

DEGREE: A Data-Efficient Generation-Based Event Extraction Model

Anonymous ACL submission

Abstract

Due to the high cost of human annotations, learning a *data-efficient* event extraction model that can be trained with *only a few* labeled examples has become a crucial challenge. In this paper, we focus on *low-resource end-to-end* event extraction. We propose DEGREE, a model that formulates event extraction as a conditional generation problem. Given a passage and a manually designed prompt, DEGREE learns to summarize the event happening in the passage into a natural sentence that follows a predefined pattern. The final event structure predictions are then extracted from the generated sentence with a deterministic algorithm. DEGREE has the following advantages to learn well with less training data. First, with our design of prompts, DEGREE obtains semantic guidance by leveraging *label semantics* and thus better captures the argument roles. In addition, the proposed model is capable of using additional *weakly-supervised* information, such as the description of events. Finally, learning triggers and argument roles in an *end-to-end manner* encourages the model to better utilize the shared knowledge and dependencies between them. Our experimental results and ablation studies demonstrate the strong performance of DEGREE for low-resource event extraction.

1 Introduction

Event extraction (EE) aims to extract different types of events, each of which includes a trigger and several participants (arguments) with specific roles, from the given passage. For example, in Figure 1, a *Justice:Execute* event is triggered by the word “*execution*” and this event contains three argument roles, including an *Agent* (*Indonesia*) who carries out the execution, a *Person* been executed (*convicts*), and a *Place* where the event occurs (not mentioned in the passage). Prior works usually divide EE into two subtasks (Wadden et al., 2019; Lin et al., 2020; Fincke et al., 2021): (1) **event**

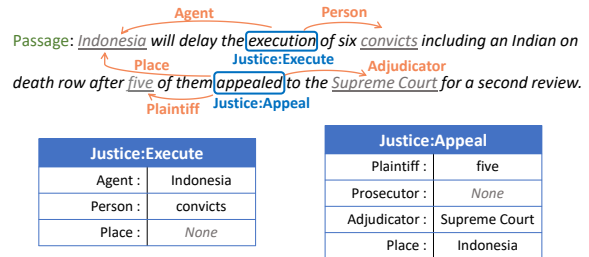


Figure 1: Two examples of events (*Justice:Execute* and *Justice:Appeal*) extracted from the given passage.

detection, which identifies the event triggers and their types, and (2) **event argument extraction**, which extracts the participants (arguments) of the event and their roles when given an event type and the corresponding event trigger.

Several previous EE approaches rely on a large amount of annotated data for training (Nguyen and Grishman, 2015; Nguyen et al., 2016; Du and Cardie, 2020; Paolini et al., 2021). However, these high-quality event annotations are expensive to be obtained. For example, the ACE 2005 corpus (Doddington et al., 2004), one of the most common EE datasets, requires two rounds of annotations by linguistics experts. The high annotation costs make these models hard to be extended to new domains and new event types. Therefore, how to learn a *data-efficient* EE model trained with *only a few* annotated examples is a crucial research question.

In this paper, we focus on *low-resource* event extraction, where only a small amount of training examples are available during training. As illustrated in Figure 2, we propose DEGREE (**D**ata-**E**fficient **G**eneRative **E**vent **E**xtraction), a generation-based model that takes a passage and a manually designed prompt as the input, and learns to summarize the passage into a natural sentence following a predefined template. The event triggers and arguments can then be extracted from the generated sentence by using a deterministic algorithm.

DEGREE enjoys the following three advantages to learn well with less training data. First, the gen-

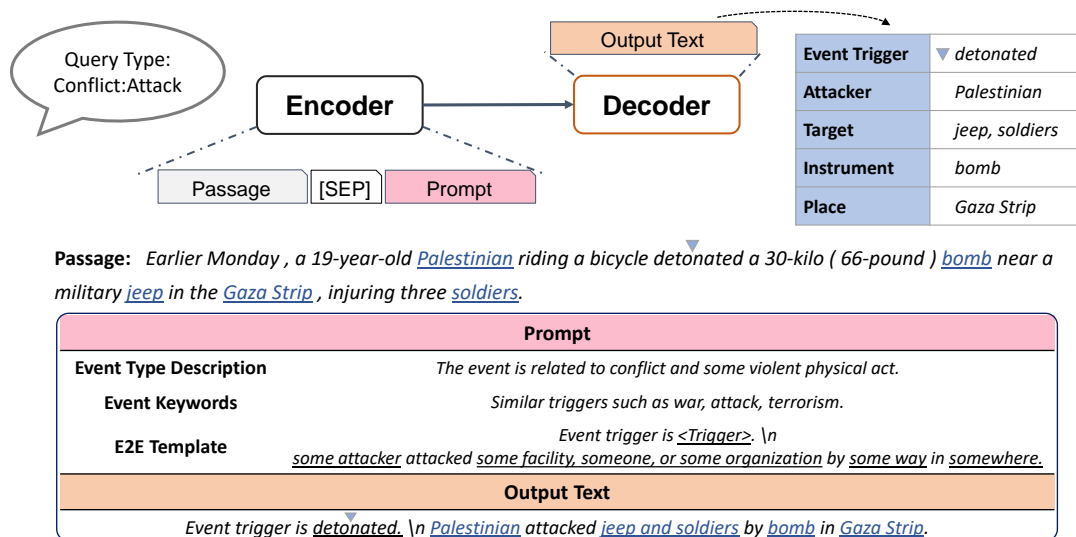


Figure 2: An illustration of DEGREE for predicting a *Contact:Attack* event. The input of DEGREE consists of the given passage and our design prompt that contains a event type description, some event keywords, and a E2E template. DEGREE is trained to generate an output to fill in the placeholders (underlined words) in the E2E template with triggers and arguments. The final event prediction is then decoded from the generated output.

eration framework provides *label semantics* with the help of the designed template in the prompts. As the example in Figure 2 shows, the word “*some-where*” in the prompt guides the model to predict words being similar to *location* for the role *Place*. Also, the word “*attacked*” in the prompt depicts the relationship between the role *Attacker* and the role *Target*. With these kinds of guidance, DEGREE can make accurate predictions without many training examples. Second, the prompts can be further extended to include additional weakly-supervised information about the task, such as the description of the event and similar keywords.¹ This information facilitates DEGREE to learn under the low-resource situation. Finally, DEGREE is designed for end-to-end event extraction and can solve event detection and event argument extraction at the same time. Utilizing the shared knowledge and dependencies between the two tasks makes DEGREE more data-efficient.

Prior approaches on EE usually have only one or two above-mentioned advantages. For example, previous classification-based models (Nguyen et al., 2016; Wang et al., 2019; Yang et al., 2019; Wadden et al., 2019; Lin et al., 2020) are hard to handle label semantics and utilize the weakly-supervised information. Recently proposed generation-based models solve event extraction in a pipeline fashion; therefore, they cannot leverage

¹These resources are usually readily available. In our experiments, we take the weak supervision signals from the annotation guideline, which is provided along with the dataset.

the shared knowledge between subtasks (Paolini et al., 2021; Li et al., 2021). In addition, their generated outputs are not natural sentences, which hinders the utilization of label semantics (Paolini et al., 2021; Lu et al., 2021). As a result, DEGREE can achieve significantly better performance than prior approaches on low-resource event extraction, as we will demonstrate in Section 3.

Our contributions can be summarized as follows:

- We propose DEGREE, a generation-based end-to-end event extraction model that learns well with less data by better incorporating label semantics and shared knowledge (Section 2).
- We conduct experiments on ACE 2005 (Doddington et al., 2004) and ERE-EN (Song et al., 2015) to demonstrate the strong performance of DEGREE in the low-resource setting (Section 3).
- We present comprehensive ablation studies in both the low-resource setting and high-resource setting to better understand the advantages and the disadvantages of our model (Section 4).

2 Data-Efficient Event Extraction

We introduce DEGREE, a generation-based model for low-resource event extraction. Unlike previous works (Wadden et al., 2019; Lin et al., 2020), which separate event extraction into two pipelined tasks (event detection and event argument extraction), DEGREE is designed for the end-to-end event extraction and makes trigger predictions and argument predictions at the same time.

