

# GADePo: Graph-Assisted Declarative Pooling Transformers for Document-Level Relation Extraction

Anonymous ACL submission

## Abstract

Document-level relation extraction typically relies on text-based encoders and hand-coded pooling heuristics to aggregate information learned by the encoder. In this paper, we leverage the intrinsic graph processing capabilities of the Transformer model and propose replacing hand-coded pooling methods with new tokens in the input, which are designed to aggregate information via explicit graph relations in the computation of attention weights. We introduce a joint text-graph Transformer model and a graph-assisted declarative pooling (GADePo) specification of the input, which provides explicit and high-level instructions for information aggregation. GADePo allows the pooling process to be guided by domain-specific knowledge or desired outcomes but still learned by the Transformer, leading to more flexible and customisable pooling strategies. We evaluate our method across diverse datasets and models and show that our approach yields promising results that are consistently better than those achieved by the hand-coded pooling functions.

## 1 Introduction

Document-level relation extraction is an important task in natural language processing, which involves identifying and categorising meaningful relationships between entities within a document, as exemplified in Figure 1. This task is foundational to many applications, including knowledge base population and completion (Banko et al., 2007; Ji et al., 2020), information retrieval and extraction (Manning et al., 2008; Theodoropoulos et al., 2021), question answering (Chen et al., 2017; Feng et al., 2022) and sentiment analysis (Pang and Lee, 2008), to name a few.

Standard methods that approach this challenge generally employ pretrained text-based encoders (Devlin et al., 2019; Beltagy et al., 2019; Zhuang et al., 2021; Cui et al., 2021), which are responsible for capturing the nuances of information con-

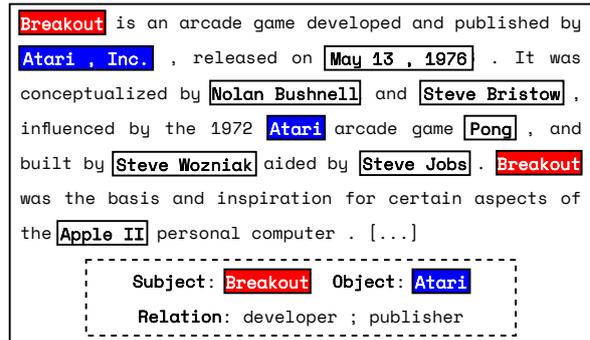


Figure 1: Document from the Re-DocRED (Tan et al., 2022b) dataset involving multiple entities and labels. Subject entity **Breakout** (red) and object entity **Atari** (blue) express relations "developer" and "publisher". Other entities are indicated as **Mention** (white).

tained in the entity mentions and their contextual surroundings. Previous successful methods often then use hand-coded pooling heuristics to aggregate the information learned by the encoder, with some aimed at creating entity representations, while others directly exploiting the pattern of attention weights to capture context aware relations between entity mentions (Zhou et al., 2021; Xiao et al., 2022; Tan et al., 2022a; Ma et al., 2023). These pooling heuristics can be very effective at leveraging the information in a pretrained encoder. However, as shown in Conneau et al. (2017); Jia et al. (2019); Reimers and Gurevych (2019); Choi et al. (2021), the selection of an appropriate pooling function can be model-dependent, task-specific, resource-intensive and time-consuming to determine, thereby limiting flexibility.

In this paper, we address these issues with a new approach where we leverage the intrinsic graph processing capabilities of the Transformer model (Vaswani et al., 2017), leveraging insights from the work of Mohammadshahi and Henderson (2020); Henderson (2020); Mohammadshahi and Henderson (2021); Henderson et al. (2023). They argue that attention weights and graph relations are functionally equivalent and show how to incorporate

068 structural dependencies between input elements  
069 by simply adding relation features to the attention  
070 functions. Transformers easily learn to integrate  
071 these relation features into their pretrained atten-  
072 tion functions, resulting in very successful graph-  
073 conditioned models (Mohammadshahi and Hen-  
074 derson, 2021; Miculicich and Henderson, 2022;  
075 Mohammadshahi and Henderson, 2023). Given  
076 this effective method for integrating explicit graphs  
077 with pretrained attention functions, we propose to  
078 use the attention function itself for aggregation. We  
079 replace the rigid pooling methods with new tokens  
080 which act as aggregation nodes, plus explicit graph  
081 relations which steer the aggregation.

082 We introduce a joint text-graph Transformer  
083 model and a graph-assisted declarative pooling  
084 (GADePo) method<sup>1</sup> that leverages these special  
085 tokens and graph relations, to provide an explicit  
086 high-level declarative specification for the infor-  
087 mation aggregation process. By integrating these  
088 graphs in the attention functions of a pretrained  
089 model, GADePo exploits the pretrained embed-  
090 dings and attention patterns but still has the flex-  
091 ibility of being trained on data. This enables the  
092 pooling to be guided by domain-specific knowl-  
093 edge or desired outcomes but still learned by the  
094 Transformer, opening up a more customisable but  
095 still data-driven relation extraction process.

096 We evaluate our method across diverse datasets  
097 and models commonly employed in document-  
098 level relation extraction tasks, and show that our  
099 approach yields promising results that are consis-  
100 tently better than those achieved by the hand-coded  
101 pooling functions.

102 **Contributions** We propose a new method for  
103 exploiting pretrained Transformer models which  
104 replaces hand-coded aggregation functions with ex-  
105 plicit graph relations and aggregation nodes. We  
106 introduce a novel form of joint text-graph Trans-  
107 former model. We evaluate our approach across  
108 various datasets and models, showing that it yields  
109 promising results that are consistently better than  
110 those achieved by hand-coded pooling functions.

## 111 2 Related Work

112 In recent studies, the scope of relation extraction  
113 has been expanded to include not only individ-  
114 ual sentences but entire documents. This exten-  
115 sion, known as document-level relation extraction,

116 presents a more realistic and challenging scenario  
117 as it seeks to extract relations both within sentences  
118 and across multiple sentences (Yao et al., 2019).  
119 Transformer-based (Vaswani et al., 2017) models  
120 have shown great potential in addressing this task.

121 Wang et al. (2019) and Tang et al. (2020) show  
122 that the BiLSTM-based (Hochreiter and Schmidhu-  
123 ber, 1997) baselines lack the capacity to model  
124 complex interactions between multiple entities.  
125 They propose a more robust approach, which con-  
126 sists of using the pretrained BERT (Devlin et al.,  
127 2019) model and a two-step prediction process, i.e.,  
128 first identifying if a link between two entities exists,  
129 followed by predicting the specific relation type.

130 GAIN (Zeng et al., 2020) leverages BERT as a  
131 text encoder and GCNs (Kipf and Welling, 2017)  
132 to process two types of graphs, one at mention level  
133 and another at entity level, showing notable perfor-  
134 mance in inter-sentence and inferential scenarios.

135 Mohammadshahi and Henderson (2020, 2021)  
136 propose the G2GT model and show how to lever-  
137 age the intrinsic graph processing capabilities of  
138 the Transformer model by incorporating structural  
139 dependencies between input elements as features  
140 input to the self-attention weight computations.

141 SSAN (Xu et al., 2021) leverages this idea and  
142 considers the structure of entities. It employs a  
143 transformation module that creates attentive biases  
144 from this structure to regulate the attention flow  
145 during the encoding phase.

146 DocuNet (Zhang et al., 2021) reformulates the  
147 task as a semantic segmentation problem. It em-  
148 ploys a U-shaped segmentation module and an en-  
149 coder module to capture global interdependencies  
150 and contextual information of entities, respectively.

