argument with highest importance value to the set of key arguments defined by *IMP*; (4) *DIS* means eliminating sentence where dependency distances between any two key arguments are greater than 2. We randomly select 100 sentences from each dataset, and annotators are asked to determine whether each sentence implies a given event.

No.	Strategy	Sent.	Type	Post.
T1	<i>ALL</i>	203	9	98%
T2	<i>IMP</i>	108K	24	22%
T3	<i>IMP+DIS</i>			37%
T4	<i>IMP&amp;TIME</i>			81%
T5	<i>IMP&amp;TIME+DIS</i>			89%

Table 2: Statistics of the datasets built with different strategies. *Sent.* is the number of found sentences. *Type* the number of different CVT types found within each dataset. *Post.* is the percent of sentences mentioning the given events explicitly.

As shown in Table 2, it is not surprising that the most strict strategy, *T1*, guarantees the quality of the generated data, while we can merely obtain 203 sentences covering 9 types of events, which is insufficient for further applications. *T2* relaxes *T1* by allowing the absence of non-key arguments, which expands the resulting dataset, but introduces more noise, indicating that *T2* is inappropriate to be used as a soft constraint. Compared with *T2*, the significant quality improvement by *T4* proves that time-related arguments within CVT schemas are critical to imply an event occurrence. Among all strategies, data obtained by *T5* achieves the highest precision, while still accounting for 7,180 sentences, showing that it is feasible to automatically collect quality training data for event extraction without either human-designed event schemas or **extra** human annotations.

### Baselines

We compare our proposed models with three baselines, including traditional feature-based methods and neural network models. All baselines are trained following the two-step pipeline, i.e., event detection and argument detection. For neural network methods, we train a BLSTM model that takes word embeddings as input, and outputs the label with the maximum probability among all possible labels. For feature-based methods, we apply Conditional Random Field (Lafferty, McCallum, and Pereira 2001) (with the CRF++

toolkit (Kudo 2005) ) and Maximum Entropy (Berger, Pietra, and Pietra 1996) (Le Zhang’s MaxEnt toolkit) to explore a variety of elaborate features, such as lexical, syntactic and entity-related features, according to the state-of-art feature-based ACE event extraction system (Li, Ji, and Huang 2013). Note that during argument detection stage, we add the label of each word output by event detection as a supplementary feature.

### Automatic Evaluation

As shown in Table 3, traditional feature-based models perform worst in both event detection and argument detection. One of the main reasons is the absence of explicit event trigger annotations in our dataset, which makes it impossible to include trigger-related features, e.g., trigger-related dependency features and positions of triggers. Although traditional models can achieve higher precisions, they only extract a limited number of events, resulting in low recalls. Neural-network methods perform much better than feature-based models, especially recalls, since they can make better use of word semantic features. Specifically, BLSTM can capture longer range dependencies and richer contextual information, instead of neighboring word features only. And the CRF component brings an averagely 4% improvement in all metrics, and by adding the ILP-based post inference module, our full model, BLSTM-CRF-ILP<sub>multi</sub>, achieves the best performance among all models.

**Effect of CRF Layer** Every model with a CRF layer over its BLSTM output layer is superior to the one with a BLSTM layer only. Compared with the BLSTM, BLSTM-CRF achieves higher precisions and recalls in all subtasks by significantly reducing the invalid labeling sequences (e.g., *I-arg* appears right after *O*). During prediction, instead of tagging each token independently, BLSTM-CRF takes into account the constraints between neighboring labels, and potentially increases the cooccurrences of key arguments regarding the same event type.

**Effect of Post Inference** As shown in Table 3, ILP-based post inference considerably improves the overall system performance, especially in event type classification. With the help of constraint **C4**, some dubious key arguments can be inferred through other key arguments from their contexts. Compared with BLSTM-CRF, BLSTM-CRF-ILP<sub>1</sub> produces an F1 gain of 7.4% in event type classification, 1.8% in event detection, and 4.6% in argument detection.

**Multi-type Event Extraction** We further investigate the effect of BLSTM-CRF-ILP<sub>multi</sub>, which is the only model that can deal with the multi-type event mentions. As shown in Table 3, the proposed strategy in BLSTM-CRF-ILP<sub>multi</sub> helps detect more event mentions, contributing to the increase of recalls, and F1 scores with a little drop of precisions. Evaluated on 38 sentences containing multi-type event mentions, the F1 scores of BLSTM-CRF-ILP<sub>multi</sub> in event type classification, event detection and argument detection are 70.7%, 58.4% and 26.9%, respectively.

