

EventBERT: Incorporating Event-based Semantics for Natural Language Understanding

Anni Zou^{1,2,3}, Zhuosheng Zhang^{1,2,3}, Hai Zhao^{1,2,3,*}

¹ Department of Computer Science and Engineering, Shanghai Jiao Tong University

² Key Laboratory of Shanghai Education Commission for Intelligent Interaction and Cognitive Engineering, Shanghai Jiao Tong University

³ MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University
{annie0103, zhangzs}@sjtu.edu.cn, zhaohai@cs.sjtu.edu.cn

Abstract

Natural language understanding tasks require a comprehensive understanding of natural language and further reasoning about it, on the basis of holistic information at different levels to gain comprehensive knowledge. In recent years, pre-trained language models (PrLMs) have shown impressive performance in natural language understanding. However, they rely mainly on extracting context-sensitive statistical patterns without explicitly modeling linguistic information, such as semantic relationships entailed in natural language. In this work, we propose EventBERT, an event-based semantic representation model that takes BERT as the backbone and refines with event-based structural semantics in terms of graph convolution networks. EventBERT benefits simultaneously from rich event-based structures embodied in the graph and contextual semantics learned in pre-trained model BERT. Experimental results on the GLUE benchmark show that the proposed model consistently outperforms the baseline model.

1 Introduction

Recent years have witnessed deep pre-trained language models (PrLM) such as ELMo (Peters et al., 2018), BERT (Devlin et al., 2019), XLNet (Yang et al., 2019) and ERNIE (Sun et al., 2020) significantly prospering the performance of a wide range of natural language understanding (NLU) tasks. The remarkable advancements brought by PrLM have shown the effectiveness of leveraging contextualized representation. However, they mainly rest on extracting context-sensitive statistical patterns without explicitly modeling linguistic information such as semantic relationships in natural language.

It is clear that natural language itself abounds with ample, multi-level linguistic information. Although PrLMs like BERT implicitly represent linguistic knowledge more or less (Rogers et al., 2020), studies disclose that linguistic knowledge is far from fully absorbed (Ettinger, 2020; Rogers et al., 2020). Therefore, there emerges a series of derivatives of PrLM intending to fuse explicit linguistic knowledge so as to acquire better language representation, including syntactic (Bai et al., 2021; Xu et al., 2021; Zhang et al., 2020b) and semantic information (Zhang et al., 2020a; Guo et al., 2020b; Guan et al., 2021).

In cognition practice, human needs to distill semantics of different levels to gain a comprehensive understanding, whereas neural language models learn semantic representation to deal with downstream tasks (Geeraerts and Cuyckens, 2007). Thus, effective learning of semantic knowledge plays a crucial role in NLU tasks and has gained growing attention recently. For instance, Zhang et al. (2020a) proposed SemBERT, which directly connects multiple predicate-argument structures acquired by semantic role labeler (SRL) to get the joint representation.

The essence of SRL (Shi and Lin, 2019) lies in that every sentence possesses multiple predicate-specific structures which can represent different frames of events, while semantic roles express the abstract role that arguments of a predicate can take in the event. Besides, the events inside a sentence have interactions with each other that serve together to present the overall semantic knowledge. As shown

*Corresponding author. This work was supported in part by the Key Projects of National Natural Science Foundation of China under Grants U1836222 and 61733011.

in Figure 1, SRL parses every sentence with multiple predicate-specific structures which can serve as events inferring *who did what to whom, when and why*. Each event has an inner structure centered on the predicate to which several arguments are associated such as *Hoy*[ARG0], *the woman's age*[ARG1] and *Tuesday*[ARGM-TMP] connected to *confirmed*[V]. Meanwhile, the multiple events work together to give a comprehensive meaning of a sentence, like the events centered on *said*, *confirmed* and *left*. With regard to delving into the inner interactions between the events and effectively capturing multiple objects, we are motivated to build a graph to reveal the intrinsic structures between and inside the events.

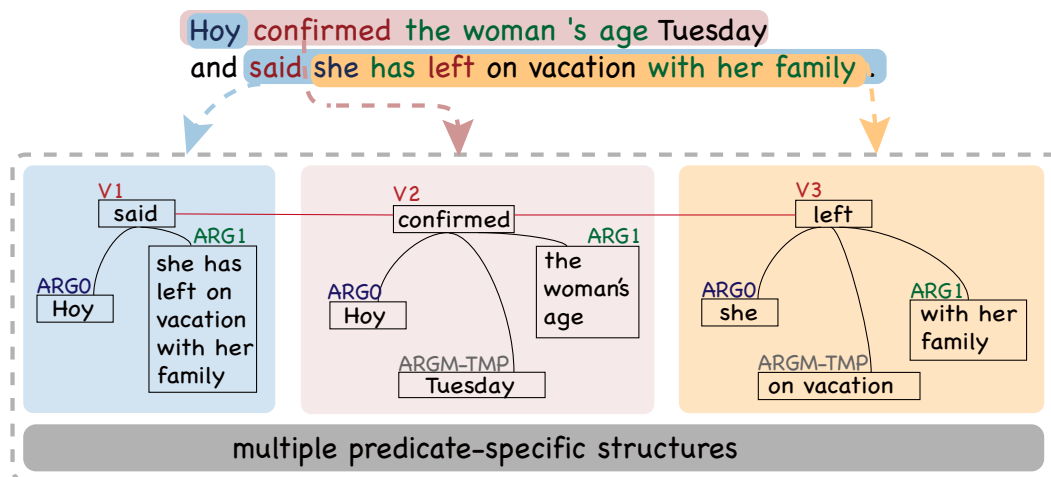


Figure 1: An example showing how SRL parses sentences and the intuition of constructing event-based graph.