151 PL-Marker (Ye et al., 2022) introduces a method  
152 that takes into account the interplay between spans  
153 via a neighbourhood-oriented and subject-oriented  
154 packing approach, highlighting the importance of  
155 capturing the interrelation among span pairs in re-  
156 lation extraction tasks.

157 SAIS (Xiao et al., 2022) explicitly models key  
158 information sources such as relevant contexts and  
159 entity types. It improves extraction quality and  
160 interpretability, while also boosting performance  
161 through evidence-based data augmentation and en-  
162 semble inference.

163 KD-DocRE (Tan et al., 2022a) proposes a semi-  
164 supervised framework with three key components.  
165 Firstly, an axial attention module enhances per-  
166 formance in handling two-hop relations by captur-  
167 ing the interdependence of entity pairs. Secondly,

<sup>1</sup>Code will be made available upon publication.

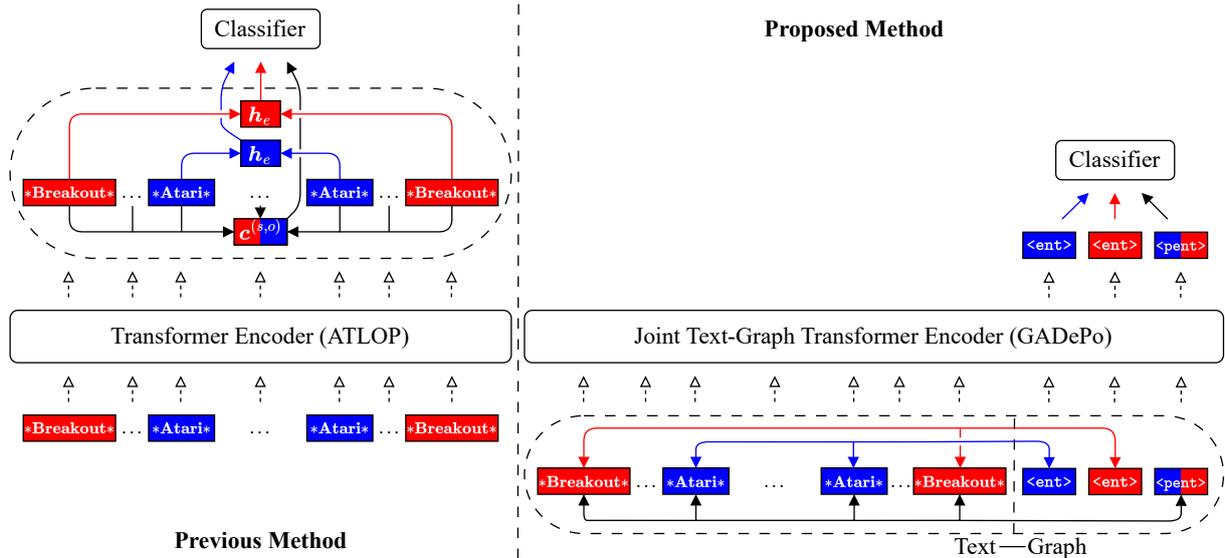


Figure 2: Comparison between the previous method ATLOP (left) and the proposed method GADePo (right), illustrating the document in Figure 1 containing two entities (red and blue), each with two mentions. In ATLOP, the mentions’ encoder outputs are aggregated into entity representations  $h_e$ , and the encoder’s attention weights are used to identify which outputs to aggregate for entity-pair representations  $c^{(s,o)}$ . In GADePo, the textual input is extended to include the graph special tokens  $\langle \text{ent} \rangle$  for entity representations and  $\langle \text{pent} \rangle$  for entity-pair representations, and explicit directional graph relations specify their associated mentions. A joint text-graph Transformer model is then used to encode this declarative pooling specification and compute the relevant aggregations.

an adaptive focal loss solution addresses the class imbalance issue. Lastly, the framework employs knowledge distillation to improve robustness and overall effectiveness by bridging the gap between human-annotated and distantly supervised data.

DREEAM (Ma et al., 2023) is a method designed to enhance document-level relation extraction by addressing memory efficiency and annotation limitations in evidence retrieval. It employs evidence as a supervisory signal to guide attention and introduces a self-training strategy to learn evidence retrieval without requiring evidence annotations.

SAIS (Xiao et al., 2022), KD-DocRE (Tan et al., 2022a), and DREEAM (Ma et al., 2023) have been built upon the foundations of ATLOP (Zhou et al., 2021). ATLOP introduces two innovative techniques, adaptive thresholding, and localised context pooling, to address challenges in multi-label and multi-entity problems. Adaptive thresholding employs a learnable entities-dependent threshold, replacing the global threshold used in previous approaches for multi-label classification (Peng et al., 2017; Christopoulou et al., 2019; Nan et al., 2020; Wang et al., 2020). Localised context pooling leverages the attention patterns of a pretrained language model to identify and extract relevant context crucial for determining the relation between entities, using specific hand-coded pooling functions.

### 3 Background

The foundational work of ATLOP (Zhou et al., 2021) has been the basis of many State-of-the-Art (SotA) models (Xiao et al., 2022; Tan et al., 2022a; Ma et al., 2023). Given the problems with hand-coded pooling functions, discussed in Section 1, we aim to provide a new baseline that can serve as the foundation for future SotA models. For this reason, we evaluate our proposed models by comparing them to this established baseline. Our goal is to demonstrate that our method not only achieves results comparable to or better than ATLOP, but also offers a novel approach which addresses its limitations. To provide a better understanding of ATLOP and its components, we present a detailed breakdown in the left portion of Figure 2, which we elaborate on in this section.

#### 3.1 Problem Formulation

The document-level relation extraction task involves analysing a document  $D$  that contains a set of entities  $\mathcal{E}_D = \{e_i\}_{i=1}^{|\mathcal{E}_D|}$ . The main objective is to determine the presence or absence of various relation types between all entity pairs  $(e_s, e_o)_{s,o \in \mathcal{E}_D, s \neq o}$ , where the subject and object entities are denoted as  $e_s$  and  $e_o$ , respectively. A key aspect to consider is that an entity can appear multiple times in the document, resulting in a cluster of

multiple mentions  $\mathcal{M}_e = \{m_i\}_{i=1}^{|\mathcal{M}_e|}$  for each entity  $e$ . The set of relations is defined as  $\mathcal{R} \cup \emptyset$ , where  $\emptyset$  represents the absence of a relation, often referred to as "no-relation". Given the clusters of mentions  $\mathcal{M}_{e_s}$  and  $\mathcal{M}_{e_o}$ , the task consists of a multi-label classification problem where there can be multiple relations between entities  $e_s$  and  $e_o$ .

### 3.2 Previous Method: ATLOP

**Text Encoding** A special token  $*$  is added at the start and end of every mention. Tokens  $\mathcal{T}_D = \{t_i\}_{i=1}^{|\mathcal{T}_D|}$  are encoded via a Pretrained Language Model (PLM) as follows:

$$\mathbf{H}, \mathbf{A} = PLM(\mathcal{T}_D), \quad (1)$$

where  $\mathbf{H} \in \mathbb{R}^{|\mathcal{T}_D| \times d}$  and  $\mathbf{A} \in \mathbb{R}^{|\mathcal{T}_D| \times |\mathcal{T}_D|}$  represent the token embeddings and the average attention weights of all attention heads, respectively, extracted from the last layer of the PLM.