2.1 DEGREE

We formulate event extraction as a conditional generation problem. As illustrated in Figure 2, given a passage and our designed prompt, DEGREE generates an output following a particular format. The final predictions of event triggers and argument roles can be then parsed from the generated output with a deterministic algorithm. Compared to the previous classification-based models (Wang et al., 2019; Yang et al., 2019; Wadden et al., 2019; Lin et al., 2020), the generation framework provides a flexible way to include additional information and guidance. By designing appropriate prompts, we encourage DEGREE to better capture the dependencies between entities and therefore reduce the number of needed training examples.

The desired prompt not only provides information but also defines the output format. As shown in Figure 2, it contains the following components:

- **Event type definition** describes the definition for the given event type.² For example, “*The event is related to conflict and some violent physical act.*” describes a *Conflict:Attack* event.
- **Event keywords** presents some words that are semantically related to the given event type. For instance, *war*, *attack*, and *terrorism* are three event keywords for the *Conflict:Attack* event. In practice, we collect three words that appear as the triggers in the example sentences from the annotation guidelines.
- **E2E template** defines the expected output format and can be separated into two parts. The first part is called ED template, which is designed as “*Event trigger is <Trigger>*”, where “<Trigger>” is a special token serving as a placeholder. The second part is the EAE template, which differs based on the given event type. For example, in Figure 2, the EAE template for the *Conflict:Attack* event is “*some attacker attacked some facility, someone, or some organization by some way in somewhere*”. Each underlined string starting with “*some-*” serves as a placeholder corresponding to an argument role in the *Conflict:Attack* event. For instance, “*some way*” corresponds to the role *Instrument* and “*somewhere*” corresponds to the role *Place*. Notice that every event type has its own EAE template. The full list of EAE templates and the constructing details can be found in Appendix A.

²The definition can be derived from the annotation guidelines, which are provided along with the datasets.

Training. The training objective of DEGREE is to generate an output that replaces the placeholders in E2E template with the gold labels. Take Figure 2 as an example, DEGREE is expected to replace “<Trigger>” with the gold trigger (*detonated*), replace “*some attacker*” with the gold argument for role *Attacker (Palestinian)*, and replace “*some way*” with the gold argument for role *Instrument (bomb)*. If there are multiple arguments for the same role, they are concatenated with “*and*”; if there is no predicted argument for one role, the model should keep the corresponding placeholder (i.e., “*some-*” in the E2E template). For the case that there are multiple triggers for the given event type, DEGREE will generate the E2E template multiple times such that each E2E template corresponds to each trigger and its argument roles. We put more training details in Appendix B.

Inference. We enumerate all event types and generate an output for each event type. Then, we compare the generated output with the placeholders in E2E template to determine the predicted trigger spans and predicted argument spans. Finally, we apply string matching to convert the word spans to the offsets in the passage. If the predicted span appears in the passage multiple times, we choose all that match for trigger predictions and choose the one being closest to the given trigger span for argument predictions.

Discussion. Notice that the E2E template plays an important role for DEGREE. First, it serves as the control signal and defines the expected output format. Second, it provides label semantics to help DEGREE make accurate predictions. Those placeholders (words starting with “*some-*”) in the E2E template give DEGREE some hints about the entity types of arguments. For instance, when seeing “*somewhere*”, DEGREE tends to generate a location rather than a person. In addition, the words other than “*some-*” describe the relationships between roles. For example, DEGREE knows the relation between the role *Attacker* and the role *Target* (who is attacking and who is attacked) because of the word “*attacked*” in E2E template. This guidance makes DEGREE learn the dependencies between entities well with less training data.

Unlike previous generation-based approaches (Paolini et al., 2021; Li et al., 2020), we intentionally write the E2E templates in natural sentences. This not only utilizes label semantics better but also

232	makes the model easier to leverage the knowledge	in the ED template with event triggers.	281
233	from the pre-trained decoder. In Section 4, we will		
234	provide experiments to demonstrate the advantage	DEGREE(EAE). The prompt of DEGREE(EAE)	282
235	of using natural sentences.	contains the following components:	283
236	We want to point out one advantage of using	• Event type definition is the same as the one	284
237	generation-based models under the low-resource	for DEGREE.	285
238	scenario compared to previous classification-based	• Query trigger is a string that indicates the trig-	286
239	event extraction models — generation-based mod-	ger word for the given event type. For example,	287
240	els do not require named entity annotations (Sha	“The event trigger word is detonated” points out	288
241	et al., 2018; Lin et al., 2020). The pre-trained de-	that “detonated” is the given trigger.	289
242	coder inherently identifies reasonable entity spans,	• EAE template is an event-type-specific tem-	290
243	which makes generation-based models become a	plate mentioned previously. It is actually the	291
244	good choice when annotations are expensive.	second part of the E2E template. The full list	292
245		of EAE templates can be found in Appendix A.	293
246	Cost of template constructing. DEGREE does	Similar to DEGREE, the goal for DEGREE(EAE) is	294
247	require human effort to design the templates; how-	to generate an outputs that replace the placeholders	295
248	ever, writing those templates is much easier and	in EAE template with event arguments.	296
249	more effortless than collecting complicated event	In Section 3, we will compare DEGREE with	297
250	annotations. As shown in Appendix A, we keep	DEGREE(PIPE) to study the benefit of dealing with	298
251	the EAE templates as simple and short as possi-	event extraction in an end-to-end manner under the	299
252	ble. Therefore, it takes only about one minute for	low-resource setting.	300
253	people who are not linguistic experts to compose		
254	a template. In fact, several prior works (Liu et al.,	3 Experiments	301
255	2020; Du and Cardie, 2020; Li et al., 2020) also use		
256	constructed templates as weakly-supervised signals	We conduct experiments for <i>low-resource</i> event	302
257	to improve models. In Section 4, we will study how	extraction to study how DEGREE performs.	303
258	different templates affect the performance.		
259		3.1 Experimental Settings	304
260	2.2 DEGREE in Pipeline Framework		
261	DEGREE is flexible and can be easily modified	Datasets. We consider ACE 2005 (Doddington	305
262	to DEGREE(PIPE), which first focuses event de-	et al., 2004) and follow the pre-processing in Wad-	306
263	tection (ED) and then solves event argument ex-	den et al. (2019) and Lin et al. (2020), resulting	307
264	traction (EAE). DEGREE(PIPE) consists of two	in two variants: ACE05-E and ACE05-E⁺ . Both	308
265	models: (1) DEGREE(ED), which aims to exact	contain 33 event types and 22 argument roles. In	309
266	event triggers for the given event type, and (2) DE-	addition, we consider ERE-EN (Song et al., 2015)	310
267	GREE(EAE), which identifies argument roles for	and adopt the pre-processing in Lin et al. (2020),	311
268	the given event type and the corresponding trig-	which keeps 38 event types and 21 argument roles.	312
269	ger. DEGREE(ED) and DEGREE(EAE) are similar		
270	to DEGREE but with different prompts and output	Data split for low-resource setting. We gener-	313
271	formats. We describe the difference as follows.	ate different proportions (1%, 2%, 3%, 5%, 10%,	314
272		20%, 30%, and 50%) of training data to study the	315
273	DEGREE(ED). The prompt of DEGREE(ED)	influence of the size of training set and use the origi-	316
274	contains the following components:	nal dev set and test set for evaluation. Appendix C	317
275		lists more details about the split generating process	318
276	• Event type definition is the same as the ones	and the data statistics.	319
277	for DEGREE.		
278	• Event keywords is the same as the one for DE-	Evaluation metrics. We consider the same crite-	320
279	GREE.	ria in prior works (Wadden et al., 2019; Lin et al.,	321
280	• ED template is designed as “Event trigger is	2020). (1) Trigger F1-score: an trigger is correctly	322
	<Trigger>”. It is actually the first part of the	identified (Tri-I) if its offset matches the gold one;	323
	E2E template.	it is correctly classified (Tri-C) if its event type also	324
		matches the gold one. (2) Argument F1-score: an	325
	Similar to DEGREE, the objective of DEGREE(ED)	argument is correctly identified (Arg-I) if its offset	326
	is to generate an output that replaces “<Trigger>”		