Inspired by the above ideas, we propose EventBERT: an event-based semantic representation model which takes BERT as the backbone and refines with event-based structural semantics. Our EventBERT benefits simultaneously from rich event-based structures embodied in the graph and contextual semantics learned in the pre-trained BERT.

Our proposed model works in three steps: it first applies an off-the-shelf SRL toolkit to parse every sentence with semantic role labels; then it constructs event-based graphs and employs Graph Convolutional Networks (GCNs) (Schlichtkrull et al., 2018) to propagate and aggregate information from neighboring nodes on the graph; at last, it combines the contextualized representation acquired by BERT encoder together with the graph-level representation to obtain an event-based contextualized representation.

The key contributions of our work are summarized as follows:

- 1) We extract event-based semantic knowledge from SRL to enrich language representation.
- 2) We employ GCNs to construct sentence-level graphs which better reveal interactions inside and between the events in a sentence.

2 Related Work

2.1 Semantics in Language Representation

Recent studies show that current prominent pre-trained language models have already incorporated semantic information to some extent (Clark et al., 2019), yet such implicit semantic information is far from enough for comprehensive natural language understanding (Ettinger, 2020). Thus there emerges a research line that focuses on fusing semantic information into contextualized language representation. ERNIE2.0 (Sun et al., 2020) adopts three-stage masking in which entity-level masking helps to obtain a word representation containing richer semantic information. SemBERT (Zhang et al., 2020a) makes use of PropBank (Palmer et al., 2005) to fuse semantic role tags into language representation. FMSR (Guo et al., 2021) utilizes FrameNet (Baker et al., 1998) to extract multi-level semantic information within sentences. SS-MRC (Guo et al., 2020a) takes advantage of syntax and frame semantics in an attempt to carve out information from two complementary perspectives to obtain richer language representation.

Besides simply employing semantic knowledge, other recent works shift the focus to exploring deeper structural semantics. Guan et al. (2021) leverage frame semantics and graph neural networks to model sentences from both intra-sentence level and inter-sentence level. Wu et al. (2021) introduce SIFT to inject predicate-argument semantic dependencies into pre-trained language models via R-GCNs. Xie et al. (2022) introduce structured knowledge through multi-tasking to get a unified model, which inspires the potential of leveraging structural information. Unlike previous works that attempt to capture shallow semantic structures by semantic tags, our model digs deeper into semantics itself and aims to find the structured event-based information behind semantics, thus unveiling richer structural-semantic information inside the sentence.

2.2 Graph Modeling for Language Understanding

As natural language itself abounds with dependencies and intricate relations between different levels of language units, graph neural networks (GNNs), which model the units as nodes in the graph and learn the weight via the message passing between nodes of the graph (Scarselli et al., 2008; Kipf and Welling, 2016; Velickovic et al., 2017), stand out by explicitly and intuitively capturing the relations. Besides, a number of extensions to the original graph neural networks have been developed, the most notable of which include graph convolutional networks (GCNs) (Kipf and Welling, 2016), graph attention networks (GANs) (Velickovic et al., 2017) and the models from Li et al. (2015) and Pham et al. (2017) utilizing gating mechanisms to facilitate optimization.

In response to the outstanding performance of GCNs, several efforts have been made in recent years to improve performance on natural language understanding using GCNs, including GraphRel (Fu et al., 2019) which considers the interaction between named entities and relations via relation-weighted GCNs to better extract relations, NumNet (Ran et al., 2019) which utilizes a numerically-aware graph to perform numerical reasoning, DFGN (Qiu et al., 2019) which dynamically builds the entity graph by adding the edges with co-occurrence relations, HGN (Fang et al., 2019) which creates a hierarchical graph by constructing nodes on different levels of granularity and social information reasoning (Li and Goldwasser, 2019) which uses GCNs to capture the documents' social context.

Moreover, R-GCNs (Schlichtkrull et al., 2018) have shown effectiveness in relational graph modeling. For example, Entity-GCN (De Cao et al., 2019) employs R-GCNs to link mentions of candidate answers for multi-document question answering. DFGN (Qiu et al., 2019) dynamically builds the entity graph by adding the edges with co-occurrence relations and softly masking out irrelevant entities. DGM (Ouyang et al., 2021) constructs two discourse graphs and uses R-GCNs to fully capture interactions among the elements. Ma et al. (2022) employs R-GCNs to enhance reference dependencies for dialogue disentanglement. In contrast with previous works, our work proposes a sentence-level graph that is finely designed to mine the relationships between multiple elements in a sentence, extract rich structural semantics and facilitate information flow over the graph as well.

3 Model

Figure 2 gives an overview of our proposed EventBERT, which consists of two major components:

1. Context Encoder which acquires deep and contextualized representations for raw input sequences by following BERT architecture;
2. Event-based Encoder which obtains richer structural-semantic representation by modeling event-based intra-sentence graphs.

We omit the details of BERT which is widely used and ubiquitous and leave readers to resort to Devlin et al. (2019) for more information.