**Entity Embedding (EE)** For each individual entity  $e$  with mentions  $\mathcal{M}_e = \{m_i\}_{i=1}^{|\mathcal{M}_e|}$ , an entity embedding  $\mathbf{h}_e \in \mathbb{R}^d$  is computed as follows:

$$\mathbf{h}_e = \log \sum_{i=1}^{|\mathcal{M}_e|} \exp(\mathbf{H}_{m_i}), \quad (2)$$

where  $\mathbf{H}_{m_i} \in \mathbb{R}^d$  is the embedding of the special token  $*$  at the starting position of mention  $m_i$ . The choice of the *logsumexp* pooling function is based on the research conducted by Jia et al. (2019). Their study offers empirical evidence that supports the use of this pooling function over others, as it facilitates accumulating weak signals from individual mentions, thanks to its smoother characteristics.

**Localised Context Embedding (LCE)** ATLOP introduces the concept of localised context embedding to accommodate the variations in relevant mentions and context for different entity pairs  $(e_s, e_o)$ . Since the attention mechanism in the PLM captures the importance of each token within the context, it can be used to determine the context relevant for both entities. The importance of each token can be computed from the cross-token dependencies matrix  $\mathbf{A}$  obtained in Equation 1. When evaluating entity  $e_s$ , the importance of individual tokens is determined by examining the cross-token dependencies across all mentions associated with  $e_s$ , denoted as  $\mathcal{M}_{e_s}$ . Initially, ATLOP collects and averages the attention  $\mathbf{A}_{m_i} \in \mathbb{R}^{|\mathcal{T}_D|}$  at the special token  $*$  preceding each mention  $m_i \in \mathcal{M}_{e_s}$ . This

process results in  $\mathbf{a}_s \in \mathbb{R}^{|\mathcal{T}_D|}$ , which represents the importance of each token concerning entity  $e_s$  (and analogously  $\mathbf{a}_o$  for  $e_o$ ). Subsequently, the importance of each token for a given entity pair  $(e_s, e_o)$ , denoted as  $\mathbf{q}^{(s,o)} \in \mathbb{R}^{|\mathcal{T}_D|}$ , is computed using  $\mathbf{a}_s$  and  $\mathbf{a}_o$  as follows:

$$\mathbf{q}^{(s,o)} = \frac{\mathbf{a}_s \circ \mathbf{a}_o}{\mathbf{a}_s^\top \mathbf{a}_o}, \quad (3)$$

where  $\circ$  represents the Hadamard product. Consequently,  $\mathbf{q}^{(s,o)}$  represents a distribution that indicates the importance of each token for both tokens in  $(e_s, e_o)$ . Finally, the localised context embedding is computed as follows:

$$\mathbf{c}^{(s,o)} = \mathbf{H}^\top \mathbf{q}^{(s,o)}, \quad (4)$$

So  $\mathbf{c}^{(s,o)} \in \mathbb{R}^d$  corresponds to a weighted average over all token embeddings that are important for both  $e_s$  and  $e_o$ .

**Relation Classification and Loss Function** The representations  $\mathbf{h}_{e_s}$ ,  $\mathbf{h}_{e_o}$  and  $\mathbf{c}^{(s,o)}$  are input to a relation classifier, and the full model is fine-tuned to predict the relation labels for  $(e_s, e_o)$ . The relation classifier and its loss function are detailed in Appendix Subsection A.1.

## 4 Proposed Method: GADePo

We propose to avoid the reliance on the EE (i.e.,  $\mathbf{h}_e$ ) and LCE (i.e.,  $\mathbf{c}^{(s,o)}$ ) heuristic aggregation functions by leveraging Transformers' attention functions to do aggregation. Given the observation of Henderson (2020); Mohammadshahi and Henderson (2020, 2021); Henderson et al. (2023) that attention weights and graph relations are functionally equivalent, we introduce the inductive biases of EE and LCE directly into the model's input as graph relations.

Our proposed graph-assisted declarative pooling (GADePo) method replaces the hand-coded aggregation functions EE and LCE with a declarative graph specification. By using the intrinsic graph processing capabilities of the Transformer model, the specified graph serves as an explicit high-level directive for the information aggregation process of the Transformer. By inputting the graph relations to the Transformer's self-attention layers, GADePo enables the aggregation to be steered by domain-specific knowledge or desired outcomes, while still allowing it to be learned by the Transformer, opening up the possibility for a more tailored and customised yet data-driven relation extraction.

Our GADePo model is illustrated in the right portion of Figure 2. We address both EE and LCE with the introduction of two special tokens, <ent> (i.e., entity) and <pent> (i.e., pair entity), and two explicit graph relations of types <ent>  $\longleftrightarrow$  \* and <pent>  $\longleftrightarrow$  \* in both directions, where \* represents the special token at the starting position of a specific mention. The set of relations is specified as  $c_{ij} \in \mathcal{C}$  which each identify the relation label from  $i$  to  $j$ . Each of these relation labels is associated with an embedding vector of dimension  $d$ , as are the special token inputs <ent> and <pent>. These two special tokens are added to the PLM’s vocabulary of input tokens, while relation label embeddings are input to the self-attention functions for every pair of related tokens. These new embeddings represent learnable parameters that are trained during the PLM fine-tuning on the downstream tasks. As reported in Appendix Subsection A.2, GADePo adds a negligible number of extra parameters, namely only the special token inputs and the graph directional relation inputs.

**Special Token <ent>** To tackle the EE pooling function, we add to the input tokens  $\mathcal{T}_D$  as many <ent> special tokens as entities in the document. This way each entity  $e$  has a corresponding entity token <ent> in the input. We connect each <ent> token with its corresponding cluster of mentions  $\mathcal{M}_e = \{m_i\}_{i=1}^{|\mathcal{M}_e|}$ , and vice-versa. The two graph relations we use are thus <ent>  $\rightarrow$  \* and \*  $\rightarrow$  <ent>, where \* represents the special token at the starting position of mention  $m_i$ . Each <ent> token receives the same <ent> embedding, with no positional encoding, since each one collectively represents a set of mentions from different positions in the input graph. These identical inputs are only disambiguated through the connections to and from mentions expressed as the <ent>  $\rightarrow$  \* and \*  $\rightarrow$  <ent> graph relations. These relations tell the self-attention mechanism to use the <ent> token to aggregate information from the associated mentions, and thus the <ent> tokens have a direct correspondence to the computed  $h_e$  in Equation 2.

**Special Token <pent>** ATLOP performs information filtering by calculating via Equation 4 a localised context embedding (LCE)  $c^{(s,o)}$  that is dependent on the cross-token attention matrix  $A$  output by the PLM. The intuition behind it is that the dependencies between different tokens are encoded as attention weights. We propose a straight-

forward adjustment of the input graph used for the EE pooling to effectively model and capture these dependencies. To address the LCE pooling function, we add to the input tokens  $\mathcal{T}_D$  as many <pent> special tokens as the number of all possible pairs of entities. Each special token <pent> thus refers to a pair of entities  $(e_s, e_o)$ . We connect each <pent> token with each mention in the two clusters of mentions  $\mathcal{M}_{e_s} = \{m_i\}_{i=1}^{|\mathcal{M}_{e_s}|}$  and  $\mathcal{M}_{e_o} = \{m_i\}_{i=1}^{|\mathcal{M}_{e_o}|}$  and vice-versa. Since the attention weights used in LCE are computed from these mention embeddings, we expect that they are sufficient for the Transformer to learn to find the relevant contexts. The two graph relations we use are thus <pent>  $\rightarrow$  \* and \*  $\rightarrow$  <pent>. Analogously to the <ent> tokens, the <pent> tokens all receive the same <pent> embedding, with no positional embeddings, and thus are only disambiguated by their different <pent>  $\rightarrow$  \* and \*  $\rightarrow$  <pent> graph relations. These relations tell the <pent> token to pay attention to its associated mentions, which in turn allows it to find the relevant context shared by these mentions. Thus, each <pent> token can be seen as having a direct correspondence to the computed  $c^{(s,o)}$  in Equation 4.