Trigger Classification F1-Score (%)																			
Model	Type	ACE05-E						ACE05-E ⁺						ERE-EN					
		1%	3%	5%	10%	20%	30%	1%	3%	5%	10%	20%	30%	1%	3%	5%	10%	20%	30%
BERT_QA	Cls	20.5	40.2	42.5	50.1	61.5	61.3	-	-	-	-	-	-	-	-	-	-	-	-
OneIE	Cls	38.5	52.4	59.3	61.5	<u>67.6</u>	<u>67.4</u>	39.0	52.5	60.6	58.1	<u>66.5</u>	66.4	11.0	36.9	<u>46.7</u>	48.8	51.8	53.5
Text2Event	Gen	14.2	35.2	46.4	47.0	55.6	60.7	15.7	38.4	43.9	46.3	56.5	62.0	6.3	25.6	33.5	42.4	46.7	50.1
TANL	Gen	34.1	48.1	53.4	54.8	61.8	61.6	30.3	50.9	53.1	55.7	60.8	61.7	5.7	30.8	43.4	45.9	49.0	49.3
DEGREE(PIPE)	Gen	<u>55.1</u>	62.8	<u>63.8</u>	66.1	64.4	64.4	56.4	<u>62.5</u>	<u>61.1</u>	<u>62.3</u>	<u>62.5</u>	67.1	32.7	<u>44.5</u>	41.6	<u>50.6</u>	51.1	53.5
DEGREE	Gen	55.4	<u>62.1</u>	65.8	<u>65.8</u>	68.3	68.2	<u>49.5</u>	63.5	<u>62.3</u>	68.5	67.6	<u>66.9</u>	<u>27.9</u>	45.5	47.0	53.0	<u>51.7</u>	53.5

Argument Classification F1-Score (%)																			
Model	Type	ACE05-E						ACE05-E ⁺						ERE-EN					
		1%	3%	5%	10%	20%	30%	1%	3%	5%	10%	20%	30%	1%	3%	5%	10%	20%	30%
BERT_QA	Cls	4.7	14.5	26.9	27.6	36.7	38.8	-	-	-	-	-	-	-	-	-	-	-	-
OneIE	Cls	9.4	22.0	26.8	26.8	42.7	47.8	10.4	20.6	29.7	35.5	<u>46.7</u>	48.0	2.6	20.3	29.7	35.1	40.7	43.0
Text2Event	Gen	3.9	12.2	19.1	24.9	32.3	39.2	5.7	16.5	21.3	26.4	35.2	42.1	2.3	15.2	23.6	28.7	35.7	38.7
TANL	Gen	8.5	17.2	24.7	29.0	34.0	39.2	8.6	22.3	30.4	29.2	34.6	39.0	1.4	20.2	29.5	30.1	35.6	36.9
DEGREE(PIPE)	Gen	<u>13.1</u>	<u>26.1</u>	<u>27.6</u>	42.1	40.7	44.0	<u>16.0</u>	<u>26.4</u>	29.9	<u>39.5</u>	41.3	<u>48.5</u>	<u>12.2</u>	29.7	<u>31.4</u>	<u>39.4</u>	<u>41.9</u>	42.2
DEGREE	Gen	21.7	30.1	35.5	<u>41.6</u>	46.2	48.7	18.7	34.0	35.7	43.6	48.9	51.2	14.5	<u>28.9</u>	33.4	41.7	42.9	45.5

Table 1: Trigger classification F1-scores and argument classification F1-scores for low-resource event extraction. Highest scores are in bold and the second best scores are underlined. ‘‘Cls’’ and ‘‘Gen’’ represent classification-based models and generation-based models, respectively. If the model is a pipelined model, then its argument predictions are based on its predicted triggers. DEGREE achieves a much better performance than other baselines. The performance gap becomes more significant for the extremely low-resource situation.

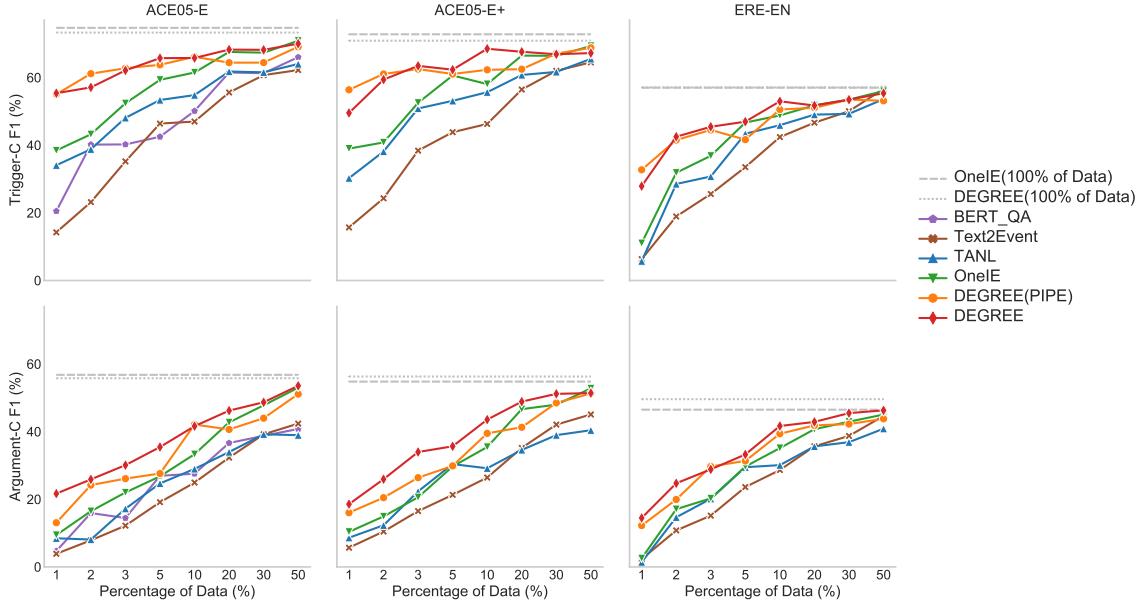


Figure 3: Trigger classification F1-scores and argument classification F1-scores for low-resource event extraction. DEGREE achieves a much better performance than other baselines. The performance gap becomes more significant for the extremely low-resource situation.

and event type match the gold ones; it is correctly classified (Arg-C) if its role matches as well.

Compared baselines. We consider the following classification-based models: (1) **OneIE** (Lin et al., 2020), the current state-of-the-art (SOTA) EE model trained with designed global features. (2) **BERT_QA** (Du and Cardie, 2020), which views EE tasks as a sequence of extractive question answering problems. Since it learns a classifier to indicate the position of the predicted span, we view it as a classification model. We also consider

the following generation-based models: (3) **TANL** (Paolini et al., 2021), which treats EE tasks as translation tasks between augmented natural languages. (4) **Text2Event** (Lu et al., 2021), a sequence-to-structure model that convert the input passage to a tree-like event structure. Notice that the outputs of both generation-based baselines are not natural sentences. Therefore, it is more difficult for them to utilize the label semantics. All the implementation details can be found in Appendix D. It is worth noting that we train OneIE with named entity an-

notations, as the original papers suggest, while the other models are trained without entity annotations.

3.2 Main Results

Table 1 shows the trigger classification F1-scores and the argument classification F1-scores across three datasets with different proportions of training data. The results are visualized in Figure 3. Since our task is *end-to-end* event extraction, the argument classification F1-score is the more important metric that we considered when comparing models.

From the figure and the table, we can observe that both DEGREE and DEGREE(PIPE) outperform all other baselines when using less than 10% of the training data. The performance gap becomes much more significant under the extremely low data situation. For example, when only 1% of training data is available, DEGREE and DEGREE(PIPE) achieve more than 15 points of trigger classification F1-scores improvement and more than 5 points of argument classification F1-scores. This demonstrates the effectiveness of our design. The generation-based model with carefully designed prompts is able to utilize the label semantics and the additional weakly-supervised signals, thus, helps the learning under the low-resource regime.

Another interesting finding is that DEGREE and DEGREE(PIPE) seem to be more beneficial to argument prediction than trigger prediction. For instance, OneIE, the strongest baseline, requires 20% of training data to achieve competitive performance on trigger prediction to DEGREE and DEGREE(PIPE); however, it requires about 50% of training data to achieve competitive performance on argument prediction. The reason is that the ability to capture dependencies becomes more important for argument prediction than trigger prediction since arguments are usually strongly dependent on each other compared to triggers. Therefore, the improvements of our models for argument prediction are more significant.

Finally, we observe that DEGREE is slightly better than DEGREE(PIPE) under the low-resource setting. We hypothesize that DEGREE jointly predicts triggers and arguments and therefore can better take advantage of the output dependencies.