3.1 Context Encoder

The raw input sentence $X = \{x_1, \dots, x_n\}$ is a sequence of words in length n . It is first tokenized to a sequence of sub-words with [SEP] inserted at the end as the end marker and [CLS] inserted at the beginning to get a sentence-level representation: $X' = \{token_1, \dots, token_m\}$. Then we pass it

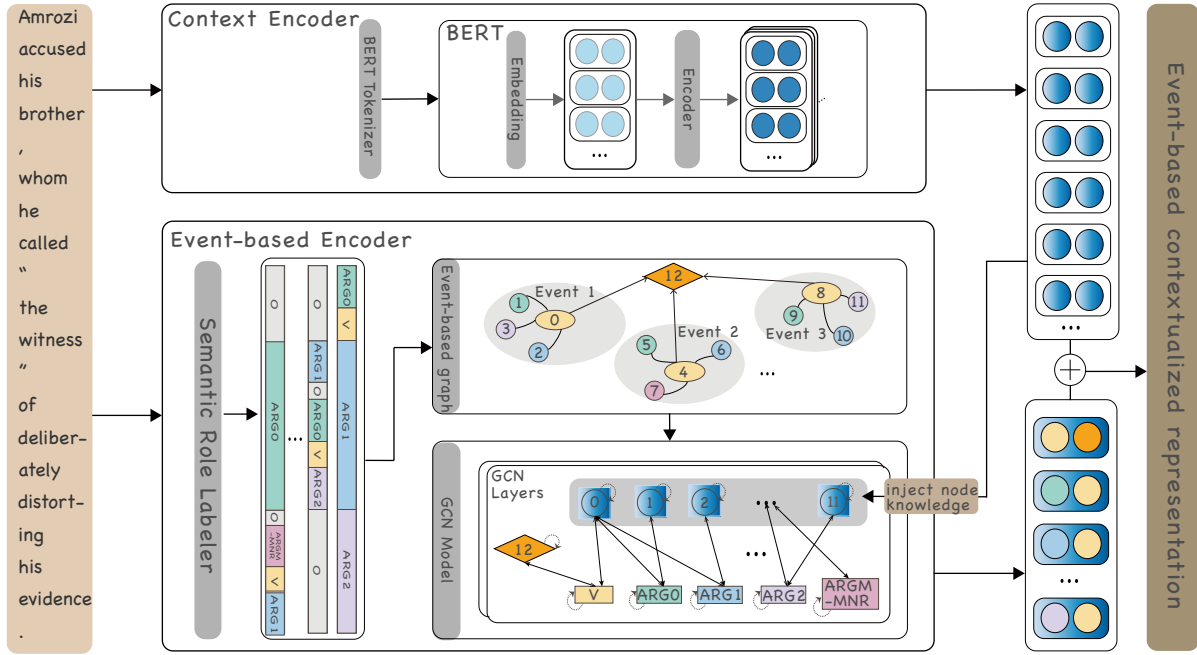


Figure 2: The overall structure of EventBERT.

through the embedding block and encoder block of BERT to produce a context-informed representation $C = \{c_1, \dots, c_m\} \in \mathbb{R}^{m \times d_{hs}}$ using the equation below:

$$C = BERT(X'), \quad (1)$$

where m denotes the length of sentence on sub-word level and d_{hs} stands for the dimension of hidden states.

3.2 Event-based Encoder

3.2.1 Semantic Role Labeler

The raw input sentence is simultaneously fed into Semantic Role Labeler (Shi and Lin, 2019) to fetch multiple predicate-specific structures tagged by PropBank semantic roles:

$$T = \{t_1, \dots, t_d\}, \quad (2)$$

where d is the number of semantic structures for one sentence. Notably, t_i can be represented under the format $\{tag_1^i, tag_2^i, \dots, tag_n^i\}$ and every tag span in t_i is recorded with its corresponding index in the context for further alignment.

3.2.2 Graph Construction

Figure 3 shows the process of graph construction: the predicates in the original input text are firstly extracted and an event subgraph is constructed with each predicate as the center; then a super event node (SEN) is applied to link all the predicates to collect the integral event information within the aggregated sentence; the Levi graph is finally constructed with reference to the method of Levi (1942), which is used to prepare the next stage of further computational operations on the graph.

For each sentence with the argument-predicate roles, we construct an event-based graph $G = (\mathcal{V}, \mathcal{E}, \mathcal{R})$ with span-level nodes $v_i \in \mathcal{V}$ and labeled edges $(v_i, r, v_j) \in \mathcal{E}$, where $r \in \mathcal{R}$ a relation type. Since every sentence has several semantic structures, here we take one structure as example and show the modeling method. Given $Seq_{tag} = \{tag_1, tag_2, \dots, tag_n\}$ a word-level tag sequence,

1. We first transform it to a span-level sequence $Seq'_{tag} = \{tag'_1, tag'_2, \dots, tag'_l\}$ by aggregating the same neighboring tags with $l \leq n$ representing the length of tags on span-level;