Model	Event Classification			Argument Detection			Event Detection		
	P	R	F	P	R	F	P	R	F
CRF									
MaxEnt									
BLSTM									
BLSTM-CRF									
BLSTM-CRF-ILP <sub>1</sub>									
BLSTM-CRF-ILP <sub>multi</sub>									

Table 3: Overall system performance of automatic evaluations (%).

## Manual Evaluation

We randomly sample 150 sentences from the test set. Annotators are asked to annotate all arguments to each sentence following two steps. First, determine whether a given sentence is positive or negative, and assign event types to positive sentences. Next, label all related arguments and their roles according to the event types for positive sentences. Two annotators will independently annotate each sentence, and discuss to achieve an agreement. The inter-annotator agreement is 87% for event types and 79% for arguments.

Model	EC	AD	ED
CRF	21.2	13.3	5.30
MaxEnt	17.7	11.7	5.44
BLSTM	79.8	64.3	41.2
BLSTM-CRF	81.3	67.9	43.2
BLSTM-CRF-ILP <sub>1</sub>	85.0	70.4	43.8
BLSTM-CRF-ILP <sub>multi</sub>	<b>85.3</b>	<b>70.8</b>	<b>44.3</b>

Table 4: The F1 scores of different systems on the manually annotated data. EC, AD, ED denote the event type classification, argument detection and event detection, respectively.

Table 4 shows the system performances in the manual evaluation. We can draw similar conclusions about the comparison of performances between different models as automatic evaluation. It is clear that BLSTM-CRF-ILP<sub>multi</sub> is the most effective model in event extraction as it achieves the highest F1 scores in both manual and automatic evaluation.

<b>S5:</b> That night, in an apparent bid to kill Amos, the car instead runs over the sheriff, leaving Chief Deputy Wade Parent (played by James Brolin) in charge.			
<b>Event Type</b>	film.performance (Wrong labeled in data generation)		
<b>Arguments</b>	actor	character	film
	James Brolin	Wade Parent	the car
<b>S6:</b> Nicholas Hammond (born May 15, 1950) is an American actor and writer who is perhaps best known for his roles as Friedrich von Trapp in the film The Sound of Music, and as Peter Parker/Spider-Man on the CBS television series The Amazing Spider-Man.			
<b>Event Type</b>	film.performance		
<b>Arguments</b>	actor	character	film
	Nicholas Hammond	Friedrich von Trapp	The Sound of Music
<b>Event Type</b>	tv.regular_tv_appearance (Missing in generated data)		
<b>Arguments</b>	actor	character	series
	Nicholas Hammond	Peter Parker/Spider-Man	The Amazing Spider-Man

Figure 3: Example outputs of BLSTM-CRF-ILP<sub>multi</sub>.

Moreover, we manually check the top 5 event types whose

EC F1 scores by BLSTM-CRF-ILP<sub>multi</sub> differ greatly between automatic evaluation and manual evaluation, summarized in Table 5. We find that most of the performance differences are caused by data generation. Figure 4 examples two types of errors in data construction. Several automatically labeled sentences do not imply any event while still matching all key properties of some CVT entries. For an example, although the phrase *the car* in S5 matches a film name, it does not indicate this film, and there is no explicit evidence indicating that an actor starred in this film. This is a bottleneck of our data generation strategy. During manual evaluation, we find 16 negative sentences which are mistakenly labeled as positive, and our model manages to rectify 6 of them.

Remarkably, our BLSTM-CRF-ILP<sub>multi</sub> model can help find more CVT instances that are not referenced in Freebase. There are two events mentioned in S6, while the arguments of the second event do not match any existing CVT instances in Freebase, which should be populated into Freebase. This phenomenon suggests that learning from distant supervision provided by Freebase, our model can help complete and update properties of Freebase instances in return.

Event type	P	R	F
olympics.medal_honor	↓ 25.0%	↓ 5.0%	↓ 13.8%
film.performance	↓ 21.4%	↑ 3.1%	↓ 10.3%
business.acquisition	→	↓ 7.1%	↓ 5.4%
tv.appearance	↓ 9.5%	↑ 3.0%	↓ 3.1%
film.release	↓ 7.7%	↑ 5.6%	↓ 0.55%

Table 5: The difference of EC F1 scores (by BLSTM-CRF-ILP<sub>multi</sub>) between automatic and manual evaluation for top 5 event types.