3.3 High-Resource Event Extraction

While we focus on data-efficient learning for low-resource event extraction, to better understand the advantages and disadvantages of our model and make sure that it is indeed more data-efficient,

Model	Type	ACE05-E		ACE05-E ⁺		ERE-EN	
		Tri-C	Arg-C	Tri-C	Arg-C	Tri-C	Arg-C
dbRNN*	Cls	69.6	50.1	-	-	-	-
DyGIE++	Cls	70.0	50.0	-	-	-	-
Joint3EE*	Cls	69.8	52.1	-	-	-	-
BERT_QA*	Cls	72.4	53.3	-	-	-	-
MQAEE*	Cls	71.7	53.4	-	-	-	-
OneIE*	Cls	74.7	56.8	72.8	54.8	57.0	46.5
TANL	Gen	68.4	47.6	68.6	46.0	54.7	43.2
Text2Event*	Gen	71.9	53.8	<u>71.8</u>	54.4	59.4	48.3
BART-Gen*	Gen	71.1	53.7	-	-	-	-
DEGREE(PIPE)	Gen	72.2	<u>55.8</u>	71.7	56.8	<u>57.8</u>	50.4
DEGREE	Gen	<u>73.3</u>	<u>55.8</u>	70.9	<u>56.3</u>	57.1	<u>49.6</u>

Table 2: Results for high-resource event extraction. Highest scores are in bold and the second best scores are underlined. *We report the numbers from the original paper. DEGREE has a competitive performance to the SOTA model (OneIE) and outperform other baselines.

Model	Type	ACE05-E		ACE05-E ⁺		ERE-EN	
		Arg-I	Arg-C	Arg-I	Arg-C	Arg-I	Arg-C
DyGIE++	Cls	66.2	60.7	-	-	-	-
BERT_QA*	Cls	68.2	65.4	-	-	-	-
OneIE	Cls	73.2	69.3	73.3	70.6	75.3	70.0
TANL	Gen	65.9	61.0	66.3	62.3	75.6	69.6
BART-Gen*	Gen	69.9	66.7	-	-	-	-
DEGREE(EAE)	Gen	76.0	73.5	75.2	73.0	80.2	76.3

Table 3: Results for high-resource event argument extraction. Models predict arguments based on the given gold triggers. Best scores are in bold. *We report the numbers from the original paper. DEGREE(EAE) achieves a new state-of-the-art performance on event argument extraction.

rather than simply a stronger model, we additionally study DEGREE in the high-resource setting for controlled comparisons.

Compared baselines. In addition to the previously mentioned EE models: **OneIE** (Lin et al., 2020), **BERT_QA** (Du and Cardie, 2020), **TANL** (Paolini et al., 2021), and **Text2Event** (Lu et al., 2021), we also consider the following baselines focusing on the high-resource setting. **dbRNN** (Sha et al., 2018) is classification-based model that adds dependency bridges for event extraction. **DyGIE++** (Wadden et al., 2019) is a classification-based model with span graph propagation technique. **Joint3EE** (Nguyen and Nguyen, 2019) is classification-based model jointly trained with entity, trigger, and argument annotations. **MQAEE** (Li et al., 2020) converts EE to a series of argument extraction question answering problems. **BART-Gen** (Li et al., 2021) is a generation-based model focusing on *only* event argument extraction.³ Appendix D shows the implementation details.

³We follow the original paper and use TAPKEY as their event detection model.

Results for event extraction. Table 2 shows the results of high-resource event extraction. In terms of trigger predictions (Tri-C), DEGREE and DEGREE(PIPE) outperforms all the baselines except for OneIE, the current state-of-the-art model. For argument predictions (Arg-C), our models have slightly better performance than OneIE in two out of the three datasets. When enough training examples are available, models can learn more sophisticated features from data, which do not necessarily follow the learned dependencies. Therefore, the advantage of DEGREE over DEGREE(PIPE) becomes less obvious. This result justifies our hypothesis that DEGREE has better performance for the *low-resource setting* because of its ability to better capture dependencies.

Results for event argument extraction. In Table 3, we additionally study the performance for event argument extraction task, where the model makes argument predictions *with the gold trigger provided*. Interestingly, DEGREE(EAE) achieves pretty strong performance and outperforms other baselines with a large margin. Combining the results in Table 2, we hypothesize that event argument extraction is a more challenging task than event trigger detection and it requires more training examples to learn well. Hence, our proposed model, which takes the advantage of using label semantics to better capture dependencies, achieves a new state-of-the-art for event argument extraction.

4 Ablation Study

In this section, we present comprehensive ablation studies to justify our design. To better understand the contribution of each component in the designed prompt and their effects on the different tasks, we ablate DEGREE(EAE) and DEGREE(ED) for both low-resource and high-resource situations.

Impacts of components in prompts. Table 4 lists the performance changes when removing the components in the prompts for event detection on ACE05-E. The performance decreases whenever removing any one of event type definition, event keywords, and ED template. The results suggest that three components are all necessary.

Table 5 demonstrates how different components in prompts affect the performance of event argument extraction on ACE05-E. Removing any one of event type definition, query trigger, and EAE template leads to performance drops,

Model	10% Data		100% Data	
	Tri-I	Tri-C	Tri-I	Tri-C
Full DEGREE(ED)	69.3	66.1	75.4	72.2
- w/o Event type definition	67.9	64.4	73.5	70.1
- w/o ED template	68.8	65.8	74.0	70.5
- w/o Event keywords	68.2	64.0	73.5	69.1
- only Event type definition	66.3	63.5	72.6	68.9
- only Event keywords	69.2	63.8	70.8	66.2

Table 4: Ablation study for the components in the prompt on event detection with ACE05-E.

Model	10% Data		100% Data	
	Arg-I	Arg-C	Arg-I	Arg-C
Full DEGREE(EAE)	63.3	57.3	76.0	73.5
- w/o Event type definition	60.3	54.4	74.5	71.1
- w/o EAE template	57.0	51.9	73.8	70.4
- w/o Query trigger	55.2	49.9	71.4	69.0
- only Query trigger	51.9	48.1	71.2	69.4
- only EAE template	51.2	46.9	71.4	68.6
- only Event type definition	46.7	42.3	71.4	68.2

Table 5: Ablation study for the components in the prompt on event argument extraction with ACE05-E.

which validates their necessity. We observe that query trigger plays the most important role among the three and when less training data is given, the superiority of leveraging any of these weakly-supervised signal becomes more obvious.

Effects of different template designs. To verify the importance of using natural sentences as outputs, we study three variants of EAE templates:

- **Natural sentence.** Our proposed templates described in Section 2, e.g., “*somebody was born in somewhere.*”, where “*somebody*” and “*some-where*” are placeholders that can be replaced by the corresponding arguments.
- **Natural sentence with special tokens.** It is similar to the natural sentence one except for using role-specific special tokens instead of “some-” words. For example, “*<Person> was born in <Place>.*” We consider this to study the *label semantics* of roles.
- **HTML-like sentence with special tokens.** To study the importance of using natural sentence, we also consider HTML-like sentence, e.g., “*<Person> </Person> <Place> </Place>.*” The model aims to put argument predictions between the corresponding HTML tags.

The results of all variants of EAE templates on ACE05-E are shown in Table 6. We notice that writing templates in a natural language style get better performance, especially when only a few data is available (10% of data). This shows our design’s capability to leverage pre-trained knowledge

Model	10% Data		100% Data	
	Arg-I	Arg-C	Arg-I	Arg-C
OneIE	48.3	45.4	73.2	69.3
BART-Gen	-	-	69.9	66.7
Natural sentence	63.3	57.3	76.0	73.5
Natural sentence w/ special tokens	59.8	55.5	74.7	72.3
HTML-like sentence w/ special tokens	60.8	51.9	74.6	71.4

Table 6: The performances of DEGREE(EAE) on ACE05-E with different types of templates.

Model	10% Data		100% Data	
	Arg-I	Arg-C	Arg-I	Arg-C
OneIE	48.3	45.4	73.2	69.3
BART-Gen	-	-	69.9	66.7
DEGREE(EAE)	63.3	57.3	76.0	73.5
DEGREE(EAE) + variant template 1	61.6	55.5	73.4	70.4
DEGREE(EAE) + variant template 2	63.9	56.9	75.5	72.5

Table 7: Study on the effect of different template constructing rules. Experiments is conducted on ACE05-E.

in the generation process. Additionally, there are over 1 F1 score performance drops when replacing natural language placeholders with special tokens. This confirms that leveraging label semantics for different roles is beneficial.