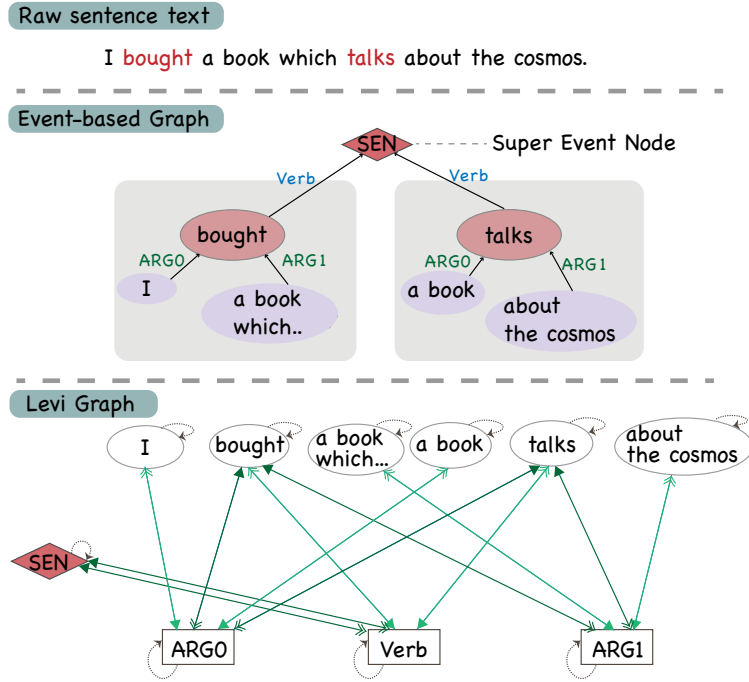


Figure 3: The process of graph construction: from raw sentence text to event-based graph and corresponding Levi graph.

2. Then, we add a Super Event Node ($v = SEN$) to seize global graph information;
3. After that, we add other nodes and edges to G based on the following process:
 - (a) we first find tag'_p which corresponds to predicate ($Verb$ in e'),
 - (b) we add a node $v = n_p$ and a directed edge $e = (n_p, Verb, SEN)$ with $r = Verb$,
 - (c) for the rest tags referring to arguments of the predicate, tag'_q for example, we add a node $v = n_q$ and a directed edge linking to the predicate $e = (n_q, tag'_q, n_p)$ with relation $r = tag_q$;
4. Finally, the corresponding Levi graph (Levi, 1942) is extended from G to $G_L = (\mathcal{V}_L, \mathcal{E}_L, \mathcal{R}_L)$. For nodes \mathcal{V}_L , we add the nodes representing relations to the original: $\mathcal{V}_L = \mathcal{V} \cup \mathcal{R}$. For edges \mathcal{E}_L , we transform each edge $e = (n_q, tag'_q, n_p)$ in G into two corresponding edges: $e_1 = (n_q, tag'_q)$ and $e_2 = (tag'_q, n_p)$ in G_L . For \mathcal{R}_L , we follow the setting of Ouyang et al. (2021) and refine it to five types: *default-in*, *default-out*, *reverse-in*, *reverse-out*, *self* according to the direction of edges towards the relation vertices, as is shown in Table 1.

Table 1: Relation types in our extended Levi graph

| \mathcal{R}_L in Levi graph | Illustration |
|-------------------------------|--|
| <i>default-in</i> | the propagation path pointing to the node as the end point |
| <i>default-out</i> | the propagation path pointing to the node as the starting point |
| <i>reverse-in</i> | the propagation path in the opposite direction of <i>default-in</i> |
| <i>reverse-out</i> | the propagation path in the opposite direction of <i>default-out</i> |
| <i>self</i> | the propagation paths pointing to the node itself |

3.2.3 Event-based Contextualized Representation

We adopt Relational Graph Convolutional Networks (R-GCNs) (Schlichtkrull et al., 2018) to implement explicit event graphs since traditional Graph Convolutional Networks (GCNs) cannot handle graphs con-

taining edge features with multiple relations. For predicate and argument nodes, we inject the corresponding span-level encoding results obtained from Context Encoder in Section 3.1. For relation nodes, we regard the relations as embeddings and use a lookup table to get the initial representation. Given that the initial representation of each node v_i is h_i^0 , the propagation process can be written as:

$$h_i^{(l+1)} = \text{ReLU} \left(\sum_{r \in \mathcal{R}_L} \sum_{v_j \in \mathcal{N}_r(v_i)} \frac{1}{c_{i,r}} w_r^{(l)} h_j^{(l)} \right), \quad (3)$$

where $h_i^{(l)} \in \mathbb{R}^{d^{(l)}}$ is the hidden state of node v_i in layer l with $d^{(l)}$ being the dimensionality of this layer's representations, $\mathcal{N}_r(v_i)$ denotes the set of neighbor indices of node v_i under the relation r , $c_{i,r}$ is a problem-specific normalization constant equal to $|\mathcal{N}_i^r|$, $w_r^{(l)}$ is the learnable parameters of layer l .

Since the importance of these relations cannot be treated the same, for example, the relation *Verb* is much more important than the relation *ARG2*, we introduce the gating mechanism (Marcheggiani and Titov, 2017). The basic idea is to compute a value between 0 and 1 for message passing control as is shown in Equation 4. Finally, the propagation process of R-GCNs under the gating mechanism is as follows:

$$g_j^{(l)} = \text{Sigmoid} \left(h_j^{(l)} W_{r,g}^{(l)} \right) \quad (4)$$

$$h_i^{(l+1)} = \text{ReLU} \left(\sum_{r \in \mathcal{R}_L} \sum_{v_j \in \mathcal{N}_r(v_i)} g_j^{(l)} \frac{1}{c_{i,r}} w_r^{(l)} h_j^{(l)} \right), \quad (5)$$

where $W_{r,g}^{(l)}$ is the learnable parameter under the l -th level relation type r .