All equations relative to the relation classification and the corresponding loss function reported in Appendix Subsection A.1 remain valid as we merely substitute the hand-coded computations of  $h_e$  and  $c^{(s,o)}$  with the embeddings of <ent> and <pent>, respectively.

**Text-Graph Encoding** We follow Mohammadshahi and Henderson (2020, 2021); Henderson et al. (2023) in leveraging the intrinsic graph processing capabilities of the Transformer model by incorporating graph relations as relation embeddings input to the self-attention function. For every pair of input tokens  $ij$ , the pre-softmax attention weight  $e_{ij} \in \mathbb{R}$  is computed from both the respective token embeddings  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$ , and an embeddings of the graph relation  $\mathbf{c}_{ij}$  between the  $i$ -th and  $j$ -th tokens. However, we change the attention weight computation to:

$$e_{ij} = \frac{\mathbf{x}_i \mathbf{W}_Q \text{diag}(\text{LN}(\mathbf{c}_{ij} \mathbf{W}_C)) (\mathbf{x}_j \mathbf{W}_K)^\top}{\sqrt{d}}, \quad (5)$$

where  $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d}$  represent the query and key matrices, respectively.  $\mathbf{c}_{ij} \in \{0, 1\}^{|\mathcal{C}|}$  represents a 0/1 encoded label of the graph relation between the  $i$ -th and  $j$ -th input elements, and  $\mathbf{W}_C \in \mathbb{R}^{|\mathcal{C}| \times d}$  represents the relations’ embedding

Model	Aggregation	Re-DocRED		HacRED		
		Ign $F_1$	$F_1$	$P$	$R$	$F_1$
ATLOP*	$h_e$	75.27	75.92	<b>76.27</b>	76.83	76.55
GADePo (ours)	<ent>	<b>75.55</b>	<b>76.38</b>	74.13	<b>79.46</b>	<b>76.70</b>
ATLOP <sup>*,<math>\diamond</math></sup>	$h_e ; c^{(s,o)}$	76.82	77.56	77.89	76.55	77.21
ATLOP*	$h_e ; c^{(s,o)}$	77.62	78.38	76.36	78.86	77.59
GADePo (ours)	<ent> ; <pent>	<b>77.70</b>	<b>78.40</b>	<b>78.27</b>	<b>79.03</b>	<b>78.65</b>

Table 1: Comparative analysis between the previous method ATLOP and the proposed method GADePo on the test set. ATLOP\* indicates our reimplement of the previous method. For Re-DocRED and HacRED we report in percentage the results obtained by Tan et al. (2022b) (ATLOP\*) and Cheng et al. (2021) (ATLOP<sup>◊</sup>), respectively. The results are reported in terms of  $F_1$  scores, Precision ( $P$ ), and Recall ( $R$ ), following the same metrics reported in prior research specific to each dataset. Ign  $F_1$  denotes the  $F_1$  score that excludes relational facts shared between the training and evaluation sets. We also comply with the standard practice where test scores are determined based on the best checkpoint from five training runs with distinct random seeds.

matrix, so  $c_{ij}W_C$  is the embedding of the relation between  $i$  and  $j$ . Finally, LN stands for the *LayerNorm* operation and *diag* returns a diagonal matrix.

Compared to the standard attention function, where  $e_{ij} = x_iW_Q(x_jW_K)^T/\sqrt{d}$ , the relation embedding determines a weighting of the different dimensions. This is a novel way to condition on the relation embedding compared to the original formulation, which only models query-relation interactions (Mohammadshahi and Henderson, 2020). This change is motivated by our task requiring a more flexible formulation which models query-relation-key interactions via a multiplicative mechanism, without requiring a full  $d \times d$  matrix of bilinear parameters. This way, a key will be relevant to a query only when both agree on the relation. In preliminary experiments, we explored various methods for biasing attention and found that the formulation presented in Equation 5 produced the best results.

## 5 Experiments

### 5.1 Datasets and Models

**Re-DocRED** (Tan et al., 2022b) is a revisited version of the DocRED (Yao et al., 2019) dataset. It is built from English Wikipedia and Wikidata and contains both distantly-supervised and human-annotated documents with named entities, coreference data, and intra- and inter-sentence relations, supported by evidence. It requires analysing multiple sentences to identify entities, establish their relationships, and integrate information from the entire document. We comply with the model used by the authors and employ the RoBERTa<sub>LARGE</sub> (Zhuang et al., 2021) model in our experiments.

**HacRED** (Cheng et al., 2021) is a large-scale, high-quality Chinese document-level relation extraction dataset, with a special focus on practical hard cases. As the authors did not provide specific information about the model used in their study, we conducted our experiments using the Chinese BERT<sub>BASE</sub> with whole word masking model (Cui et al., 2021).

**Datasets statistics** Re-DocRED and HacRED exhibit notable distinctions in their statistics, as summarised in Table 2. Re-DocRED comprises a larger number of facts, entities per document, and relations compared to HacRED. This indicates a potentially richer and more extensive dataset in terms of factual information and relationship types. However, HacRED contains more documents and may present a broader range of scenarios for relation extraction, including more challenging cases, as it has been specifically created with a focus on practical hard cases.

Statistic	Re-DocRED	HacRED
Facts	120,664	65,225
Relations	96	26
Documents	4,053	9,231
Average Entities	19.4	10.8

Table 2: Re-DocRED and HacRED human-annotated datasets statistics.

### 5.2 Results and Discussion

We follow the standard practice from prior research and report the results of our experiments on the Re-DocRED and HacRED datasets in Table 1 and Figure 4. For all datasets and models, we provide our reimplement of the ATLOP baseline (indicated as ATLOP\*), which achieves or surpasses pre-



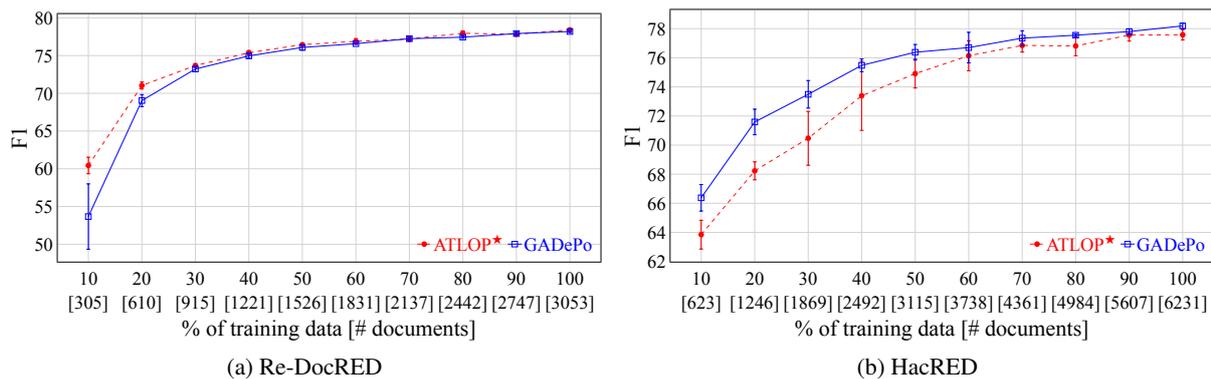


Figure 4: Performance of ATLOP\* ( $h_e ; c^{(s,o)}$ ) and GADePo (<ent> ; <pent>) on the development set under varying data availability conditions on Re-DocRED (4a) and HacRED (4b). The  $x$ -axis represents the percentage and number of documents from the training dataset, while the  $y$ -axis displays the  $F_1$  score in percentage. Each point on the graph represents the mean value, while error bars indicate the standard deviation derived from five distinct training runs with separate random seeds.