Sensitivity to template design. Finally, we study how sensitive our model is to the template. In addition to the original design of templates for event argument extraction, we compose other two sets of templates with different constructing rules (e.g., different word choices and different orders of roles). Table 7 shows the results of using different sets of templates. We observe a performance fluctuation when using different templates, which indicates that the quality of templates does affect the performance to a certain degree. Therefore, we need to be cautious when designing templates. However, even though our model could be sensitive to the template design, it still outperforms OneIE and BART-Gen, which are the best classification-based model and the best generation-based baseline, respectively.

5 Related Work

Fully supervised event extraction. Event extraction has been studied for over a decade (Ahn, 2006; Ji and Grishman, 2008) and most traditional event extraction works follow the fully supervised setting (Nguyen et al., 2016; Sha et al., 2018; Nguyen and Nguyen, 2019; Yang et al., 2019; Lin et al., 2020; Liu et al., 2020; Li et al., 2020). Many of them use classification-based models and use pipeline-style frameworks to extract events (Nguyen et al., 2016; Yang et al., 2019; Wadden et al., 2019). To better leverage shared knowledge in event triggers and arguments, some works

propose to incorporate global features to jointly decide triggers and arguments (Lin et al., 2020; Li et al., 2013; Yang and Mitchell, 2016). Recently, few generation-based event extraction models have been proposed. TANL (Paolini et al., 2021) treats event extraction as translation tasks between augmented natural languages. Their predicted target—augmented language embed labels into the input passage via using brackets and vertical bar symbols, hindering the model from fully leveraging label semantics. BART-Gen (Li et al., 2021) is also a generation-based model focusing on document-level event argument extraction. Yet, similar to TANL, they solve event extraction with a pipeline, which prevents knowledge sharing across subtasks. All these fully supervised methods can achieve substantial performance with a large amount of annotated data. However, their designs are not specific for low-resource scenarios, hence, these models can not enjoy all the benefits that DEGREE obtains for low-resource event extraction at the same time, as we mentioned in Section 1.

Low-resource event extraction. It has been a rising interest in event extraction under less data scenario. Liu et al. (2020) uses a machine reading comprehension formulation to conduct event extraction in a low-resource regime. Text2Event (Lu et al., 2021), a sequence-to-structure generation paradigm, first presents events in a linearized format, and then trains a generative model to generate the linearized event sequence. Text2Event’s unnatural output format hinders the model from fully leveraging pre-trained knowledge. Hence, their model falls short on the cases with only extremely low data being available (as shown in Section 3).

Another thread of works are using meta-learning to deal with the less label challenge (Deng et al., 2020; Shen et al., 2021; Cong et al., 2021). However, their methods can only be applied to event detection, which differs from our main focus on studying end-to-end event extraction.

6 Conclusion

In this paper, we present DEGREE, a data-efficient generation-based event extraction model. DEGREE requires less training data because it better utilizes label semantics as well as weakly-supervised information, and captures better dependencies by jointly predicting triggers and arguments. Our experimental results and ablation studies show the superiority of DEGREE for low-resource event extraction.

585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639

References

David Ahn. 2006. The stages of event extraction. In *Proceedings of the Workshop on Annotating and Reasoning about Time and Events*.

Xin Cong, Shiyao Cui, Bowen Yu, Tingwen Liu, Yubin Wang, and Bin Wang. 2021. Few-shot event detection with prototypical amortized conditional random field. In *Findings of the Association for Computational Linguistics: ACL/IJCNLP*.

Shumin Deng, Ningyu Zhang, Jiaojian Kang, Yichi Zhang, Wei Zhang, and Huajun Chen. 2020. Meta-learning with dynamic-memory-based prototypical network for few-shot event detection. In *The Thirteenth ACM International Conference on Web Search and Data Mining (WSDM)*.

George R. Doddington, Alexis Mitchell, Mark A. Przybocki, Lance A. Ramshaw, Stephanie M. Strassel, and Ralph M. Weischedel. 2004. The automatic content extraction (ACE) program - tasks, data, and evaluation. In *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC)*.

Xinya Du and Claire Cardie. 2020. Event extraction by answering (almost) natural questions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Steven Fincke, Shantanu Agarwal, Scott Miller, and Elizabeth Boschee. 2021. Language model priming for cross-lingual event extraction. *arXiv preprint arXiv:2109.12383*.

Heng Ji and Ralph Grishman. 2008. Refining event extraction through cross-document inference. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics (ACL)*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*.

Fayuan Li, Weihua Peng, Yuguang Chen, Quan Wang, Lu Pan, Yajuan Lyu, and Yong Zhu. 2020. Event extraction as multi-turn question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings (EMNLP)*.

Qi Li, Heng Ji, and Liang Huang. 2013. Joint event extraction via structured prediction with global features. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (ACL)*.

Sha Li, Heng Ji, and Jiawei Han. 2021. Document-level event argument extraction by conditional generation. In *Proceedings of the 2021 Conference of the North*

American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT).

Ying Lin, Heng Ji, Fei Huang, and Lingfei Wu. 2020. A joint neural model for information extraction with global features. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*.

Jian Liu, Yubo Chen, Kang Liu, Wei Bi, and Xiaojiang Liu. 2020. Event extraction as machine reading comprehension. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.

Yaojie Lu, Hongyu Lin, Jin Xu, Xianpei Han, Jialong Tang, Annan Li, Le Sun, Meng Liao, and Shaoyi Chen. 2021. Text2event: Controllable sequence-to-structure generation for end-to-end event extraction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL/IJCNLP)*.

Thien Huu Nguyen, Kyunghyun Cho, and Ralph Grishman. 2016. Joint event extraction via recurrent neural networks. In *The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*.

Thien Huu Nguyen and Ralph Grishman. 2015. Event detection and domain adaptation with convolutional neural networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics (ACL)*.

Trung Minh Nguyen and Thien Huu Nguyen. 2019. One for all: Neural joint modeling of entities and events. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference (AAAI)*.

Giovanni Paolini, Ben Athiwaratkun, Jason Krone, Jie Ma, Alessandro Achille, Rishita Anubhai, Cícero Nogueira dos Santos, Bing Xiang, and Stefano Soatto. 2021. Structured prediction as translation between augmented natural languages. In *9th International Conference on Learning Representations (ICLR)*.

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 the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI)*.

695 Shirong Shen, Tongtong Wu, Guilin Qi, Yuan-Fang Li,
696 Gholamreza Haffari, and Sheng Bi. 2021. Adaptive
697 knowledge-enhanced bayesian meta-learning for few-
698 shot event detection. In *Findings of the Association
699 for Computational Linguistics: ACL/IJCNLP*.

700 Zhiyi Song, Ann Bies, Stephanie M. Strassel, Tom
701 Riese, Justin Mott, Joe Ellis, Jonathan Wright, Seth
702 Kulick, Neville Ryant, and Xiaoyi Ma. 2015. From
703 light to rich ERE: annotation of entities, relations,
704 and events. In *Proceedings of the The 3rd Workshop
705 on EVENTS: Definition, Detection, Coreference, and
706 Representation, (EVENTS@HLP-NAACL)*.

707 David Wadden, Ulme Wennberg, Yi Luan, and Han-
708 naneh Hajishirzi. 2019. Entity, relation, and event
709 extraction with contextualized span representations.
710 In *Proceedings of the 2019 Conference on Empirical
711 Methods in Natural Language Processing and the 9th
712 International Joint Conference on Natural Language
713 Processing (EMNLP-IJCNLP)*.

714 Xiaozhi Wang, Ziqi Wang, Xu Han, Zhiyuan Liu, Juanzi
715 Li, Peng Li, Maosong Sun, Jie Zhou, and Xiang Ren.
716 2019. HMEAE: hierarchical modular event argument
717 extraction. In *Proceedings of the 2019 Conference
718 on Empirical Methods in Natural Language Process-
719 ing and the 9th International Joint Conference on
720 Natural Language Processing (EMNLP-IJCNLP)*.