With R-GCNs model, we obtain a graph-level semantic representation:

$$R = \{r_1, \dots, r_f\} \in \mathbb{R}^{f \times d_{hs}} \quad (6)$$

where f is the number of nodes in the graph and d_{hs} is the same dimension as the representation C in Equation 1 obtained from the context encoder.

At last, we concatenate R with the contextual sub-word-level representation C provided by Context Encoder and generate an event-based contextualized representation taking the mean value of both sub-word-level and graph-level information, which is then used as the new sequence representation for downstream tasks following the same way of Devlin et al. (2019).

4 Experiments

4.1 Setup

4.1.1 Datasets

We build EventBERT on the BERT backbone and fine-tune the model on GLUE (General Language Understanding Evaluation) benchmark (Wang et al., 2018) to evaluate the performance, which includes two single-sentence tasks CoLA (Warstadt et al., 2018), SST-2 (Socher et al., 2013)), three similarity and paraphrase tasks MRPC (Dolan and Brockett, 2005), STS-B (Cer et al., 2017), QQP (Chen et al., 2018), three inference tasks MNLI (Nangia et al., 2017), QNLI (Rajpurkar et al., 2016), RTE (Bentivogli et al., 2009). We exclude the controversial and problematic dataset WNLI (Levesque et al., 2012).

4.1.2 Evaluation Metrics

According to Wang et al. (2018), different datasets in GLUE correspond to different evaluation metrics, which include accuracy (acc), Matthew's correlation (mc) and Pearson correlation (pc). Among the eight datasets, STS-B is reported by Pearson correlation, CoLA is reported by Matthew's correlation, and other tasks are reported by accuracy.

| Model | CoLA (mc) | SST-2 (acc) | MNLI (acc) | QNLI (acc) | RTE (acc) | MRPC (acc) | QQP (acc) | STS-B (pc) | Avg - |
|----------------------------|--------------|----------------|---------------|---------------|--------------|---------------|--------------|---------------|------------|
| <i>Base-size</i> | | | | | | | | | |
| BERT _{BASE} | 58.4 | 92.8 | 83.2 | 88.6 | 68.5 | 86.0 | 86.5 | 87.8 | 81.5 |
| EventBERT _{BASE} | 59.6 | 93.3 | 83.9 | 91.8 | 69.7 | 89.7 | 89.8 | 88.9 | 83.3(↑1.8) |
| <i>Large-size</i> | | | | | | | | | |
| BERT _{LARGE} | 60.3 | 93.1 | 85.2 | 91.5 | 70.3 | 88.5 | 90.2 | 89.3 | 83.6 |
| EventBERT _{LARGE} | 63.1 | 94.0 | 85.3 | 92.6 | 71.4 | 89.5 | 90.6 | 89.5 | 84.5(↑0.9) |

Table 2: Comparisons between our models and baseline models on GLUE dev set. STS-B is reported by Pearson correlation, CoLA is reported by Matthew’s correlation, and other tasks are reported by accuracy.

4.1.3 Implementation Details

For the experiments, we use an initial learning rate in $\{1e-5, 2e-5, 3e-5\}$ with warm-up rate of 0.1 and L2 weight decay of 0.01. The batch size is selected in $\{16, 32\}$. The maximum number of epochs is set in $[2, 5]$ depending on tasks. Texts are tokenized with maximum length of 256 for the tasks. We use 2 layers of R-GCNs in our model.

4.2 Results

Table 2 presents the results on the GLUE benchmark, which show that EventBERT achieves consistent gains over all the subtasks under both base and large models.

The results indicate that our model performs better on longer sentences as shown in Section 5.3. Furthermore, our analysis shows that EventBERT can effectively benefit from the fine-grained graph-like event-based structures, as illustrated in case studies in Section 5.4. The results also disclose that modeling intrinsic structures between and inside events is crucial for language understanding.

In addition, the experimental results show that EventBERT has a significant performance gain on small datasets such as CoLA and MRPC, which indicates that semantic information involving event modeling is more advantageous and competitive in smaller datasets. In practice or industry, large-scale annotated data is rare and scarce due to the high cost and required expensive human resources, so language models that dominate in small-scale datasets are more valuable and important for most NLP tasks.

5 Analysis

5.1 Ablation Study

We conduct the ablation study to investigate the effects of the gating mechanism and the addition of global nodes in the event-based encoder module. Results in Table 3 show that both the gating mechanism and global nodes are non-trivial.

5.2 Methods of Aggregation

During the period of concatenating and aggregating the graph level semantic representation R and the contextual representation C , we further analyze the influence of different methods of aggregation such as max-pooling and mean-pooling by comparing the models with the same hyper-parameters on three datasets CoLA, MRPC and RTE respectively. Results in Table 3 demonstrate that employing mean-pooling presents better performance.

| Model | CoLA (mc) | MRPC (acc) | RTE (acc) |
|----------------------------|--------------|---------------|--------------|
| <i>Ablation study</i> | | | |
| EventBERT _{base} | 59.6 | 89.7 | 69.7 |
| w/o gating | 58.6 | 86.8 | 69.0 |
| w/o global node | 58.4 | 87.0 | 67.9 |
| <i>Aggregation methods</i> | | | |
| BERT _{base} | 58.4 | 86.0 | 68.5 |
| w/ max-pooling | 59.1 | 86.8 | 68.2 |
| w/ mean-pooling | 59.6 | 89.7 | 69.7 |

Table 3: Ablation study and comparison of aggregation methods on three datasets.