<pent> into ATLOP\* and GADePo, respectively, highlight the significant contributions of these features. GADePo outperforms ATLOP\* with an  $F_1$  score of 78.65% compared to 77.59%. This larger improvement on HacRED suggests that GADePo is better at handling challenging cases, which is not surprising given its greater flexibility over the fixed pooling functions of ATLOP.

**Data Ablation** To evaluate the models’ sensitivity to dataset size, the performance evaluation depicted in Figure 4 compares ATLOP\* ( $h_e ; c^{(s,o)}$ ) and GADePo (<ent> ; <pent>) on the development set, considering different levels of training data availability on the Re-DocRED and HacRED datasets. Accuracies generally converge as the dataset sizes increase, but on the challenging cases of HacRED, GADePo maintains a substantial advantage across the full range. On Re-DocRED, GADePo catches up with and slightly outperforms ATLOP\* as data size increases. This lower performance on smaller datasets is presumably because GADePo must learn how to exploit the graph relations to the special tokens <ent> and <pent> and pool information through them, whereas for ATLOP this pooling is hand-coded. On the Re-DocRED dataset, ATLOP\* appears to have relatively consistent variance, while GADePo exhibits higher variance in the smaller training sets, while on the HacRED dataset, GADePo is significantly more stable for smaller datasets.

The data ablation analysis shows that the performance of hand-coded pooling functions can be dataset-specific, which restricts their adaptability. In contrast, GADePo consistently outperforms its

hand-coded counterparts on larger datasets, and matches them on all but some smaller datasets, presumably due to its flexibility. This pattern suggests that GADePo has a greater potential for optimisation, particularly on larger datasets. This is supported by GADePo’s better performance on HacRED, which is both larger and designed to be more challenging than Re-DocRED.

## 6 Conclusion

In this paper we proposed a novel approach to document-level relation extraction, challenging the conventional reliance on hand-coded pooling functions for information aggregation. Our method leverages the power of Transformer models by incorporating explicit graph relations as instructions for information aggregation. By combining graph processing with text-based encoding, we introduced the graph-assisted declarative pooling (GADePo) specification, which allows for more flexible and customisable specification of pooling strategies which are still learned from data.

We conducted evaluations using diverse datasets and models commonly employed in document-level relation extraction tasks. The results of our experiments demonstrated that our approach achieves promising performance that is comparable to or better than that of hand-coded pooling functions. This suggests that our method can serve as a viable basis for other relation extraction methods, providing a more adaptable and tailored approach. In particular, recent methods have improved performance by exploiting information about evidence, which can naturally be incorporated in our graph-based approach.

## 618 Limitations

619 While the proposed GADePo model offers a  
620 promising and innovative approach to relation ex-  
621 traction, there are issues which the current study  
622 does not address. According to the data in Ap-  
623 pendix Table 2, the average number of entities per  
624 document across datasets is approximately 15. This  
625 means that, on average, there will be an additional  
626 15 <ent> tokens and 105 <pent> tokens. Given  
627 that the maximum allowable input length for the  
628 models is 512 tokens, the inclusion of these extra  
629 tokens results in roughly a 3% and 20% increase  
630 in the overall input length for <ent> and <pent>,  
631 respectively. It’s evident that the majority of the  
632 increase in input length is due to the quadratic num-  
633 ber of <pent> special tokens, but we believe that  
634 an appropriate pruning strategy could easily reduce  
635 this number to linear in the number of entities with-  
636 out degrading accuracy. One such pruning strategy  
637 could involve an <ent>-only model with a binary  
638 classifier which is trained to predict pairs of related  
639 entities. This model could then be used to prune  
640 the set of candidate entity pairs for the final relation  
641 classification, with <pent> tokens being instanti-  
642 ated only for these candidate pairs. We have chosen  
643 to leave this approach as a potential avenue for fu-  
644 ture work, opting instead to focus on demonstrating  
645 the promise of the current simpler formulation.

## 646 Ethics Statement

647 We do not anticipate any ethical concerns related to  
648 our work, as it primarily presents an alternative ap-  
649 proach to a previously proposed method. Our main  
650 contribution lies in introducing a novel method-  
651 ology for relation extraction. In our experiments,  
652 we use the same datasets and pretrained models as  
653 previous research, all of which are publicly avail-  
654 able. However, it is important to acknowledge that  
655 these datasets and models may still require further  
656 examination for potential fairness issues and the  
657 knowledge they encapsulate.

## 658 References

659 Michele Banko, Michael J. Cafarella, Stephen Soder-  
660 land, Matthew Broadhead, and Oren Etzioni. 2007.  
661 Open information extraction from the web. In  
662 *CACM*.

663 Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [SciB-](#)  
664 [ERT: A pretrained language model for scientific text](#).  
665 In *Proceedings of the 2019 Conference on Empirical*  
666 *Methods in Natural Language Processing and the*

*9th International Joint Conference on Natural Lan- 667*  
*guage Processing (EMNLP-IJCNLP)*, pages 3615– 668  
3620, Hong Kong, China. Association for Computa- 669  
tional Linguistics. 670

Danqi Chen, Adam Fisch, Jason Weston, and Antoine 671  
Bordes. 2017. [Reading Wikipedia to answer open-](#) 672  
[domain questions](#). In *Proceedings of the 55th Annual* 673  
*Meeting of the Association for Computational Lin-* 674  
*guistics (Volume 1: Long Papers)*, pages 1870–1879, 675  
Vancouver, Canada. Association for Computational 676  
Linguistics. 677

Qiao Cheng, Juntao Liu, Xiaoye Qu, Jin Zhao, Jiaqing 678  
Liang, Zhefeng Wang, Baoxing Huai, Nicholas Jing 679  
Yuan, and Yanghua Xiao. 2021. [HacRED: A large-](#) 680  
[scale relation extraction dataset toward hard cases in](#) 681  
[practical applications](#). In *Findings of the Association* 682  
*for Computational Linguistics: ACL-IJCNLP 2021*, 683  
pages 2819–2831, Online. Association for Computa- 684  
tional Linguistics. 685

Hyunjin Choi, Judong Kim, Seongho Joe, and 686  
Youngjune Gwon. 2021. [Evaluation of bert and albert](#) 687  
[sentence embedding performance on downstream nlp](#) 688  
[tasks](#). *2020 25th International Conference on Pattern* 689  
*Recognition (ICPR)*, pages 5482–5487. 690

Fenia Christopoulou, Makoto Miwa, and Sophia Ana- 691  
niadou. 2019. [Connecting the dots: Document-level](#) 692  
[neural relation extraction with edge-oriented graphs](#). 693  
In *Proceedings of the 2019 Conference on Empirical* 694  
*Methods in Natural Language Processing and the* 695  
*9th International Joint Conference on Natural Lan-* 696  
*guage Processing (EMNLP-IJCNLP)*, pages 4925– 697  
4936, Hong Kong, China. Association for Computa- 698  
tional Linguistics. 699

Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc 700  
Barrault, and Antoine Bordes. 2017. [Supervised](#) 701  
[learning of universal sentence representations from](#) 702  
[natural language inference data](#). In *Proceedings of* 703  
*the 2017 Conference on Empirical Methods in Nat-* 704  
*ural Language Processing*, pages 670–680, Copen- 705  
hagen, Denmark. Association for Computational Lin- 706  
guistics. 707

Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, and 708  
Ziqing Yang. 2021. [Pre-training with whole word](#) 709  
[masking for chinese bert](#). *IEEE/ACM Trans. Audio,* 710  
*Speech and Lang. Proc.*, 29:3504–3514. 711

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 712  
Kristina Toutanova. 2019. [BERT: Pre-training of](#) 713  
[deep bidirectional transformers for language under-](#) 714  
[standing](#). In *Proceedings of the 2019 Conference of* 715  
*the North American Chapter of the Association for* 716  
*Computational Linguistics: Human Language Tech-* 717  
*nologies, Volume 1 (Long and Short Papers)*, pages 718  
4171–4186, Minneapolis, Minnesota. Association for 719  
Computational Linguistics. 720

William Falcon and The PyTorch Lightning team. 2019. 721  
[PyTorch Lightning](#). 722

723	Yue Feng, Zhen Han, Mingming Sun, and Ping Li.	Alireza Mohammadshahi and James Henderson.	779
724	2022. <a href="#">Multi-hop open-domain question answering over structured and unstructured knowledge</a> . In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 151–156, Seattle, United States. Association for Computational Linguistics.	2020. <a href="#">Graph-to-graph transformer for transition-based dependency parsing</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3278–3289, Online. Association for Computational Linguistics.	780
725			781
726			782
727			783
728			784
729	James Henderson. 2020. <a href="#">The unstoppable rise of computational linguistics in deep learning</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6294–6306, Online. Association for Computational Linguistics.	Alireza Mohammadshahi and James Henderson. 2021. <a href="#">Recursive non-autoregressive graph-to-graph transformer for dependency parsing with iterative refinement</a> . <i>Transactions of the Association for Computational Linguistics</i> , 9:120–138.	785
730			786
731			787
732			788
733			789
734	James Henderson, Alireza Mohammadshahi, Andrei Coman, and Lesly Miculicich. 2023. <a href="#">Transformers as graph-to-graph models</a> . In <i>Proceedings of the Big Picture Workshop</i> , pages 93–107, Singapore. Association for Computational Linguistics.	Alireza Mohammadshahi and James Henderson. 2023. <a href="#">Syntax-aware graph-to-graph transformer for semantic role labelling</a> . In <i>Proceedings of the 8th Workshop on Representation Learning for NLP (RepL4NLP 2023)</i> , pages 174–186, Toronto, Canada. Association for Computational Linguistics.	790
735			791
736			792
737			793
738			794
739	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. <i>Neural Computation</i> , 9:1735–1780.		795
740			
741			
742	Shaoxiong Ji, Shirui Pan, E. Cambria, Pekka Martinen, and Philip S. Yu. 2020. A survey on knowledge graphs: Representation, acquisition, and applications. <i>IEEE Transactions on Neural Networks and Learning Systems</i> , 33:494–514.	Guoshun Nan, Zhijiang Guo, Ivan Sekulic, and Wei Lu. 2020. <a href="#">Reasoning with latent structure refinement for document-level relation extraction</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1546–1557, Online. Association for Computational Linguistics.	796
743			797
744			798
745			799
746			800
747	Robin Jia, Cliff Wong, and Hoifung Poon. 2019. <a href="#">Document-level n-ary relation extraction with multi-scale representation learning</a> . In <i>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)</i> , pages 3693–3704, Minneapolis, Minnesota. Association for Computational Linguistics.	Bo Pang and Lillian Lee. 2008. Opinion mining and sentiment analysis. <i>Found. Trends Inf. Retr.</i> , 2:1–135.	802
748			803
749			
750			
751		Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. <a href="#">PyTorch: An Imperative Style, High-Performance Deep Learning Library</a> . In <i>Advances in Neural Information Processing Systems 32</i> , pages 8024–8035. Curran Associates, Inc.	804
752			805
753			806
754			807
755	Thomas N. Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In <i>International Conference on Learning Representations (ICLR)</i> .		808
756			809
757			810
758			811
759	Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. 2020. <a href="#">On the variance of the adaptive learning rate and beyond</a> . In <i>International Conference on Learning Representations</i> .	Nanyun Peng, Hoifung Poon, Chris Quirk, Kristina Toutanova, and Wen-tau Yih. 2017. <a href="#">Cross-sentence n-ary relation extraction with graph LSTMs</a> . <i>Transactions of the Association for Computational Linguistics</i> , 5:101–115.	815
760			816
761			817
762			818
763			819
764	Youni Ma, An Wang, and Naoaki Okazaki. 2023. <a href="#">DREEAM: Guiding attention with evidence for improving document-level relation extraction</a> . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 1971–1983, Dubrovnik, Croatia. Association for Computational Linguistics.	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-BERT: Sentence embeddings using Siamese BERT-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	820
765			821
766			822
767			823
768			824
769			825
770			826
771	Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. <i>Introduction to Information Retrieval</i> . Cambridge University Press.		827
772			
773			
774	Lesly Miculicich and James Henderson. 2022. <a href="#">Graph refinement for coreference resolution</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2732–2742, Dublin, Ireland. Association for Computational Linguistics.	Qingyu Tan, Ruidan He, Lidong Bing, and Hwee Tou Ng. 2022a. <a href="#">Document-level relation extraction with adaptive focal loss and knowledge distillation</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 1672–1681, Dublin, Ireland. Association for Computational Linguistics.	828
775			829
776			830
777			831
778			832
			833