721 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
722 Chaumond, Clement Delangue, Anthony Moi, Pier-
723 ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,
724 Joe Davison, Sam Shleifer, Patrick von Platen, Clara
725 Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven
726 Le Scao, Sylvain Gugger, Mariama Drame, Quentin
727 Lhoest, and Alexander Rush. 2020. Huggingface’s
728 transformers: State-of-the-art natural language pro-
729 cessing. In *Proceedings of the 2020 Conference on
730 Empirical Methods in Natural Language Processing:
731 System Demonstrations*.

732 Bishan Yang and Tom M. Mitchell. 2016. Joint extrac-
733 tion of events and entities within a document context.
734 In *The 2016 Conference of the North American Chap-
735 ter of the Association for Computational Linguistics:
736 Human Language Technologies (NAACL-HLT)*.

737 Sen Yang, Dawei Feng, Linbo Qiao, Zhigang Kan, and
738 Dongsheng Li. 2019. Exploring pre-trained language
739 models for event extraction and generation. In *Pro-
740 ceedings of the 57th Conference of the Association
741 for Computational Linguistics (ACL)*.

A EAE Template Constructing

Our strategy to create an EAE template is first identifying all valid argument roles for the event type,⁴ such as *Attacker*, *Target*, *Instrument*, and *Place* roles. Then, for each argument role, according to the semantics of the role type, we select natural and fluent words to form its placeholder (e.g., *some way* for *Instrument*). This design aims to provide a simple way to help the model learn both the roles' label semantics and the *event structure*. Finally, we create a natural language sentence that connects all these placeholders. Notice that we try to keep the template as simple and short as possible. Table 8 lists all designed EAE templates for ACE05-E and ACE05-E⁺. The EAE templates of ERE-EN can be found in Table 9.

B Training Details of Proposed Model

Given a passage, its annotated event types are considered as positive event types. During training, we additionally sample m event types that are not related to the passage as the negative examples, where m is a hyper-parameter. In our experiments, m is usually set to 13 or 15.

For all of DEGREE, DEGREE(ED), and DEGREE(EAE), we fine-tune the pre-trained BART-large (Lewis et al., 2020) with Huggingface package (Wolf et al., 2020). The number of parameters is around 406 millions. We train DEGREE with our machine that equips 128 AMD EPYC 7452 32-Core Processor, 4 NVIDIA A100 GPUs, and 792G RAM. We consider AdamW optimizer (Loshchilov and Hutter, 2019) with learning rate set to 10^{-5} and the weight decay set to 10^{-5} . We set the batch size to 6 for DEGREE(EAE) and 32 for DEGREE(ED) and DEGREE. The number of training epochs is 45. It takes around 2 hours, 18 hours, 22 hours to train DEGREE(EAE), DEGREE(ED), and DEGREE, respectively.

We do hyper-parameter search on m , the number of negative examples, from $\{3, 5, 7, 10, 13, 15, 18, 21\}$, and our preliminary trials shows that m less than 10 are usually less useful. For the learning rate and the weight decay, we tune it based on our preliminary experiment for event argument extraction from $\{10^{-5}, 10^{-4}\}$, while they are both fixed to 10^{-5} for all the experiments.

⁴The valid roles for each event type are predefined in the event ontology for each dataset, or can be decided by the user of interest.

C Datasets

We consider ACE 2005⁵ (Doddington et al., 2004) and ERE⁶ (Song et al., 2015). Both consider *LDC User Agreement for Non-Members*⁷ as the licenses. Both datasets are created for entity, relation, and event extraction while our focus is only event extraction in this paper. In the original ACE 2005 dataset, it contains data for English, Chinese, and Arabic and we only take the English data for our experiment. In the original ERE dataset, it contains data for English, and Chinese and we only take the English data for our experiment as well.

Because both datasets contain event like *Justice:Execute* and *Life:Die*, it is possible that some offensive words (e.g., killed) would appear in the passage. Also, some real names may appear in the passage as well (e.g., Palestinian president, Mahmoud Abbas). How to accurately identify these kinds of information is part of the goal of the task. Therefore, we do not take any changes on the datasets for protecting or anonymizing.

We split the training data based on documents, which is a more realistic setup compared to splitting data by instance. Table 10 lists the statistics of ACE05-E, ACE05-E⁺, and ERE-EN. Specifically, we try to make each proportion of data contain as many event types as possible.

D Implementation Details

This section describes the implementation details for all baselines we use. We run the experiments with three different random seeds and report the best value.

- **DyGIE++**: we use their released pre-trained model⁸ for evaluation.
- **OneIE**: we use their provided code⁹ to train the model with default parameters.
- **BERT_QA**: we use their provided code¹⁰ to train the model with default parameters.
- **TANL**: we use their provided code¹¹ to train the

⁵<https://catalog.ldc.upenn.edu/LDC2006T06>

⁶<https://catalog.ldc.upenn.edu/LDC2020T19>

⁷<https://catalog.ldc.upenn.edu/license/ldc-non-members-agreement.pdf>

⁸<https://github.com/dwadden/dygiepp>

⁹<http://blender.cs.illinois.edu/software/oneie/>

¹⁰<https://github.com/xinyadu/eeqa>

¹¹<https://github.com/amazon-research/tanl>

828 model. We conduct the experiments with two
829 variations: (1) using their default parameters,
830 and (2) using their default parameters but with
831 more training epochs. We observe that the sec-
832 ond variant works better. As a result, we report
833 the number obtained from the second setting.

- 834 • **Text2Event**: we use their official code¹² to train
835 the model with the provided parameter setting.
- 836 • **dbRNN**: we directly report the experimental
837 results from their paper.
- 838 • **Joint3EE**: we directly report the experimental
839 results from their paper.
- 840 • **MQAEE**: we directly report the experimental
841 results from their paper.
- 842 • **BART-Gen**: we report the experimental results
843 from their released appendix¹³.

¹²<https://github.com/luyaojie/Text2Event>

¹³https://github.com/raspberryyice/gen-arg/blob/main/NAACL_2021_Appendix.pdf

Event Type	EAE Template
Life:Be-Born	somebody was born in somewhere.
Life:Marry	somebody got married in somewhere.
Life:Divorce	somebody divorced in somewhere.
Life:Injure	somebody or some organization led to some victim injured by some way in somewhere.
Life:Die	somebody or some organization led to some victim died by some way in somewhere.
Movement:Transport	something was sent to somewhere from some place by some vehicle. somebody or some organization was responsible for the transport.
Transaction:Transfer-Ownership	someone got something from some seller in somewhere.
Transaction:Transfer-Money	someone paid some other in somewhere.
Business:Start-Org	somebody or some organization launched some organization in somewhere.
Business:Merge-Org	some organization was merged.
Business:Declare-Bankruptcy	some organization declared bankruptcy.
Business:End-Org	some organization dissolved.
Conflict:Attack	some attacker attacked some facility, someone, or some organization by some way in somewhere.
Conflict:Demonstrate	some people or some organization protest at somewhere.
Contact:Meet	some people or some organization met at somewhere.
Contact:Phone-Write	some people or some organization called or texted messages at somewhere.
Personnel:Start-Position	somebody got new job and was hired by some people or some organization in somewhere.
Personnel:End-Position	somebody stopped working for some people or some organization at somewhere.
Personnel:Nominate	somebody was nominated by somebody or some organization to do a job.
Personnel:Elect	somebody was elected a position, and the election was voted by some people or some organization in somewhere.
Justice:Arrest-Jail	somebody was sent to jailed or arrested by somebody or some organization in somewhere.
Justice:Release-Parole	somebody was released by some people or some organization from somewhere.
Justice:Trial-Hearing	somebody, prosecuted by some other, faced a trial in somewhere. The hearing was judged by some adjudicator.
Justice:Charge-Indict	somebody was charged by some other in somewhere. The adjudication was judged by some adjudicator.
Justice:Sue	somebody was sued by some other in somewhere. The adjudication was judged by some adjudicator.
Justice:Convict	somebody was convicted of a crime in somewhere. The adjudication was judged by some adjudicator.
Justice:Sentence	somebody was sentenced to punishment in somewhere. The adjudication was judged by some adjudicator.
Justice:Fine	some people or some organization in somewhere was ordered by some adjudicator to pay a fine.
Justice:Execute	somebody was executed by somebody or some organization at somewhere.
Justice:Extradite	somebody was extradited to somewhere from some place. somebody or some organization was responsible for the extradition.
Justice:Acquit	somebody was acquitted of the charges by some adjudicator.
Justice:Pardon	somebody received a pardon from some adjudicator.
Justice:Appeal	some other in somewhere appealed the adjudication from some adjudicator.