5.3 Effectiveness of semantic structures

In order to dig deeper into the rationale behind the effectiveness of the model, we select two datasets QNLI and MRPC, representing large-scale and small-scale datasets respectively. We statistically calculate the accuracy of the corresponding models on different word-level sequence length intervals for EventBERT and baseline. Figure 4 shows that our model outperforms the baseline especially when the sequence is relatively long and our model performs better on longer sentences compared with shorter ones, which implies that modeling intrinsic semantic structures is potential to guide the model to learn richer structural semantics more than contextualized information. Thus, the analysis of word sequence lengths shows that EventBERT performs better on data with longer sequence lengths, which indicates that event-level modeling is promising and competitive for understanding long texts. Under many practical situations where available data are long texts, the idea of extracting event-level structural-semantic information is promising in many NLP tasks.

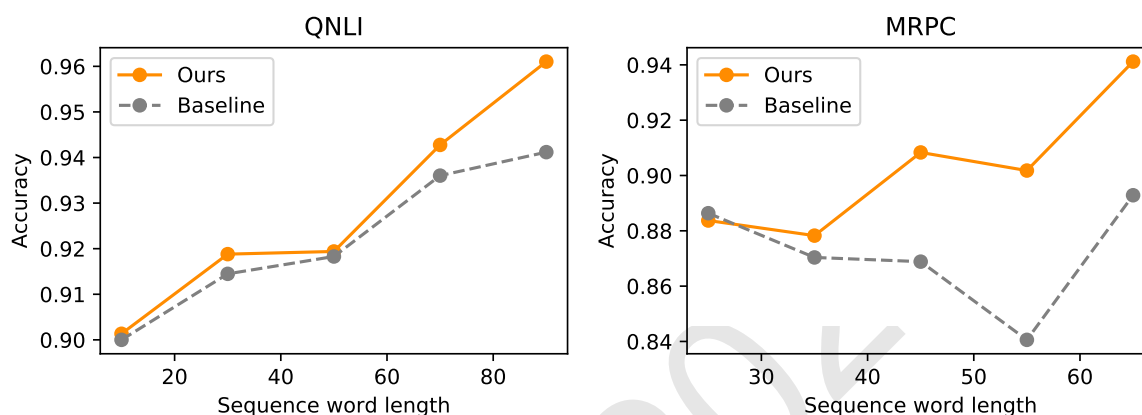


Figure 4: Accuracy of different sequence word lengths on QNLI and MRPC.

5.4 Interpretability: Case Study

We select three cases in Classification, Sentence Similarity and Language Inference from SST-2, MRPC and QNLI respectively which are shown in Figure 5, aiming to further explore the mechanism. It can be seen that our model can perceive explicit structural meaning to better understand the language. We will analyze each of the three cases in detail so as to analyze the advantages of EventBERT more intuitively.

5.4.1 Classification

In the case from SST-2, our model succeeds in capturing and understanding the event *Friel and william's exceptional performances[ARG0] anchored[V] the film's power[ARG1]*, whereas the baseline does not manage to capture this meaning, thus leading to the failure.

5.4.2 Sentence Similarity

The case from MRPC demonstrates that our model grabs the distinct semantic structures centered on *is* and *has* and thus gives the right answer *not equivalent*. The event centered on the predicate *donate* belongs to the same structure, which contains the arguments *ARG0*, *ARG1* and *ARGM-TMP* having the same contents (i.e., *the woman donated blood*). Nevertheless, the remaining events which center on the predicate *is* and the predicate *has* in the sentence pair are semantically different as one structure includes the arguments *ARG1* and *ARG2* while the other contains only *ARG0* and *ARG1*.

In Sentence Similarity tasks, two sentences in a sentence pair are likely to have one or several events in common, such as the event centered *donate* in this case. However, a subtle difference in a key element in the semantic structure of the sentence may also lead to a very different semantics of the whole sentence, such as the events centered on *is* and *has*. Our proposed model EventBERT precisely appreciates the

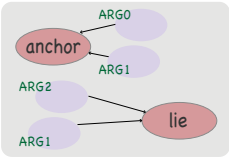
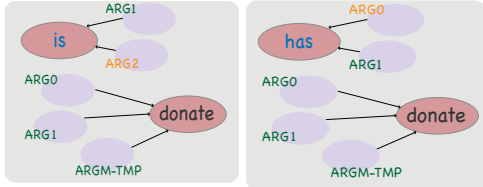
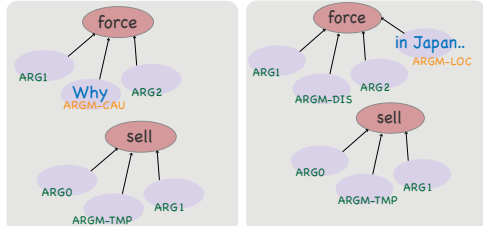
| Task | Example | Graph | Results |
|---------------------|--|--|---|
| Classification | anchored by friel and williams 's exceptional performances , the film 's power lies in its complexity . |  | EventBERT: positive ✓ Baseline: negative ✗ |
| Sentence Similarity | A: The Calgary woman , who is in her twenties, donated blood on Aug. 7 . B: The woman -- who has no symptoms of illness -- donated blood Aug. 7 |  | EventBERT: not equivalent ✓ Baseline: equivalent ✗ |
| Language Inference | A: Why was ABC forced to sell its interests in international networks in the 70s? B: As a result, ABC was forced to sell all of its interests in international networks, mainly in Japan and Latin America, in the 1970s. |  | EventBERT: not entailment ✓ Baseline: entailment ✗ |