## Tables as Indirect Supervision

To investigate the applicability of our approach to other structured tables besides Freebase CVT tables, we automatically build a new dataset with the supervision provided by Wikipedia tables, which may characterize events about acquisition, winning of the Olympics games, and winning prestigious awards in entertainment (Table 6).

We train our BLSTM-CRF-ILP<sub>multi</sub> on this dataset and evaluate it on 100 manually annotated sentences. We can see that without extra human annotations, our model can learn to extract events from the training data weakly supervised by Wikipedia tables. Given a specific event type, as long as we can acquire tables implying events of such type, it is possible



Event type	Entr.	Sent.	EC	AD	ED
Acquisition	690	414	87.0	72.0	69.6
Olympics	2503	1460	77.2	64.5	38.6
Awards	3039	2217	95.0	82.8	58.6

Table 6: Statistics of the Table dataset and performance of our model. *Entr.* is the number of table entries. *Sent.* is number of positive instances.

to automatically collect training data from such tables, and learn to extract structured event representations of that type.

## ACE Event Extraction

### Update Freebase with Extracted Events

#### Related Work

Most event extraction works are within the tasks defined by several evaluation frameworks (e.g., MUC (Grishman and Sundheim 1996), ACE (Doddington et al. 2004), ERE (Song et al. 2015) and TAC-KBP (Mitamura et al. 2015)), all of which can be considered as a template-filling-based extraction task. These frameworks focus on limited number of event types, which are designed and annotated by human experts and hard to generalize to other domains. Furthermore, existing extraction systems, which usually adopt a supervised learning paradigm, have to rely on those high-quality training data within those frameworks, thus hard to move to more domains in practice, regardless of feature-based (Gupta and Ji 2009; Hong et al. 2011; Li, Ji, and Huang 2013) or neural-network-based methods (Chen et al. 2015; Nguyen, Cho, and Grishman 2016).

Besides the works focusing on small human-labeled corpus, Huang et al. (2016) propose a novel Liberal Event Extraction paradigm which automatically discovers event schemas and extract events simultaneously from any unlabeled corpus. In contrast, we propose to exploit existing structured knowledge bases, e.g., Freebase, to automatically discover types of events as well as their corresponding argument settings, without expert annotation, and further automatically construct training data, with the essence of distant supervision (Mintz et al. 2009).

Distant supervision (DS) has been widely used in binary relation extraction, where the key assumption is that sentences containing both the subject and object of a  $\langle \text{subj}, \text{rel}, \text{obj} \rangle$  triple can be seen as its support, and further used to train a classifier to identify the relation *rel*. However, this assumption does not fit to our event extraction scenario, where an event usually involves several arguments and it is hard to collect enough training sentences with all arguments appearing in, as indicated by the low coverage of *T1*. We therefore investigate different hypotheses for event extraction within the DS paradigm and propose to utilize time and syntactic clues to refine the DS assumption for better data quality. We further relieve the reliance on event trigger annotations by previous event extractors, and define a novel event extraction paradigm with key arguments to characterize an event type.

## Conclusions and Future Work

In this paper, we propose a novel event extraction paradigm without expert-designed event templates by leveraging structured knowledge bases or tables to automatically acquire event schema and corresponding training data. We propose a BLSTM-CRF model with ILP-based post inference to extract multi-typed event mentions without explicit trigger annotations. Experimental results on both manual and automatic evaluations show that it is possible to learn to identify both typed events and typed roles with indirect supervision from Freebase or Wikipedia tables. In the future, we will first further investigate how to better characterize key arguments for a certain event type, and extend our work to update a structured knowledge base according to news events.