Table 8: All EAE templates for ACE05-E and ACE05-E⁺.

Event Type	EAE Template
Life:Be-Born	somebody was born in somewhere.
Life:Marry	somebody got married in somewhere.
Life:Divorce	somebody divorced in somewhere.
Life:Injure	somebody or some organization led to some victim injured by some way in somewhere.
Life:Die	somebody or some organization led to some victim died by some way in somewhere.
Movement:Transport-Person	somebody was moved to somewhere from some place by some way. somebody or some organization was responsible for the movement.
Movement:Transport-Artifact	something was sent to somewhere from some place. somebody or some organization was responsible for the transport.
Business:Start-Org	somebody or some organization launched some organization in somewhere.
Business:Merge-Org	some organization was merged.
Business:Declare-Bankruptcy	some organization declared bankruptcy.
Business:End-Org	some organization dissolved.
Conflict:Attack	some attacker attacked some facility, someone, or some organization by some way in somewhere.
Conflict:Demonstrate	some people or some organization protest at somewhere.
Contact:Meet	some people or some organization met at somewhere.
Contact:Correspondence	some people or some organization contacted each other at somewhere.
Contact:Broadcast	some people or some organization made announcement to some publicity at somewhere.
Contact:Contact	some people or some organization talked to each other at somewhere.
Manufacture:Artifact	something was built by somebody or some organization in somewhere.
Personnel:Start-Position	somebody got new job and was hired by some people or some organization in somewhere.
Personnel:End-Position	somebody stopped working for some people or some organization at somewhere.
Personnel:Nominate	somebody was nominated by somebody or some organization to do a job.
Personnel:Elect	somebody was elected a position, and the election was voted by somebody or some organization in somewhere.
Transaction:Transfer-Ownership	The ownership of something from someone was transferred to some other at somewhere.
Transaction:Transfer-Money	someone paid some other in somewhere.
Transaction:Transaction	someone give some things to some other in somewhere.
Justice:Arrest-Jail	somebody was sent to jailed or arrested by somebody or some organization in somewhere.
Justice:Release-Parole	somebody was released by somebody or some organization from somewhere.
Justice:Trial-Hearing	somebody, prosecuted by some other, faced a trial in somewhere. The hearing was judged by some adjudicator.
Justice:Charge-Indict	somebody was charged by some other in somewhere. The adjudication was judged by some adjudicator.
Justice:Sue	somebody was sued by some other in somewhere. The adjudication was judged by some adjudicator.
Justice:Convict	somebody was convicted of a crime in somewhere. The adjudication was judged by some adjudicator.
Justice:Sentence	somebody was sentenced to punishment in somewhere. The adjudication was judged by some adjudicator.
Justice:Fine	some people or some organization in somewhere was ordered by some adjudicator to pay a fine.
Justice:Execute	somebody was executed by somebody or some organization at somewhere.
Justice:Extradite	somebody was extradicted to somewhere from some place. somebody or some organization was responsible for the extradition.
Justice:Acquit	somebody was acquitted of the charges by some adjudicator.
Justice:Pardon	somebody received a pardon from some adjudicator.
Justice:Appeal	somebody in somewhere appealed the adjudication from some adjudicator.

Table 9: All EAE templates for ERE-EN.

Dataset	Split	#Docs	#Sents	#Events	#Event Types	#Args	#Arg Types
ACE05-E	Train (full)	529	17172	4202	33	4859	22
	Train (1%)	5	103	47	14	65	16
	Train (2%)	10	250	77	17	104	16
	Train (3%)	15	451	119	23	153	17
	Train (5%)	25	649	212	27	228	21
	Train (10%)	50	1688	412	28	461	21
	Train (20%)	110	3467	823	33	936	22
	Train (30%)	160	5429	1368	33	1621	22
	Train (50%)	260	8985	2114	33	2426	22
	Dev	28	923	450	21	605	22
Test	40	832	403	31	576	20	
ACE05-E ⁺	Train (full)	529	19216	4419	33	6607	22
	Train (1%)	5	92	49	15	75	16
	Train (2%)	10	243	82	19	129	16
	Train (3%)	15	434	124	24	203	19
	Train (5%)	25	628	219	27	297	21
	Train (10%)	50	1915	428	29	629	21
	Train (20%)	110	3834	878	33	1284	22
	Train (30%)	160	6159	1445	33	2212	22
	Train (50%)	260	10104	2231	33	3293	22
	Dev	28	901	468	22	759	22
Test	40	676	424	31	689	21	
ERE-EN	Train (full)	396	14736	6208	38	8924	21
	Train (1%)	4	109	61	14	78	16
	Train (2%)	8	228	128	21	183	19
	Train (3%)	12	419	179	26	272	19
	Train (5%)	20	701	437	31	640	21
	Train (10%)	40	1536	618	37	908	21
	Train (20%)	80	2848	1231	38	1656	21
	Train (30%)	120	4382	1843	38	2632	21
	Train (50%)	200	7690	3138	38	4441	21
	Dev	31	1209	525	34	730	21
Test	31	1163	551	33	822	21	

Table 10: Dataset statistics. Our experiments are conducted in sentences, which were split from documents. In the table, “#Docs” means the number of documents; “#Sents” means the number of sentences, “#Events” means the number of events; “#Event Types” means the number of event types in total; “#Args” means the number of argument in total; “#Arg Types” means the number of argument role types in total.

E Few-Shot and Zero-Shot Event Extraction

In order to further test our models’ generalizability, we additionally conduct zero-shot and few-shot experiments on the ACE05-E dataset with DEGREE(ED) and DEGREE(EAE).

Settings. We first select the top n common event types as “seen” types and use the rest as “unseen/rare” types, where the top common types are listed in Table 11. To simulate a zero-shot scenario, we remove all events with “unseen/rare” types from the training data. To simulate a few-shot scenario, we keep only k event examples for each “unseen/rare” type (denoted as k -shot). During the evaluation, we calculate micro F1-scores only for these “unseen/rare” types.

n	Seen Event Types for Training/Development
5	Conflict:Attack, Movement:Transport, Life:Die, Contact:Meet, Personnel:Elect
10	Conflict:Attack, Movement:Transport, Life:Die, Contact:Meet, Personnel:Elect, Life:Injure, Personnel:End-Position, Justice:Trial-Hearing, Contact:Phone-Write, Transaction:Transfer-Money

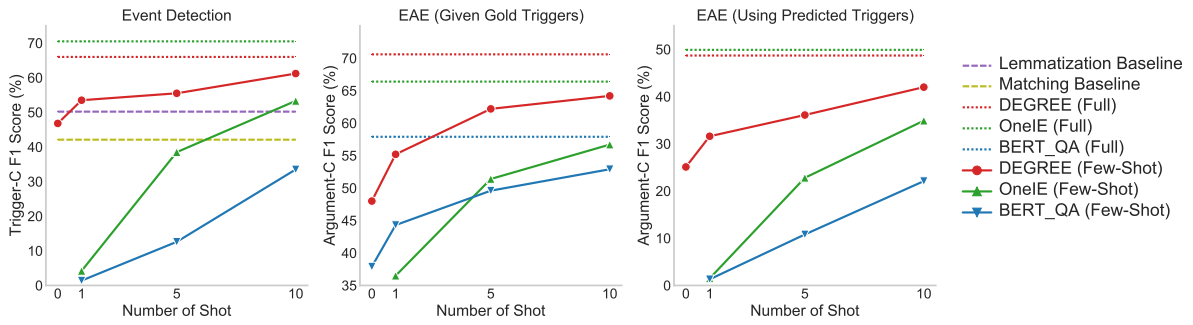
Table 11: Common event types in ACE05-E.

Compared baselines. We consider the following baselines: (1) **BERT_QA** (Du and Cardie, 2020) (2) **OneIE** (Lin et al., 2020) (3) **Matching baseline**, a proposed baseline that makes trigger predictions by performing string matching between the input passage and the event keywords. (4) **Lemma-tization baseline**, another proposed baseline that performs string matching on lemmatized input passage and the event keywords. (Note: (3) and (4) are baselines only for event detection tasks.)