Figure 5: Examples selected from the dev set of SST-2, MRPC and QNLI where baseline fails but our model succeeds.

value of abstracting structural semantics, benefiting from capturing event-based semantic knowledge to perceive the differences between sentences and thus make more accurate judgments.

5.4.3 Language Inference

Referring to the case from QNLI, as can be seen from Figure 5, the question and paragraph texts are broadly similar in terms of *sell*-centered structure, both containing the arguments labeled *ARG0*, *ARG1*, and *ARGM-TMP*. However, by means of graph modeling, it can be clearly and explicitly observed that the structures centered on *force* are distinct, with the structure in the interrogative sentence containing the argument *ARGM-CAU* and the corresponding structure in the paragraph texts containing the argument *ARGM-LOC* instead. It is worth noting that one of the most crucial steps in determining whether a paragraph entails the correct answer to a question is whether the corresponding semantic structure in paragraph texts has the span labeled with the semantic role referring to the interrogative in the question. For example, in this case, the interrogative *Why* is exactly the *ARGM-CAU* of the predicate *force*; whereas the structure centered on *force* in the paragraph lacks the corresponding argument content and is replaced by *ARGM-LOC* instead. Therefore, it can be easily inferred that the paragraph focuses on the location (i.e., *in Japan and Latin America*) while the question concentrates on the cause (i.e., *Why*), which exactly reflects that there is no answer span for the interrogative of the question.

It is known that interrogative in the question and corresponding answer span should belong to the same semantic role. EventBERT takes full advantage of extracting abstracted semantics based on predicates, thus conducting language inference tasks more efficiently.

5.5 Error Analysis

We select bad cases of the baseline model and further investigate the ones of which our EventBERT also fails to predict the correct answers. We study two cases respectively from MRPC and QNLI as is shown in Table 4. The first error is caused by EventBERT’s identification of the argument *in a written statement* of the predicate *said* in the first sentence, which is not entailed in the second sentence. However, the lack of this argument does not affect the main semantic information. The second error is due to argument reference confusion for the special predicate *is*. For instance, the interrogative *What* is labeled as *ARG2* whereas the correct answer *Hypersensitivity* is labeled as *ARG1*. From the above error cases, it may

suggest that our model needs to have a more accurate perception of semantic relationships, which is left for future studies.

| Example | EventBERT | Golden Answer |
|--|----------------|---------------|
| This decision is clearly incorrect ,” FTC Chairman Timothy Muris <i>said in a written statement</i> . The decision is ” clearly incorrect ,” FTC Chairman Tim Muris <i>said</i> . | Not equivalent | Equivalent |
| <i>What is</i> the name for a response of the immune system that damages the body’s native tissues? <i>Hypersensitivity is</i> an immune response that damages the body’s own tissues. | Not entailment | Entailment |

Table 4: Errors in predictions for cases in MRPC and QNLI dev set. The words in magenta indicate the key predicate. The words in blue indicate the key arguments referred to the predicate.

6 Conclusion

In this work, we propose EventBERT, an event-based semantic representation model that builds on BERT architecture and incorporates event-based structural semantics in terms of graph network modeling for fine-grained language representation. Experiments on a wide range of NLU tasks show the effectiveness of our model by consistently surpassing the baseline. While most existing works focus on fusing accurate semantic signals to enhance semantic information, we open up a novel perspective to model intrinsic structural semantics for deeper comprehension and inference in an intuitive and explicit way.

References

- Jiangang Bai, Yujing Wang, Yiren Chen, Yaming Yang, Jing Bai, Jing Yu, and Yunhai Tong. 2021. Syntax-bert: Improving pre-trained transformers with syntax trees. *arXiv preprint arXiv:2103.04350*.
- Collin F. Baker, Charles J. Fillmore, and John B. Lowe. 1998. The Berkeley FrameNet project. In *36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Volume 1*, pages 86–90, Montreal, Quebec, Canada, August. Association for Computational Linguistics.
- Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. 2009. The fifth pascal recognizing textual entailment challenge. In *ACL-PASCAL*.
- Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. *arXiv preprint arXiv:1708.00055*.
- Zihan Chen, Hongbo Zhang, Xiaoji Zhang, and Leqi Zhao. 2018. Quora question pairs.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does BERT look at? an analysis of BERT’s attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy, August. Association for Computational Linguistics.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. 2019. Question answering by reasoning across documents with graph convolutional networks. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2306–2317, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- William B Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *IWP2005*.
- Allyson Ettinger. 2020. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics*, 8:34–48.