834	Qingyu Tan, Lu Xu, Lidong Bing, Hwee Tou Ng, and Sharifah Mahani Aljunied. 2022b. <a href="#">Revisiting DocRED - addressing the false negative problem in relation extraction</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 8472–8487, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
835		
836		
837		
838		
839		
840		
841		
842	Hengzhu Tang, Yanan Cao, Zhenyu Zhang, Jiangxia Cao, Fang Fang, Shi Wang, and Pengfei Yin. 2020. <a href="#">Hin: Hierarchical inference network for document-level relation extraction</a> . In <i>Advances in Knowledge Discovery and Data Mining: 24th Pacific-Asia Conference, PAKDD 2020, Singapore, May 11–14, 2020, Proceedings, Part I 24</i> , pages 197–209. Springer.	
843		
844		
845		
846		
847		
848		
849	Christos Theodoropoulos, James Henderson, Andrei Catalin Coman, and Marie-Francine Moens. 2021. <a href="#">Imposing relation structure in language-model embeddings using contrastive learning</a> . In <i>Proceedings of the 25th Conference on Computational Natural Language Learning</i> , pages 337–348, Online. Association for Computational Linguistics.	
850		
851		
852		
853		
854		
855		
856	Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all you need</a> . In <i>NIPS</i> .	
857		
858		
859		
860	Difeng Wang, Wei Hu, Ermei Cao, and Weijian Sun. 2020. <a href="#">Global-to-local neural networks for document-level relation extraction</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3711–3721, Online. Association for Computational Linguistics.	
861		
862		
863		
864		
865		
866	Hong Wang, Christfried Focke, Rob Sylvester, Nilesh Mishra, and William Yang Wang. 2019. <a href="#">Fine-tune bert for docred with two-step process</a> . <i>ArXiv</i> , abs/1909.11898.	
867		
868		
869		
870	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Transformers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	
871		
872		
873		
874		
875		
876		
877		
878		
879		
880		
881		
882	Yuxin Xiao, Zecheng Zhang, Yuning Mao, Carl Yang, and Jiawei Han. 2022. <a href="#">SAIS: Supervising and augmenting intermediate steps for document-level relation extraction</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2395–2409, Seattle, United States. Association for Computational Linguistics.	
883		
884		
885		
886		
887		
888		
889		
890		
	Benfeng Xu, Quan Wang, Yajuan Lyu, Yong Zhu, and Zhendong Mao. 2021. <a href="#">Entity structure within and throughout: Modeling mention dependencies for document-level relation extraction</a> . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 35(16):14149–14157.	891 892 893 894 895 896
	Yuan Yao, Deming Ye, Peng Li, Xu Han, Yankai Lin, Zhenghao Liu, Zhiyuan Liu, Lixin Huang, Jie Zhou, and Maosong Sun. 2019. <a href="#">DocRED: A large-scale document-level relation extraction dataset</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 764–777, Florence, Italy. Association for Computational Linguistics.	897 898 899 900 901 902 903 904
	Deming Ye, Yankai Lin, Peng Li, and Maosong Sun. 2022. <a href="#">Packed levitated marker for entity and relation extraction</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4904–4917, Dublin, Ireland. Association for Computational Linguistics.	905 906 907 908 909 910
	Shuang Zeng, Runxin Xu, Baobao Chang, and Lei Li. 2020. <a href="#">Double graph based reasoning for document-level relation extraction</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1630–1640, Online. Association for Computational Linguistics.	911 912 913 914 915 916
	Ningyu Zhang, Xiang Chen, Xin Xie, Shumin Deng, Chuanqi Tan, Mosha Chen, Fei Huang, Luo Si, and Huajun Chen. 2021. <a href="#">Document-level relation extraction as semantic segmentation</a> . In <i>International Joint Conference on Artificial Intelligence</i> .	917 918 919 920 921
	Wenxuan Zhou, Kevin Huang, Tengyu Ma, and Jing Huang. 2021. <a href="#">Document-level relation extraction with adaptive thresholding and localized context pooling</a> . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> .	922 923 924 925 926
	Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. 2021. <a href="#">A robustly optimized BERT pre-training approach with post-training</a> . In <i>Proceedings of the 20th Chinese National Conference on Computational Linguistics</i> , pages 1218–1227, Huhhot, China. Chinese Information Processing Society of China.	927 928 929 930 931 932
	<b>A Appendix</b>	933
	<b>A.1 ATLOP: Relation Classification and Loss Function</b>	934 935
	<b>Relation Classification</b> To predict the relation between the subject entity $e_s$ and object entity $e_o$ , ATLOP first generates context-aware subject and object representations as follows:	936 937 938 939
	$z_s = \tanh(\mathbf{W}_s[\mathbf{h}_{e_s}; \mathbf{c}^{(s,o)}] + \mathbf{b}_s)$	(6) 940
	$z_o = \tanh(\mathbf{W}_o[\mathbf{h}_{e_o}; \mathbf{c}^{(s,o)}] + \mathbf{b}_o),$	(7) 941 942
	where $z_s, z_o \in \mathbb{R}^d$ , $[\cdot; \cdot]$ represents the concatenation of two vectors, and $\mathbf{W}_s, \mathbf{W}_o \in \mathbb{R}^{d \times 2d}$ together	943 944

with  $\mathbf{b}_s, \mathbf{b}_o \in \mathbb{R}^d$  are trainable parameters. Then, the entity pair representation is computed as:

$$\mathbf{x}^{(s,o)} = \mathbf{z}_s \otimes \mathbf{z}_o, \quad (8)$$

where  $\mathbf{x}^{(s,o)} \in \mathbb{R}^{d^2}$  and  $\otimes$  stands for the vectorised Kronecker product. Finally, relation scores are computed as:

$$\mathbf{y}^{(s,o)} = \mathbf{W}_r \mathbf{x}^{(s,o)} + \mathbf{b}_r, \quad (9)$$

where  $\mathbf{y}^{(s,o)} \in \mathbb{R}^{|\mathcal{R}|}$ , with  $\mathbf{W}_r \in \mathbb{R}^{|\mathcal{R}| \times d^2}$  and  $\mathbf{b}_r \in \mathbb{R}^{|\mathcal{R}|}$  representing learnable parameters. The probability of relation  $r \in \mathcal{R}$  between the subject and object entities is computed as follows:

$$P(r|s, o) = \sigma(\mathbf{y}^{(s,o)}), \quad (10)$$

where  $\sigma$  is the sigmoid function. To reduce the number of parameters in the classifier, a grouped function is used, which splits the embedding dimensions into  $k$  equal-sized groups and applies the function within the groups as follows:

$$\mathbf{z}_s = [\mathbf{z}_s^1; \dots; \mathbf{z}_s^k] \quad (11)$$

$$\mathbf{z}_o = [\mathbf{z}_o^1; \dots; \mathbf{z}_o^k] \quad (12)$$

$$\mathbf{x}^{(s,o)} = [\mathbf{x}^{(s,o)^1}; \dots; \mathbf{x}^{(s,o)^k}] \quad (13)$$

$$\mathbf{y}^{(s,o)} = \sum_{i=1}^k \mathbf{W}_r^i \mathbf{x}^{(s,o)^i} + \mathbf{b}_r, \quad (14)$$

where  $\mathbf{z}_s^i, \mathbf{z}_o^i \in \mathbb{R}^{d/k}$ ,  $\mathbf{x}^{(s,o)^i} \in \mathbb{R}^{d^2/k}$ , and  $\mathbf{W}_r^i \in \mathbb{R}^{|\mathcal{R}| \times d^2/k}$ . This way, the number of parameters can be reduced from  $d^2$  to  $d^2/k$ .

**Loss Function** ATLOP introduces the adaptive thresholding loss concept. This approach involves training a model to learn a hypothetical threshold class  $TH$ , which dynamically adjusts for each relation class  $r \in \mathcal{R}$ . During training, for each entity pair  $(e_s, e_o)$ , the loss enforces the model to generate scores above  $TH$  for positive relation classes  $\mathcal{R}_P$  and scores below  $TH$  for negative relation classes  $\mathcal{R}_N$ . The loss is computed as follows:

$$\mathcal{L} = - \sum_{s \neq o} \sum_{r \in \mathcal{R}_P} \frac{\exp(y_r^{(s,o)})}{\sum_{r' \in \mathcal{R}_P \cup \{TH\}} \exp(y_{r'}^{(s,o)})} - \frac{\exp(y_{TH}^{(s,o)})}{\sum_{r' \in \mathcal{R}_N \cup \{TH\}} \exp(y_{r'}^{(s,o)})} \quad (15)$$

## A.2 GADePo’s Extra Parameters

GADePo introduces few extra parameters to the PLM. The amount of parameters is reported in Table 4.

Parameter	Model	
	RoBERTa <sub>LARGE</sub>	BERT <sub>BASE</sub>
<ent>	1024	768
<pent>	1024	768
<ent> $\rightarrow$ *	24 $\times$ 1024	12 $\times$ 768
* $\rightarrow$ <ent>	24 $\times$ 1024	12 $\times$ 768
<pent> $\rightarrow$ *	24 $\times$ 1024	12 $\times$ 768
* $\rightarrow$ <pent>	24 $\times$ 1024	12 $\times$ 768
<b>Total</b>	100,352	38,400

Table 4: GADePo’s extra parameters count.