## References

- Aguilar, J.; Beller, C.; McNamee, P.; Van Durme, B.; Strassel, S.; Song, Z.; and Ellis, J. 2014. A comparison of the events and relations across ace, ere, tac-kbp, and framenet annotation standards. In *Proceedings of the Second Workshop on EVENTS: Definition, Detection, Coreference, and Representation*, 45–53.
- Ahn, D. 2006. The stages of event extraction. In *Proceedings of the Workshop on Annotating and Reasoning about Time and Events*, 1–8. Association for Computational Linguistics.
- Berger, A. L.; Pietra, V. J. D.; and Pietra, S. A. D. 1996. A maximum entropy approach to natural language processing. *Computational linguistics* 22(1):39–71.
- Bollacker, K.; Evans, C.; Paritosh, P.; Sturge, T.; and Taylor, J. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD 2008*, 1247–1250. ACM.
- Chen, Y.; Xu, L.; Liu, K.; Zeng, D.; and Zhao, J. 2015. Event extraction via dynamic multi-pooling convolutional neural networks. In *ACL 2015*, volume 1, 167–176.
- Collobert, R.; Weston, J.; Bottou, L.; Karlen, M.; Kavukcuoglu, K.; and Kuksa, P. 2011. Natural language processing (almost) from scratch. *Journal of Machine Learning Research* 12(Aug):2493–2537.
- Doddington, G. R.; Mitchell, A.; Przybicki, M. A.; Ramshaw, L. A.; Strassel, S.; and Weischedel, R. M. 2004. The automatic content extraction (ace) program-tasks, data, and evaluation. In *LREC 2004*, volume 2, 1.
- Grishman, R., and Sundheim, B. 1996. Message understanding conference-6: A brief history. In *COLING 1996*, volume 96, 466–471.
- Gupta, P., and Ji, H. 2009. Predicting unknown time arguments based on cross-event propagation. In *ACL-IJCNLP 2009*, 369–372. Association for Computational Linguistics.
- Gurobi Optimization, I. 2016. Gurobi optimizer reference manual.
- Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.
- Hong, Y.; Zhang, J.; Ma, B.; Yao, J.; Zhou, G.; and Zhu, Q. 2011. Using cross-entity inference to improve event extrac-

tion. In *ACL 2011*, 1127–1136. Association for Computational Linguistics.

Huang, L.; Cassidy, T.; Feng, X.; Ji, H.; Voss, C.; Han, J.; and Sil, A. 2016. Liberal event extraction and event schema induction. In *ACL 2016*.

Huang, Z.; Xu, W.; and Yu, K. 2015. Bidirectional lstm-crf models for sequence tagging. *arXiv preprint arXiv:1508.01991*.

Kudo, T. 2005. Crf++: Yet another crf toolkit. *Software available at <http://crfpp.sourceforge.net>*.

Lafferty, J.; McCallum, A.; and Pereira, F. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *ICML 2001*, volume 1, 282–289.

Lample, G.; Ballesteros, M.; Subramanian, S.; Kawakami, K.; and Dyer, C. 2016. Neural architectures for named entity recognition. *arXiv preprint arXiv:1603.01360*.

Li, Q.; Ji, H.; and Huang, L. 2013. Joint event extraction via structured prediction with global features. In *ACL 2013*, 73–82.

Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G. S.; and Dean, J. 2013. Distributed representations of words and phrases and their compositionality. In *NIPS 2013*, 3111–3119.

Mintz, M.; Bills, S.; Snow, R.; and Jurafsky, D. 2009. Distant supervision for relation extraction without labeled data. In *ACL-IJCNLP 2009*, 1003–1011. Association for Computational Linguistics.

Mitamura, T.; Yamakawa, Y.; Holm, S.; Song, Z.; Bies, A.; Kulick, S.; and Strassel, S. 2015. Event nugget annotation: Processes and issues. In *Proc. of the Workshop on EVENTS at the NAACL-HLT*, 66–76.

Nguyen, T. H.; Cho, K.; and Grishman, R. 2016. Joint event extraction via recurrent neural networks. In *NAACL-HLT 2016*, 300–309.

Rabiner, L. R. 1989. A tutorial on hidden markov models and selected applications in speech recognition. *Proc. of the IEEE* 77(2):257–286.

Song, Z.; Bies, A.; Strassel, S.; Riese, T.; Mott, J.; Ellis, J.; Wright, J.; Kulick, S.; Ryant, N.; and Ma, X. 2015. From light to rich ere: annotation of entities, relations, and events. In *Proc. of the Workshop on EVENTS at the NAACL-HLT*, 89–98.

Zeng, D.; Liu, K.; Chen, Y.; and Zhao, J. 2015. Distant supervision for relation extraction via piecewise convolutional neural networks. In *EMNLP 2015*, 1753–1762.