- Yuwei Fang, Siqi Sun, Zhe Gan, Rohit Pillai, Shuohang Wang, and Jingjing Liu. 2019. Hierarchical graph network for multi-hop question answering. *arXiv preprint arXiv:1911.03631*.
- Tsu-Jui Fu, Peng-Hsuan Li, and Wei-Yun Ma. 2019. Graphrel: Modeling text as relational graphs for joint entity and relation extraction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1409–1418.
- Dirk Geeraerts and Hubert Cuyckens. 2007. Introducing cognitive linguistics. In *The Oxford handbook of cognitive linguistics*.
- Yong Guan, Shaoru Guo, Ru Li, Xiaoli Li, and Hongye Tan. 2021. Frame semantic-enhanced sentence modeling for sentence-level extractive text summarization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4045–4052, Online and Punta Cana, Dominican Republic, November. Association for Computational Linguistics.
- Shaoru Guo, Yong Guan, Ru Li, Xiaoli Li, and Hongye Tan. 2020a. Incorporating syntax and frame semantics in neural network for machine reading comprehension. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2635–2641, Barcelona, Spain (Online), December. International Committee on Computational Linguistics.
- Shaoru Guo, Ru Li, Hongye Tan, Xiaoli Li, Yong Guan, Hongyan Zhao, and Yueping Zhang. 2020b. A frame-based sentence representation for machine reading comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 891–896, Online, July. Association for Computational Linguistics.
- Shaoru Guo, Yong Guan, Ru Li, Xiaoli Li, and Hongye Tan. 2021. Frame-based multi-level semantics representation for text matching. *Knowledge-Based Systems*, 232:107454.
- Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*.
- Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In *Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning*.
- Friedrich Wilhelm Levi. 1942. *Finite geometrical systems: six public lectures delivered in February, 1940, at the University of Calcutta*. University of Calcutta.
- Chang Li and Dan Goldwasser. 2019. Encoding social information with graph convolutional networks for political perspective detection in news media. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2594–2604.
- Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. 2015. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493*.
- Xinbei Ma, Zhuosheng Zhang, and Hai Zhao. 2022. Structural characterization for dialogue disentanglement. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 285–297, Dublin, Ireland, May. Association for Computational Linguistics.
- Diego Marcheggiani and Ivan Titov. 2017. Encoding sentences with graph convolutional networks for semantic role labeling. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1506–1515, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Nikita Nangia, Adina Williams, Angeliki Lazaridou, and Samuel R Bowman. 2017. The repeval 2017 shared task: Multi-genre natural language inference with sentence representations. In *RepEval*.
- Siru Ouyang, Zhuosheng Zhang, and Hai Zhao. 2021. Dialogue graph modeling for conversational machine reading. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3158–3169, Online, August. Association for Computational Linguistics.
- Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. The Proposition Bank: An annotated corpus of semantic roles. *Computational Linguistics*, 31(1):71–106.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana, June. Association for Computational Linguistics.

- Trang Pham, Truyen Tran, Dinh Phung, and Svetha Venkatesh. 2017. Column networks for collective classification. In *Thirty-first AAAI conference on artificial intelligence*.
- Lin Qiu, Yunxuan Xiao, Yanru Qu, Hao Zhou, Lei Li, Weinan Zhang, and Yong Yu. 2019. Dynamically fused graph network for multi-hop reasoning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6140–6150, Florence, Italy, July. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *EMNLP*.
- Qiu Ran, Yankai Lin, Peng Li, Jie Zhou, and Zhiyuan Liu. 2019. Numnet: Machine reading comprehension with numerical reasoning. *arXiv preprint arXiv:1910.06701*.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8:842–866.
- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2008. The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80.
- Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. In *European semantic web conference*, pages 593–607. Springer.
- Peng Shi and Jimmy J. Lin. 2019. Simple bert models for relation extraction and semantic role labeling. *ArXiv*, abs/1904.05255.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP*.
- Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. 2020. Ernie 2.0: A continual pre-training framework for language understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8968–8975.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. *stat*, 1050:20.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium, November. Association for Computational Linguistics.
- Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. 2018. Neural network acceptability judgments. *arXiv preprint arXiv:1805.12471*.
- Zhaofeng Wu, Hao Peng, and Noah A. Smith. 2021. Infusing finetuning with semantic dependencies. *Transactions of the Association for Computational Linguistics*, 9:226–242.
- Tianbao Xie, Chen Henry Wu, Peng Shi, Ruiqi Zhong, Torsten Scholak, Michihiro Yasunaga, Chien-Sheng Wu, Ming Zhong, Pengcheng Yin, Sida I Wang, et al. 2022. Unifiedskg: Unifying and multi-tasking structured knowledge grounding with text-to-text language models. *arXiv preprint arXiv:2201.05966*.
- Zenan Xu, Daya Guo, Duyu Tang, Qinliang Su, Linjun Shou, Ming Gong, Wanjun Zhong, Xiaojun Quan, Daxin Jiang, and Nan Duan. 2021. Syntax-enhanced pre-trained model. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5412–5422, Online, August. Association for Computational Linguistics.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32.
- Zhuosheng Zhang, Yuwei Wu, Hai Zhao, Zuchao Li, Shuailiang Zhang, Xi Zhou, and Xiang Zhou. 2020a. Semantics-aware bert for language understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 9628–9635.
- Zhuosheng Zhang, Yuwei Wu, Junru Zhou, Sufeng Duan, Hai Zhao, and Rui Wang. 2020b. Sg-net: Syntax-guided machine reading comprehension. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 9636–9643.