## A.3 Training Details

We generally comply with the hyperparameters of ATLOP and set the output dimension in Equation 6 and Equation 7 to 768. We also set the block size in Equation 11 and Equation 12 to 64, i.e.,  $k = 12$ .

In all our experiments we perform early stopping on the development set based on the Ign  $F_1 + F_1$  score for DocRED and Re-DocRED, and  $F_1$  score for HacRED. The five different seeds we use are  $\{73, 21, 37, 7, 3\}$ .

We use RAdam (Liu et al., 2020) as our optimizer. On the RoBERTa<sub>LARGE</sub> based models we train for 8 epochs and set the learning rates to  $3e^{-5}$  and  $1e^{-4}$  for the PLM parameters and the new additional parameters, respectively. On the BERT<sub>BASE</sub> based models we train for 10 epochs and set the learning rates to  $1e^{-5}$  and  $1e^{-4}$  for the PLM parameters and the new additional parameters, respectively. We use a cosine learning rate decay throughout the training process.

In all our experiments the batch size is set to 4 for ATLOP and 2 for GADePo, with gradient accumulation set to 1 and 2, for ATLOP and GADePo, respectively. We clip the gradients to a max norm of 1.0. All models are trained with mixed precision.

We run our experiments on two types of GPUs, namely the NVIDIA V100 32GB for the RoBERTa<sub>LARGE</sub> based models and NVIDIA RTX 3090 24GB for the BERT<sub>BASE</sub> based models, respectively.

We use PyTorch (Paszke et al., 2019), Lightning (Falcon and The PyTorch Lightning team, 2019), and Hugging Face’s Transformers (Wolf et al., 2020) libraries to develop our models.

Model	Aggregation	Dev		Test	
		Ign $F_1$	$F_1$	Ign $F_1$	$F_1$
ATLOP*	$h_e$	$75.46 \pm 0.16$	$76.16 \pm 0.16$	75.27	75.92
GADePo (ours)	<ent>	$75.46 \pm 0.20$	$76.31 \pm 0.24$	<b>75.55</b>	<b>76.38</b>
ATLOP <sup>o</sup>	$h_e ; c^{(s,o)}$	76.79	77.46	76.82	77.56
ATLOP*	$h_e ; c^{(s,o)}$	$77.75 \pm 0.08$	$78.41 \pm 0.10$	77.62	78.38
GADePo (ours)	<ent> ; <pent>	$77.48 \pm 0.12$	$78.19 \pm 0.14$	<b>77.70</b>	<b>78.40</b>

Table 5: Results in percentage for the development and test sets of Re-DocRED. We report the results obtained by Tan et al. (2022b) (ATLOP<sup>o</sup>) on Re-DocRED. ATLOP\* indicates our reimplement of the previous method. We report the mean and standard deviation of Ign  $F_1$  and  $F_1$  on the development set, calculated from five training runs with distinct random seeds. We report the test score achieved by the best checkpoint on the development set. Ign  $F_1$  refers to the  $F_1$  score that excludes relational facts shared between the training and development/test sets.

Model	Aggregation	Dev			Test		
		$P$	$R$	$F_1$	$P$	$R$	$F_1$
ATLOP*	$h_e$	$77.37 \pm 0.22$	$77.40 \pm 0.31$	$77.39 \pm 0.13$	<b>76.27</b>	76.83	76.55
GADePo (ours)	<ent>	$72.96 \pm 0.96$	$79.22 \pm 1.20$	$75.96 \pm 0.99$	74.13	<b>79.46</b>	<b>76.70</b>
ATLOP <sup>o</sup>	$h_e ; c^{(s,o)}$	–	–	–	77.89	76.55	77.21
ATLOP*	$h_e ; c^{(s,o)}$	$77.18 \pm 0.14$	$77.98 \pm 0.66$	$77.58 \pm 0.36$	76.36	78.86	77.59
GADePo (ours)	<ent> ; <pent>	$75.98 \pm 0.94$	$80.54 \pm 0.72$	$78.19 \pm 0.19$	<b>78.27</b>	<b>79.03</b>	<b>78.65</b>

Table 6: Results in percentage for the development and test sets of HacRED. We report the results obtained by Cheng et al. (2021) (ATLOP<sup>o</sup>) on HacRED. ATLOP\* indicates our reimplement of the previous method. We report the mean and standard deviation of Precision ( $P$ ), Recall ( $R$ ) and  $F_1$  on the development set, calculated from five training runs with distinct random seeds. We report the test score achieved by the best checkpoint on the development set.

Model	Aggregation	Dev		Test	
		Ign $F_1$	$F_1$	Ign $F_1$	$F_1$
ATLOP*	$h_e$	$59.66 \pm 0.20$	$61.60 \pm 0.21$	59.22	61.37
GADePo (ours)	<ent>	$59.04 \pm 0.52$	$61.18 \pm 0.46$	<b>59.30</b>	<b>61.63</b>
ATLOP <sup>o</sup>	$h_e ; c^{(s,o)}$	$61.32 \pm 0.14$	$63.18 \pm 0.19$	61.39	63.40
ATLOP*	$h_e ; c^{(s,o)}$	$61.41 \pm 0.26$	$63.38 \pm 0.28$	<b>61.62</b>	63.72
GADePo (ours)	<ent> ; <pent>	$61.19 \pm 0.55$	$63.26 \pm 0.48$	61.52	<b>63.75</b>

Table 7: Results in percentage for the development and test sets of DocRED. We report the results obtained by Zhou et al. (2021) (ATLOP<sup>o</sup>) on DocRED. ATLOP\* indicates our reimplement of the previous method. We report the mean and standard deviation of Ign  $F_1$  and  $F_1$  on the development set, calculated from five training runs with distinct random seeds. We report the test score achieved by the best checkpoint on the development set. Ign  $F_1$  refers to the  $F_1$  score that excludes relational facts shared between the training and development/test sets.

1020

#### A.4 Additional Results

1021

**Re-DocRED and HacRED** Table 5 and Table 6 present additional results for Re-DocRED and HacRED, respectively. In addition to the results outlined in Section 5, these tables include the mean and standard deviation on the development set, calculated from five training runs with distinct random seeds, as reported in Appendix Subsection A.3.

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

**DocRED results** The DocRED (Yao et al., 2019) dataset consists of 56,354 facts, 96 relations, 5,053 documents, and 26.2 average number of entities per document. In line with the approach taken for Re-DocRED and HacRED, Table 7 and Figure 5 illustrate the results for DocRED.

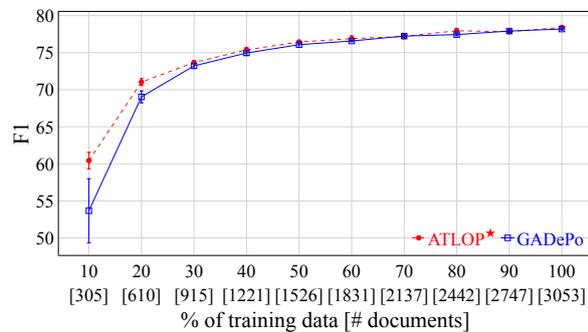


Figure 5: Performance of ATLOP\* ( $h_e ; c^{(s,o)}$ ) and GADePo (<ent> ; <pent>) on the development set under varying data availability conditions on DocRED. The  $x$ -axis represents the percentage and number of documents from the training dataset, while the  $y$ -axis displays the  $F_1$  score in percentage. Each point on the graph represents the mean value, while error bars indicate the standard deviation derived from five distinct training runs with separate random seeds.