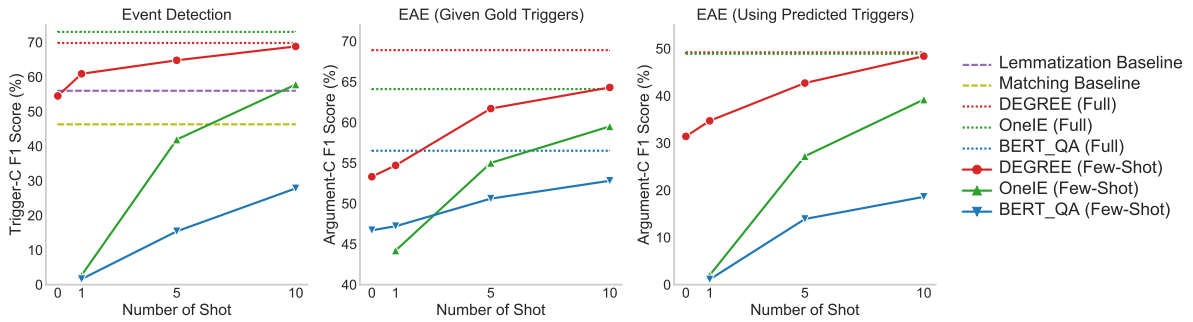
Experimental results. Figure 4, Table 12, and Table 13 show the results of $n = 5$ and $n = 10$. From the two subfigures in the left column, we see that DEGREE(ED) achieves promising results in the zero-shot setting. In fact, it performs better than BERT_QA trained in the 10-shot setting and OneIE trained in the 5-shot setting. This demonstrates the great potential of DEGREE(ED) to discover new event types. Interestingly, we observe that our two proposed baselines perform surprisingly well, suggesting that the trigger annotations in ACE05-E are actually not diverse. Despite their impressive performance, DEGREE(ED) still outperforms the matching baseline by over 4.7% absolute trigger classification F1 in both $n = 5$ and

$n = 10$ cases in zero-shot scenario. Additionally, with only one training instance for each unseen type, DEGREE(ED) can outperform both proposed baselines.

Next, we compare the results for the event argument extraction task. From the two middle subfigures, we observe that when given gold triggers, our model performs much better than all baselines with a large margin. Lastly, we train models for both trigger and argument extraction and report the final argument classification scores in the two right subfigures. We justify that our model has strong generalizability to unseen event types and it can outperform BERT_QA and OneIE even when they are both trained in 5-shot settings.



(a) Results for top common 5 event types.



(b) Results for top common 10 event types.

Figure 4: The zero/few-shot experimental results. **Left:** The result for the models on event detection task with the scores reported in trigger classification F1. **Middle:** The models are tested under the scenario of given gold trigger and evaluated with argument classification criterion. **Right:** The results for the models to perform event extraction task, which aims to predict triggers and their corresponding arguments (we report the argument classification F1).

Event Extraction									
Trigger	Argument	Common 5				Common 10			
		Tri-I	Tri-C	Arg-I	Arg-C	Tri-I	Tri-C	Arg-I	Arg-C
Matching Baseline		42.7	42.1	-	-	46.3	46.3	-	-
Lemmatization Baseline		51.5	50.2	-	-	56.6	56.0	-	-
BERT_QA 1-shot		10.0	1.4	1.3	1.3	8.2	1.6	1.1	1.1
BERT_QA 5-shot		14.0	12.6	11.1	10.8	20.8	15.4	14.6	13.9
BERT_QA 10-shot		37.8	33.5	22.9	22.1	32.0	27.8	19.5	18.6
OneIE 1-shot		4.2	4.2	1.5	1.5	4.1	2.7	2.0	2.0
OneIE 5-shot		39.3	38.5	24.8	22.8	41.9	41.9	29.7	27.2
OneIE 10-shot		54.8	53.3	36.0	34.9	61.5	57.8	41.4	39.2
DEGREE(ED) 0-shot	DEGREE(EAE) 0-shot	53.3	46.8	29.6	25.1	60.9	54.5	42.0	31.4
DEGREE(ED) 1-shot	DEGREE(EAE) 1-shot	60.1	53.3	38.8	31.6	61.2	60.9	41.1	34.7
DEGREE(ED) 5-shot	DEGREE(EAE) 5-shot	57.8	55.5	40.6	36.1	65.8	64.8	45.3	42.7
DEGREE(ED) 10-shot	DEGREE(EAE) 10-shot	63.8	61.2	46.0	42.0	72.1	68.8	52.5	48.4
OneIE (Full)		72.7	70.5	52.3	49.9	74.5	73.0	51.2	48.9
DEGREE(ED) (Full)	DEGREE(EAE) (Full)	68.4	66.0	51.9	48.7	72.0	69.8	52.5	49.2

Table 12: Full results of zero/few-shot event extraction on ACE05-E.

Event Argument Extraction									
Trigger	Argument	Common 5				Common 10			
		Tri-I	Tri-C	Arg-I	Arg-C	Tri-I	Tri-C	Arg-I	Arg-C
Gold Triggers	BERT_QA 0-shot	100.0	100.0	55.8	37.9	100.0	100.0	57.2	46.7
Gold Triggers	BERT_QA 1-shot	100.0	100.0	55.8	44.3	100.0	100.0	57.8	47.2
Gold Triggers	BERT_QA 5-shot	100.0	100.0	56.6	49.6	100.0	100.0	59.1	50.6
Gold Triggers	BERT_QA 10-shot	100.0	100.0	58.8	52.9	100.0	100.0	60.5	52.8
Gold Triggers	OneIE 1-shot	100.0	100.0	40.9	36.5	100.0	100.0	48.3	44.2
Gold Triggers	OneIE 5-shot	100.0	100.0	55.6	51.4	100.0	100.0	58.6	55.0
Gold Triggers	OneIE 10-shot	100.0	100.0	59.4	56.7	100.0	100.0	62.0	59.5
Gold Triggers	DEGREE(EAE) 0-shot	100.0	100.0	56.1	48.0	100.0	100.0	66.5	53.3
Gold Triggers	DEGREE(EAE) 1-shot	100.0	100.0	65.2	55.2	100.0	100.0	65.4	54.7
Gold Triggers	DEGREE(EAE) 5-shot	100.0	100.0	70.9	62.2	100.0	100.0	68.0	61.7
Gold Triggers	DEGREE(EAE) 10-shot	100.0	100.0	71.1	64.2	100.0	100.0	71.6	64.3
Gold Triggers	BERT_QA (Full)	100.0	100.0	63.1	57.9	100.0	100.0	62.1	56.5
Gold Triggers	OneIE (Full)	100.0	100.0	70.8	66.4	100.0	100.0	67.9	64.1
Gold Triggers	DEGREE(EAE) (Full)	100.0	100.0	74.5	70.6	100.0	100.0	73.6	68.9

Table 13: Full results of zero/few-shot event argument extraction on ACE05-E.

900 **F Limitations and Potential Risks**

901 **Limitations.** DEGREE assumes that some
902 weakly-supervised information (the description
903 of events, similar keywords, and human-written
904 templates) is accessible and not expensive. We
905 believe this assumption holds for most of common
906 NLP tasks. However, for some specific domains,
907 such as the biomedical domain, acquiring this
908 information can be a bit difficult (e.g., needs to hire
909 experts to write down templates), which increases
910 the cost of training DEGREE. In addition, our
911 proposed model is based on pre-trained language
912 models. DEGREE performs well because it is
913 able to leverage the prompts and the pre-trained
914 knowledge. However, if the downstream domain is
915 far from the pre-trained corpus, the advantage of
916 leveraging knowledge becomes restricted.

917 Due to the high cost of annotations, there are
918 not many public datasets for event extraction. DE-
919 GREE achieves a good performance on two datasets
920 (ACE 2005 and ERE-EN), which are more related
921 to news-styled passages. When considering other
922 downstream domains, it is possible that the im-
923 provement is not as significant as it is for the two
924 datasets we use in the paper. The reason is the
925 gap between the downstream domain knowledge
926 and the pre-trained knowledge, as mentioned in the
927 previous paragraph.

928 **Potential risks.** DEGREE fine-tunes the pre-
929 trained generative language model (Lewis et al.,
930 2020). Therefore, the generated output is poten-
931 tially affected by the corpus for pre-training. Al-
932 though with a low possibility, it is possible for
933 our model to accidentally generate some malicious,
934 counterfactual, and biased sentences, which may
935 cause ethics concerns. We suggest carefully exam-
936 ining those potential issues before deploying the
937 model in any real-world applications.